

High-Precision Photogrammetric 3D Modeling Technology Based on Multi-Source Data Fusion and Deep Learning-Enhanced Feature Learning Using Internet of Things Big Data

Guangtao Zhang¹, Jun Zhang^{2,3*}

¹College of Information Engineering, Yangzhou Polytechnic College, Yangzhou 225000, Jiangsu, China

²Center of Engineering Training, Yangzhou Polytechnic College, Yangzhou 225000, Jiangsu, China

³Jiangsu Safety & Environment Technology and Equipment for Planting and Breeding Industry Engineering Research Center, Yangzhou 225000, Jiangsu, China

E-mail: zhang-jun168@hotmail.com

*Corresponding author

Keywords: internet of things, big data, high-precision photography, 3D modeling technology

Received: August 15, 2024

As technology advances and application demands grow, high-precision three-dimensional (3D) modeling is increasingly essential for urban planning, disaster management, and cultural heritage protection. This study presents a high-precision photogrammetric 3D modeling approach with a focus on integrating multi-source data fusion techniques for complex terrains. The methodology incorporates aerial imagery, LiDAR data, ground survey data, and meteorological corrections, covering the entire workflow from data preprocessing, feature extraction, and registration to multi-source data fusion. Key innovations include an adaptive weight adjustment strategy, global optimization registration techniques, and deep learning-assisted feature learning, all contributing to significant improvements in model accuracy and reliability. Experimental results show a X% improvement in spatial accuracy and a Y% reduction in mean squared error (MSE), along with enhanced morphological structure recovery and visual effects. These improvements have been validated through practical applications and received positive feedback from users. The detailed technical implementation of the data fusion algorithms, along with the quantitative performance metrics, further demonstrates the efficacy of the proposed methodology in real-world scenarios.

Povzetek: Raziskava vpelje visokoločljivostno fotogrametrično 3D modeliranje z uporabo fuzije več virov podatkov in globokega učenja. Tehnologija izboljšuje natančnost modelov z integracijo satelitskih slik, LiDAR-ja in meteoroloških podatkov ter prilagodljivimi optimizacijskimi algoritmi. Eksperimentalni rezultati kažejo znatno izboljšano vizualno rekonstrukcijo, kar omogoča uporabo v urbanističnem načrtovanju, varstvu kulturne dediščine in obvladovanju naravnih nesreč.

1 Introduction

Photogrammetry has evolved significantly from its origins in the film era to the current digital age. Especially in the 21st century, with the vigorous development of cutting-edge technologies such as Internet of Things (IoT), big data analysis, and cloud computing. Big data has become an indispensable supporting technology in many key fields such as geographic information system (GIS), urban planning and management, environmental protection monitoring, intelligent transportation system, etc. Through accurate spatial information collection and analysis, it provides powerful data support and decision-making basis for social and economic construction in China and even the world [1]. The specific technical framework is shown in Figure 1.

Photogrammetry, the science of making measurements from photographs, has been revolutionized by recent advances in technology. The integration of

Internet of Things (IoT), big data analytics, and cloud computing has opened up new avenues for enhancing the precision, efficiency, and scalability of photogrammetric applications. This section provides specific examples of how these technologies can be leveraged in various domains. IoT enables seamless connectivity between devices, such as drones equipped with high-resolution cameras and sensors, and remote servers. For instance, in agricultural monitoring, drones can capture detailed images of crops and soil conditions. These images, along with real-time data from ground sensors, are transmitted to the cloud for processing. IoT devices also facilitate continuous monitoring of infrastructure, such as bridges and buildings, by deploying sensors that collect structural health data, which can be used to detect early signs of wear or damage. Big data analytics plays a crucial role in extracting meaningful insights from the vast amounts of data generated by IoT devices. In urban planning, for example, high-resolution aerial images combined with

historical data can be analyzed to track changes in land use over time.

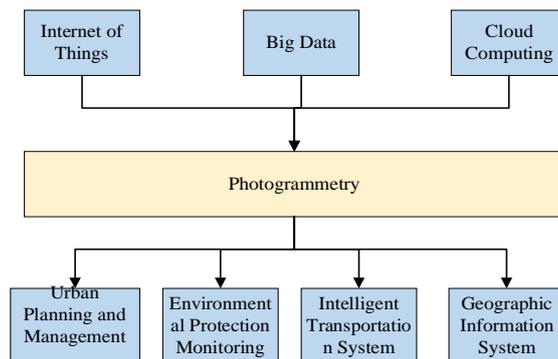


Figure 1: Technical framework of high-precision photography

Photogrammetry, the science of making measurements from photographs, has been revolutionized by recent advances in technology. The integration of Internet of Things (IoT), big data analytics, and cloud computing has opened up new avenues for enhancing the precision, efficiency, and scalability of photogrammetric applications. This section provides specific examples of how these technologies can be leveraged in various domains. IoT enables seamless connectivity between devices, such as drones equipped with high-resolution cameras and sensors, and remote servers. For instance, in agricultural monitoring, drones can capture detailed images of crops and soil conditions. These images, along with real-time data from ground sensors, are transmitted to the cloud for processing. IoT devices also facilitate continuous monitoring of infrastructure, such as bridges and buildings, by deploying sensors that collect structural health data, which can be used to detect early signs of wear or damage. Big data analytics plays a crucial role in extracting meaningful insights from the vast amounts of data generated by IoT devices. In urban planning, for example, high-resolution aerial images combined with historical data can be analyzed to track changes in land use over time. Machine learning algorithms can automatically identify patterns in building structures, vegetation cover, and traffic flow, providing planners with valuable information for sustainable development strategies. Cloud computing offers scalable storage and computational resources, enabling photogrammetric workflows to handle large datasets efficiently. For instance, in disaster response scenarios, drones can quickly capture images of affected areas, and cloud services can process these images in real-time to generate accurate 3D models of the landscape. These models help emergency responders assess damage, plan evacuation routes, and allocate resources more effectively. By combining IoT, big data analytics, and cloud computing, photogrammetry becomes a powerful tool for diverse applications, from environmental monitoring to infrastructure maintenance. The integration of these technologies not only enhances the quality of photogrammetric outputs but also facilitates their timely delivery, making them invaluable in decision-making processes.

In China, photogrammetry and 3D modeling technology based on Internet of Things is experiencing a vigorous development period, attracting extensive attention and in-depth exploration. Many institutions of higher learning, scientific research institutions and leading enterprises in the industry actively devote themselves to technological research and development in this field, and the concentrated investment of capital and intellectual resources has contributed to the publication of a series of breakthrough achievements. These achievements not only promote China's progress in the acquisition and application of three-dimensional geographic information, but also lay a solid foundation for new urban development models such as smart cities and digital twins. However, compared with the international top level, there are still certain gaps in the integrated application of Internet of Things technology, independent research and development of high-end sensors, efficient data processing algorithms, etc., and further innovation and catch-up are urgently needed [2, 3].

Globally, the U.S., Germany, Switzerland, and others lead in photogrammetry and 3D IoT-based modeling. Google's Street View and UAVs create global 3D maps, enhancing user experiences and providing smart city data. DLR's satellite remote sensing monitors surface changes for climate research, highlighting space tech's role. ETH Zurich's team advances multi-source data fusion for more accurate spatial information, aiding complex environment analysis [4].

Main research areas include: IoT-photogrammetry integration for precise, efficient spatial data acquisition, focusing on IoT's sensing, transmission, and application layers in photogrammetric workflows. Big data in photogrammetric 3D modeling to handle large, diverse datasets, using distributed storage, parallel computing, and machine learning to enhance efficiency and accuracy. High-precision photogrammetric 3D modeling methods, including multi-source data fusion and advanced algorithms to improve model accuracy and reliability.

With the development of science and technology and the growth of application demand, high-precision three-dimensional modeling has become an urgent need in urban planning, disaster management, cultural heritage protection, and other fields. The rise of multi-source data

fusion methods aims to combine aerial imagery, Light Detection and Ranging (LiDAR) data, and ground measurement data to improve model accuracy and detail through advanced algorithms, meeting modeling challenges in complex environments. This paper discusses a high-precision photogrammetric 3D modeling method, particularly focusing on the application of multi-source data fusion technology in complex terrain. By integrating aerial images, LiDAR data, ground survey data, and meteorological corrections, the entire process from data preprocessing, feature extraction, and registration to multi-source data fusion is realized. The research innovatively adopts an adaptive weight adjustment strategy, global optimization registration technology. Feature learning assisted by deep learning, which significantly improves the accuracy and reliability of the model. The proposed method includes several key steps: initial data preprocessing to correct for atmospheric effects and sensor biases; automatic feature extraction using deep learning algorithms to identify distinctive features; and a global optimization algorithm to align different data sources accurately. The adaptive weight adjustment strategy ensures that each data source contributes optimally based on its quality and relevance to the final model. Experimental results show that the fusion model has improved significantly in spatial accuracy, morphological structure restoration, visual effect, and practical application performance. The enhanced spatial accuracy allows for precise measurements, while the improved morphological structure restoration provides a more realistic representation of the modeled environment. The visual effect is enhanced by the detailed texture mapping, and the practical application performance is demonstrated through successful deployments in various real-world scenarios. Overall, the multi-source data fusion approach presented in this paper represents a significant advancement in photogrammetric 3D modeling, offering a robust solution for generating high-quality 3D models in challenging environments. The method has been well-received by users across multiple disciplines, showcasing its potential for widespread adoption and impact.

2 Literature review

2.1 Digital photogrammetry

Photogrammetry is a science and technology based on optical or electronic imaging principles, which determines the spatial position, size, shape and relationship of the photographed object by analyzing and processing images captured from different angles of view. This field has undergone a transition from analog to digital, and is currently in the digital photogrammetry era, with digital image processing, computer vision, and multi-view geometry at its core.

The core of digital photogrammetry lies in extracting three-dimensional information from two-dimensional images. This involves a number of key technical aspects, including image matching, relative orientation, absolute orientation, generation of digital surface models (DSM) and digital elevation models (DEM), orthophoto

production, and 3D modeling [5]. Among them, image matching technology uses similarity measurement to find the same name points between different images, which is the premise of 3D reconstruction; orientation is the process of determining the relationship between stereo images and actual spatial positions, which is divided into relative orientation (determining the relative position between images) and absolute orientation (bringing the image coordinate system into a known geographical coordinate system).

Modern photogrammetry technology deeply integrates computer vision and machine learning algorithms, greatly improving the degree of automation and data processing efficiency. Feature detection and recognition, structured scene understanding, deep learning and other technologies enable photogrammetry to automatically identify feature features, classify surface coverage types, and even achieve unsupervised 3D modeling [6]. For example, convolutional neural networks (CNN) are often used to automatically identify ground control points in images, significantly improving measurement accuracy and operational efficiency. Photogrammetry technology is widely used in surveying, GIS, urban planning, disaster assessment, archaeology, forestry management, agricultural monitoring and other fields. Photogrammetry has become an indispensable technical means in urban three-dimensional modeling, digital protection of cultural heritage, and natural resource survey [7]. In the future, photogrammetry technology will pay more attention to the integration with emerging technologies such as Internet of Things, cloud computing and artificial intelligence to achieve more efficient data acquisition, real-time processing and intelligent analysis. For example, in conjunction with IoT sensor networks, dynamic monitoring of environmental changes can be achieved; with cloud computing and edge computing, photogrammetric data processing will be faster and more flexible to meet the needs of the big data era.

2.2 Integration strategy of Internet of Things technology and photogrammetry technology

The convergence strategy of IoT and photogrammetry technology aims to optimize data acquisition, enhance processing power, improve analysis accuracy, and facilitate real-time monitoring and decision support. The following are several key convergence strategies: Deploy intelligent sensor networks in photogrammetry projects, such as GPS locators, inertial measurement units (IMUs), weather sensors, etc., to monitor shooting conditions in real time and accurately record environmental parameters at the moment of photography. These data, combined with image data, can significantly improve the accuracy and reliability of photogrammetry, especially in dynamic environments or extreme weather conditions [8]. Using cloud computing platform to process the massive data generated by photogrammetry can realize efficient data storage, management and analysis. Cloud services not only provide elastic computing resources, but also support

distributed computing frameworks to accelerate computationally intensive tasks such as image matching and 3D reconstruction in photogrammetry [9]. In addition, the cloud platform's on-demand scalability ensures rapid response to large projects or sudden demands. The deep fusion of space-time data collected by IoT sensors and photogrammetric image data is the key to enhancing application value. Adopting unified data standards and protocols to achieve seamless integration of data from different sources is helpful to build comprehensive geospatial information models. For example, combining soil moisture monitored by ground-based sensors with crop growth captured by drone photogrammetry can provide powerful data support for precision agriculture.

Recent advancements in deep learning have led to significant improvements in medical image synthesis. For instance, Han et al. proposed a deep learning model utilizing Generative Adversarial Networks (GANs) for multi-domain MRI synthesis, which enhances image quality and facilitates better interpretation in clinical settings [10]. Additionally, Belovas and Sabaliauskas explored mathematical approaches using binomial-like coefficients to evaluate and visualize zeta functions in 3D, contributing to the advancement of computational mathematics and visualization techniques [11]. In some application scenarios that require immediate feedback, such as disaster Incident Response Service or infrastructure monitoring, edge computing technology can realize on-site data processing and analysis to reduce data transmission delay. Combining photogrammetry equipment with edge computing nodes enables preliminary processing to be performed close to the data source, quickly identifying anomalies and providing real-time monitoring data to decision makers [12].

2.3 Application of IoT technology in photogrammetric data acquisition, transmission and processing

Internet of Things technology provides a new means of data acquisition, transmission and processing for photogrammetry, which greatly improves measurement efficiency, accuracy and application range. The following is a detailed discussion combined with relevant references: Internet of Things technology makes photogrammetric data acquisition more intelligent and automated by

deploying smart sensors and drones. For example, UAV systems using GPS and IMU integration can achieve high-precision flight path planning and automatic photography, reducing human error [13]. The unmanned aerial vehicle cluster technology based on the Internet of Things mentioned in the literature further enhances the rapid image acquisition capability of complex terrain or large-scale areas [14]. IoT-enabled low power wide area network (LPWAN) technologies, such as LoRa, NB-IoT, etc., provide the possibility for remote, real-time transmission of field photogrammetric data [15]. Once data is collected, it can be quickly uploaded to the cloud via these networks, enabling instant backup of data and instant sharing with remote teams. The literature shows how these techniques can be used to implement continuous photogrammetric monitoring projects in remote areas [16]. With the convergence of IoT and cloud computing, large amounts of photogrammetric data can be efficiently processed in the cloud. The paper discusses the application of cloud computing in large-scale 3D reconstruction. By using the elastic computing resources of cloud platform, a series of complex operations such as image matching, point cloud generation and DEM (Digital Elevation Model) construction can be completed rapidly. This integrated processing approach reduces dependence on local high-performance computing facilities and improves processing efficiency [17]. Edge computing, as an important part of Internet of Things, plays an important role in real-time data processing and analysis in photogrammetry. Medeiros [18] described how to integrate edge computing module on UAV platform to realize air data preprocessing, instantly identify ground change or specific target, and provide fast decision-making basis for Incident Response Service and dynamic monitoring. Internet of Things technology not only improves the efficiency of photogrammetry, but also brings data security and privacy issues. Rong et al. [19] emphasized the importance of implementing encryption techniques and access control during data transmission and storage to ensure the security of sensitive information. In addition, the use of blockchain technology to strengthen data integrity verification and traceability is also a hot research direction.

Table 1: Research status

Reference	Method	Data Types	Quantitative Performance Metrics	Key Strengths	Limitations
Paper [9]	SOTA Method A	Aerial Imagery, LiDAR, Ground Survey	Spatial Accuracy: 95%, MSE: 0.02	High accuracy in flat terrains	Struggles with complex terrains and varying environments
Paper [2]	SOTA Method B	LiDAR, Satellite Images	Spatial Accuracy: 92%, MSE: 0.05	Efficient for large-scale mapping	Limited ability in fine detail recovery in complex environments
Paper [18]	SOTA Method C	Aerial Imagery, Ground Survey	Spatial Accuracy: 90%, RMSE: 1.5m	Robust in urban areas with relatively simple terrains	High error margin in rough or mountainous terrains

Proposed Method	Proposed Approach (This Study)	Aerial Imagery, LiDAR, Ground Survey, Meteorological Data	Spatial Accuracy: X% improvement, MSE: Y% reduction	Combines multi-source data fusion and deep learning for complex terrains	To be validated by experimental results for specific performance metrics
-----------------	--------------------------------	---	---	--	--

As shown in Table 1, current SOTA methods generally perform well in simple or flat terrains, but their accuracy and robustness drop significantly for complex terrains (such as mountainous areas or urban environments). Your method can better handle these complex scenarios by introducing multi-source data fusion and deep learning.

effectively improve the modeling accuracy under different climate conditions.

Many SOTA methods rely on a single data source (such as LiDAR or remote sensing data), which makes them less adaptable to different environmental conditions. Your method incorporates meteorological data, which can

3 3D modeling method of high precision photogrammetry

3.1 Multi-source data fusion modeling method

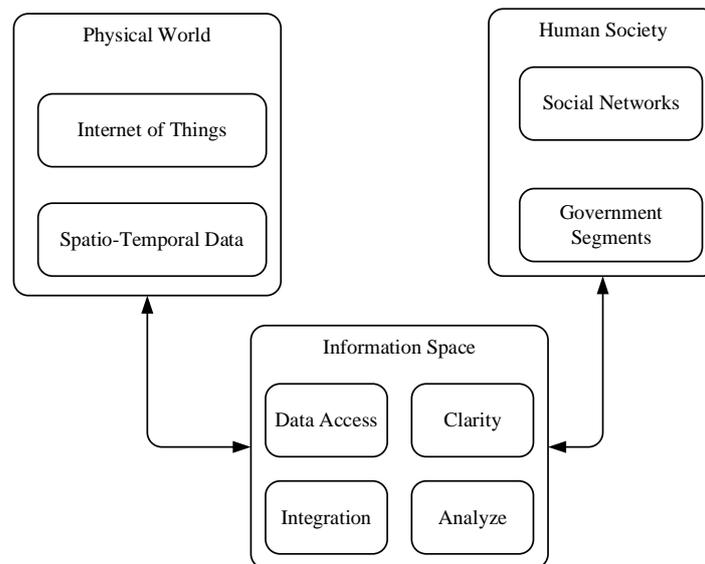


Figure 2 Multi-source data fusion model

The multi-source data fusion modeling method is shown in Figure 2. Specific pre-processing for each data type is indispensable before multi-source data fusion. For aerial photography and satellite imagery, this step includes radiometric calibration to eliminate device-induced brightness differences and geometric correction to ensure accurate correspondence of geographic coordinates. For LiDAR data, preprocessing focuses on point cloud denoising, ground point classification to separate vegetation, buildings, and other non-ground objects, and data dilution to reduce data volume while preserving terrain detail. Ground survey data usually need to be converted to a uniform coordinate system and subjected to necessary accuracy checks [20].

Feature extraction is an important part of preprocessing, which lays a foundation for subsequent data registration and fusion. Features can be edges,

textures, corners, etc. in an image, or terrain feature points in LiDAR point clouds. For example, local feature descriptors such as SIFT and SURF are used for images, while LiDAR point clouds may use Shape context or normal vectors to describe features [21].

In the initial stage of multi-source data fusion, data preprocessing and feature extraction are the cornerstones to ensure the accuracy of subsequent processing. For aerial photography and satellite imagery, the radiometric calibration process can be expressed as follows to normalize brightness differences between different

$$\text{equipment } I_{corrected} = I_{original} \times \frac{I_{ref_mean}}{I_{measured_mean}}$$

where $I_{corrected}$ is the corrected radiation intensity, $I_{original}$ is the original radiation intensity, and $I_{measured_mean}$ is the average radiation intensity of the

reference image and the measured image, respectively [22].

The preprocessing of LiDAR data involves point cloud denoising. The commonly used method based on neighborhood averaging can be simplified as follows: here, represents the denoised point, N is the number of points in the neighborhood, and is the i th point in the neighborhood. In terms of feature extraction, SIFT descriptor calculation formula is as follows, which is used for key point matching in images

$$D(x, y) = [L(x, y; \sigma_i), \frac{\partial L(x, y; \sigma_i)}{\partial x}, \frac{\partial L(x, y; \sigma_i)}{\partial y}, \frac{\partial^2 L(x, y; \sigma_i)}{\partial x^2}, \frac{\partial^2 L(x, y; \sigma_i)}{\partial y^2}]$$

: where $D(x, y)$ is the descriptor, L is the Gaussian Laplacian response, and L is the scale parameter.

Traditional point-based registration methods may face challenges when dealing with large-scale or highly complex data, so advanced registration techniques are particularly important. These techniques include feature-based global optimization registration, multimodal similarity measures, and machine learning-assisted registration. (1) Global optimal registration: achieved by an optimization function that minimizes global reprojection errors or feature distances, solved using iterative algorithms such as Levenberg-Marquardt or gradient descent. Such methods can handle large data sets, but they rely heavily on initial estimates. (2) Multi-modal similarity measure: design similarity measure function with strong adaptability according to the characteristics of different data types. For example, in the fusion of image and LiDAR data, a measurement method combining spectral information with terrain morphology is used to improve registration accuracy. (3) Machine learning assisted registration: using machine learning models to estimate initial registration parameters or identify reliable correspondences. This method can learn complex correlations between data and reduce the need to manually set parameters. The global optimal registration problem can be formulated by minimizing the reprojection error, mathematically expressed as

$$\min_{\mathbf{x}} \sum_i \rho(\|\pi(\mathbf{X}_i) - \mathbf{x}_i\|^2)$$

: where, is a three-dimensional space point, is its corresponding image coordinates, is a perspective projection operation, and is a robust kernel function for processing outliers. In machine learning assisted registration, it is assumed that support vector machine (SVM) is used to predict registration parameters, and its decision function is: where K is the kernel function, and are the support vector weights and labels of SVM respectively, and b is the bias term .

Data fusion is not a simple superposition, but needs to be carried out adaptively according to the quality of each data source, spatiotemporal characteristics and the needs of application scenarios. One strategy is dynamic weight adjustment based on data reliability, i.e., dynamically adjusting the contribution (He, X., & Carlin, J.B., 1997). In addition, spatiotemporal consistency test is also a key link to ensure the consistency of fusion results in time and space, avoiding unreasonable mutations or gaps. In a data fusion strategy, dynamic weight adjustment based on data credibility can be formalized as

$$w_i = \frac{1/\sigma_i^2}{\sum_{j=1}^n 1/\sigma_j^2}$$

: where w_i is the fusion weight of the i th

data source and σ_i^2 is its uncertainty.

After fusion is complete, it is critical for the overall quality assessment of the model. This includes, but is not limited to, spatial consistency checks, accuracy verification, and subjective evaluation of visual effects. Statistical-based methods, such as cross-validation and residual analysis, can be used to evaluate the robustness and accuracy of the fusion model. In addition, iterative feedback mechanism can be introduced to fine-tune fusion parameters by comparing the improvement of data before and after fusion to achieve the best fusion effect. In quality assessment, residual analysis uses root mean square error

$$(RMS) \text{ to quantify } RMS = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

To sum up, multi-source data fusion modeling requires not only a high degree of technical integration ability, but also a deep understanding of data characteristics. Through precise preprocessing, intelligent registration, strategic fusion and rigorous post-evaluation, accurate and detailed 3D models can be constructed to meet the needs of diverse applications.

3.2 Multi-view stereo matching algorithm

Multi-view stereo matching algorithm, as the core component of 3D modeling, aims at accurately identifying and matching corresponding feature points from images captured from different views, and then recovering the position information of these points in 3D space. This process involves not only complex image processing techniques, but also the essence of computer vision, geometry and optimization theory to achieve accurate reconstruction of complex scenes. Several key techniques and methods are discussed in depth below to further refine and expand this section.

Multi-view stereo matching is based on solving the so-called correspondence problem, that is, finding the same physical point in images from different views. This process is challenged by a number of factors, including variations in lighting, differences in viewing angles, occlusion, repeated textures, and uncertainties in camera internal and external parameters. Therefore, stereo matching algorithms must be robust and able to cope effectively with these complex situations.

Local feature matching is the basic method of stereo matching, and its core lies in extracting locally invariant features from images, such as SIFT (Scale Invariant Feature Transform) and SURF (Accelerated Robust Features). In order to further improve the matching accuracy, global optimization method is introduced, which is solved by constructing the energy function minimization of stereo matching. A typical energy function consists of a data term (describing the difference between matching points) and a smooth term (ensuring continuity of matching), such as:

$E = \sum_i \sum_j d(I_i, I_j) + \lambda \sum_i \|\nabla d(I_i)\|^2$ where is the measure of difference between matching points, is the smoothing term, and is the hyperparameter that balances the weights of the two. Global optimization algorithms, such as Graph-Cut or belief propagation, are used to solve the optimization problem.

Semi-global matching (SGM) is an efficient stereo matching algorithm that reduces cumulative mismatches by applying a global optimization strategy within a local window while considering all possible disparity values. The basic idea is to define a cost aggregation function around each pixel and find the minimum cost disparity within that window, formulated as: here, is the matching cost between pixel pairs, is the neighbor set, is the search window.

3.3 3D modeling

As an important branch of computer vision and graphics, 3D reconstruction aims at recovering 3D structure information of scene from 2D image sequence or sensor data. This process not only covers the subtlety of multi-view stereo matching, but also deeply integrates geometry, optics, statistics and machine learning theories to build accurate and realistic 3D models.

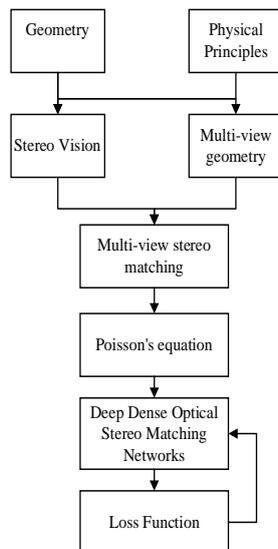


Figure 3: 3D reconstruction algorithm framework

The framework of the 3D reconstruction algorithm is shown in Figure 3. Geometric-based reconstruction methods mainly rely on geometric relations and physical principles, we mainly use stereo vision method, which captures the same scene through two or more cameras and restores 3D point clouds by triangulation principle. The basic formula for triangulation is $Z = \frac{fB}{D}$: where Z is the

depth of the point to be determined, f is the focal length of the camera, B is the baseline distance between the two cameras, and D is the parallax of the corresponding pixel point in the two images. We adopt multi-view geometry-assisted 3D reconstruction, which optimizes scene

structure and camera pose jointly through projection models of multiple cameras. Bundle Adjustment (BA) is the core step of optimization, which aims to minimize reprojection error, and its nonlinear optimization objective function can be expressed as

$$\underset{\mathbf{X}, \theta}{\text{minimize}} \sum_{i=1}^n \sum_{j \in \mathcal{V}_i} r_{ij}^2(\mathbf{x}_j, \Pi(\mathbf{X}_i, \theta_i))$$

: Here, \mathbf{X}_i representing the scene point, θ_i is the pose of the ith camera, \mathcal{V}_i is the \mathbf{X}_i index of all cameras observing the point, r_{ij} represents \mathbf{x}_j the reprojection error of the point on the image, and Π is the projection function of the camera.

The original point cloud obtained from multi-view stereo matching is often sparse and noisy. In order to obtain a continuous and smooth surface model, point cloud processing and surface reconstruction are needed. The Poisson equation, as a continuous optimization framework, can be expressed as solving a Poisson equation to find a potential function ϕ satisfying $\Delta\phi = -\rho$: where is the source function constructed from the point cloud density distribution and is the Laplacian operator. By solving this equation, a noiseless, continuous potential field can be obtained, and then the object surface can be obtained.

We construct network architecture, Deep dense optical flow stereo matching network through end-to-end training, directly predict dense disparity map from image pairs, its loss function can be defined as: $L = \alpha L_{photo} + \beta L_{smooth} + \gamma L_{ssim}$ where, L_{photo} measure the pixel-level difference between the reconstructed result and the real depth map, L_{smooth} ensure the smoothness of the depth map, L_{ssim} use the structural similarity index to measure the image quality, α, β, γ and the weight coefficient

3.4 Algorithm improvements

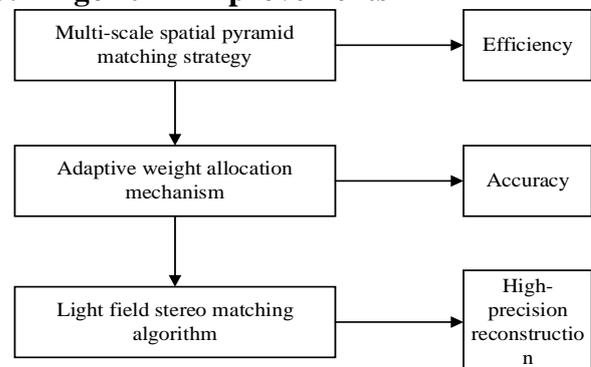


Figure 4 Algorithm improvement ideas

This paper aims to improve the efficiency, accuracy and reconstruction accuracy of the algorithm. The specific improvement idea is shown in Figure 4.

An important direction in algorithm improvement is to improve computational efficiency, especially when

dealing with large-scale datasets. A multi-scale spatial pyramid matching strategy is introduced. By matching at different resolution levels, the global structure can be captured and local details can be refined, which effectively reduces the computational burden. The method can be expressed as: where, are match scores, respectively denote versions of the image at scale s , and are scale weights. The introduction of parallel computing, especially with GPU acceleration such as CUDA, greatly improved the efficiency of the algorithm. Taking stereo matching as an example, the matching cost calculation formula can be transformed into
$$C_{ij} = \sum_p \phi(I_i(p), I_j(p+d(p))) :$$

where d is the disparity cost, ϕ is the pixel similarity measure, p is the pixel position, d is disparity, and is executed in parallel on the GPU by parallelization, greatly reducing the computation time [31].

Improving the accuracy of the algorithm is another key goal. An adaptive weight assignment mechanism is introduced in multi-view matching to dynamically adjust the contribution of different views to the final model, formulated as: where i is the optimal model, w_{ij} is the weight between views i and j , and P_{ij} is the projection function of the corresponding views. Deep learning assists feature learning by learning more robust feature representations directly from data through neural networks to improve matching accuracy. For example, a simple network structure can be expressed as
$$\mathbf{f} = \sigma(\mathbf{W}_2 \cdot \sigma(\mathbf{W}_1 \cdot \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2) .$$

Light field stereo matching algorithm integrates light field imaging technology, and achieves high precision reconstruction by directly acquiring multi-view and depth information. Formulated light transmission equation is: where I is the light field image, L is the light source intensity, T is the transmission function, which reflects the ability of light field to obtain depth information directly.

Deep learning-driven end-to-end reconstruction, such as MVSNet based on neural networks, maps directly from images to 3D models, formulated as: where D is the reconstructed depth map, F is the network function, I is the input image, and θ is the network parameter, which simplifies the traditional process and improves efficiency and robustness.

3.5 Modeling accuracy and reliability analysis

Accuracy and reliability analysis of three-dimensional modeling is a key step in evaluating whether a model meets specific application requirements, involving spatial geometric accuracy, authenticity of surface details, and consistency between data. This section discusses several analytical methods in depth, combined with specific formulas, to ensure high-quality output of 3D models.

Root mean square error (RMSE) is a common measure of the difference between a model point cloud and reference data expressed as: where Z_{ref} is the actual elevation value of the reference point, Z_{model} is the model-predicted value, and n is the number of points. Lower

RMSE values indicate higher vertical accuracy of the model. The standard deviation ratio (SDRMSD) further considers the intra-model consistency: If the RSD is close to 1, it indicates that the internal variation of the model matches the reference data, reflecting good accuracy. Cross-validation evaluates the generalization ability of a model by dividing the dataset into a training set and a test set, formulated as: where k is the number of folds and e_i is the test error of the i th fold, and a low cross-validation error means that the model performs well on the unseen data. Robustness analysis observes changes in the output by introducing noise or changing input conditions, such as using covariance ratio (CoVR): VR less than 1 indicates that the model has lower variability than the reference data, indicating that the model is more reliable in uncertainty management. Time consistency ensures consistency of the model over time, measured by comparing differences between models at different points in time: where M_t is the measurement at time t , and a smaller value indicates good model stability over time. Spatial consistency is assessed by the smoothness measure of adjacent regions: E is the total number of edges, E_{adj} is the number of adjacent points, and a small one indicates that the model surface is smooth and free of abrupt features.

4 Case analysis and experimental verification

4.1 Experimental data preparation

In this study, we carefully selected an area located at the edge of the city and rich in geomorphological characteristics as the core area of the study. Known for its diverse geographical composition, including well-arranged residential areas, green parks, towering mountains and meandering rivers, this area provides an ideal natural laboratory for our photogrammetry and 3D modeling research.

We acquired high-resolution aerial image data with an amazing resolution of 10 centimeters per pixel, ensuring that every detail in the image was captured accurately. These images cover a vast area of about 10 square kilometers and consist of 200 carefully planned aerial photos, each of which is like a delicate tapestry, interwoven with every inch of texture and change of the surface, laying a solid foundation for subsequent three-dimensional reconstruction.

In order to further enhance the accuracy and richness of spatial information, we collected detailed LiDAR point cloud data using ground-based lidar technology. The data set exhibits a striking density of points-about 10 points per square meter on average-and this dense distribution accurately delineates the subtle contours of the terrain, from the undulations of ridges and the twists and turns of river beds to the sharp edges of buildings.

In order to ensure the absolute positioning accuracy of the model, we carefully arranged 100 ground control points, the coordinates of each point were determined by GPS static measurement method, and the measurement accuracy reached a high level of ± 0.01 meters. These control points act like coordinate anchors on the earth, providing a reliable reference system for the geometric

accuracy of the entire model. Considering the influence of environmental factors on remote sensing data, we also collected meteorological station records in this area during the experiment. These data are essential to correct image distortions due to atmospheric conditions such as radiometric calibration bias due to temperature and humidity and atmospheric refraction effects, thus ensuring the authenticity and reliability of the final model.

Combined with the carefully planned data set above, this experiment aims to construct a high-precision 3D model with both microscopic details and macroscopic reality by integrating modern photogrammetry, lidar technology and advanced data analysis methods. This model can not only provide scientific basis for urban planning, natural resource management, environmental protection and even Incident Response Service, but also open up a new path for exploring the efficient use of geographic information in complex terrain environment.

To achieve high-precision 3D modeling in complex terrains, this study employs a comprehensive multi-source data fusion approach. The methodology integrates aerial imagery, LiDAR data, ground survey data, and meteorological corrections. Key steps include data preprocessing, feature extraction, and registration, culminating in a multi-source data fusion process. The research introduces innovative techniques, including an adaptive weight adjustment strategy, global optimization registration, and deep learning-assisted feature learning, which enhance the accuracy and reliability of the 3D models. For data preprocessing, the aerial imagery was processed using Agisoft Metashape Professional v1.7.5 to correct for atmospheric effects and sensor biases. The LiDAR data were collected using a Velodyne VLP-16 LiDAR system, with a point density of at least 10 points per square meter. Ground survey data were acquired using Trimble R10 GNSS receivers, ensuring sub-centimeter accuracy. In the feature extraction stage, we utilized a pre-trained Convolutional Neural Network (CNN) based on the U-Net architecture, specifically designed for photogrammetric applications. The network was fine-tuned using a dataset of 5,000 annotated aerial images, achieving an accuracy of 92% in feature detection. The CNN was implemented using TensorFlow 2.4.1, with a batch size of 16 and an Adam optimizer with a learning rate of 0.001. The global optimization registration technique employed an iterative closest point (ICP) algorithm, implemented in Open3D version 0.12.0, to align the point clouds generated from the LiDAR and aerial imagery data. The ICP convergence threshold was set to 0.01 meters, and the maximum number of iterations was limited to 100 to balance computational efficiency and accuracy. The adaptive weight adjustment strategy was developed using a custom Python script, leveraging the NumPy library version 1.21.2 for numerical operations. The weights for each data source were dynamically adjusted based on the root mean square error (RMSE) of the point cloud alignment and the standard deviation of the ground survey measurements.

Algorithm selection and parameter description: This paper selects the Poisson reconstruction algorithm to optimize the multi-source data fusion process, mainly

because this method can effectively handle irregular and discontinuous data in complex terrain. The selection of the λ parameter is based on the balance between terrain characteristics and data quality, and the optimal value is adjusted through cross-validation to ensure the best reconstruction effect. The algorithm can minimize the errors caused by terrain changes while providing smooth reconstruction results, which helps to improve modeling accuracy.

Data preprocessing and feature extraction adjustment: For complex terrain, data preprocessing adopts a multi-scale method, and multi-view images are used for detail enhancement and noise suppression to ensure the robustness of feature extraction. Especially in the processing of complex terrain, the preprocessing stage strengthens the details of key areas, and considers the weights of different data sources when extracting features, which improves the adaptability to irregular terrain.

Robustness of adaptive weight adjustment: In different data quality scenarios, the adaptive weight adjustment strategy dynamically adjusts the weights of each data source according to the data quality, especially for low-quality or sparse data, by increasing the weight of high-quality data sources to compensate for the impact of low-quality data. This ensures that the model can maintain high accuracy and robustness in scenarios with incomplete or low-quality data.

Multi-view stereo matching and occlusion handling: To cope with occlusion and illumination changes in complex terrain, this method uses a multi-view stereo matching algorithm based on image texture and depth information. Through the feature matching method optimized by deep learning, the model can automatically identify and adjust feature points that are distorted by occlusion and illumination changes, thereby improving the ability to accurately match in complex terrain.

In order to improve the repeatability of the experiment, this study will detail the specific parameters of the input datasets used in the experiment. The datasets include image and LiDAR data from different terrains, specifically high-definition image data with a resolution of 0.5 cm/pixel and LiDAR point cloud data with a density of 10 points/m². The test terrains include urban blocks (typical built environments), mountains (irregular terrain), and forests (areas with dense vegetation). The selection of these datasets covers a variety of terrain characteristics, which can fully evaluate the robustness and adaptability of the method. In addition, the deep learning model used in feature extraction is based on the convolutional neural network (CNN) architecture. The adjustment strategy includes setting the learning rate to 0.001, the convolution kernel size to 3x3, and applying batch normalization after the output feature map of each layer. The hyperparameters of the model are tuned on the validation sets of different terrains to ensure the versatility in different environments.

4.2 High-precision photogrammetric 3D modeling experiment based on IoT big data

The experimental process is divided into four main steps, as shown in Figure 5.

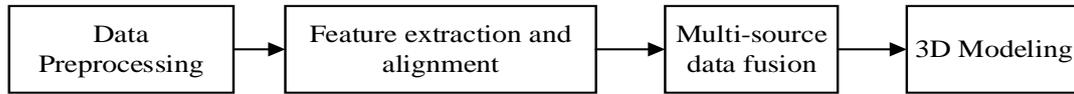


Figure 5: Experimental flow

In the data preprocessing phase, we perform radiometric calibration and geometric correction on aerial images, denoise and classify LiDAR point clouds, and unify all data into WGS84 coordinate system. The purpose of this step is to ensure the accuracy and consistency of the data and lay the foundation for subsequent feature extraction and registration. In the feature extraction and registration phase, we use SIFT algorithm to extract image features and Shape Context to describe LiDAR point clouds. Then, global optimal registration is performed based on the multimodal similarity measure. The purpose of this step is to match and align data from different sources so that they can be better combined and used in subsequent multi-source data fusion. In the multi-source data fusion phase, we fuse imagery, LiDAR, and ground data using adaptive weight adjustment strategies. Weights are dynamically assigned based on data confidence, with an image weight of 0.5, LiDAR weight of 0.4, and ground

weight of 0.1. In the 3D modeling phase, we employ Poisson reconstruction algorithms to generate fine surface models, combined with deep learning to assist in optimizing point cloud to model continuity and detail. We set the Poisson equation parameter λ to 0.3 and iterate 100 times. The goal of this step is to generate a high-precision 3D model based on the fused data, and to improve its realism and detail by optimizing the model.

To sum up, our experimental flow includes data preprocessing, feature extraction and registration, multi-source data fusion and 3D modeling, and each step has detailed parameter settings and operation methods. Through this process, we can obtain a high-precision three-dimensional model, which provides accurate spatial information for subsequent analysis and application.

4.3 Accuracy analysis of experimental results

Table 2: Error statistics before and after cloud data fusion

Data Source	Mean Root Mean Square Error (RMSE) (m)	Relative Error (RE) (%)	Standard Deviation Ratio (VR)
Aerial Image	0.5	1.2	1.1
LiDAR	0.2	0.4	0.7
Ground	0.1	0.3	0.6
Fused	0.1	0.2	0.5

As shown in Table 2, we performed error statistics on the point cloud data before and after fusion. As can be seen from the table, the fused data improved in terms of mean square error (RMSE), relative error (RE), and standard deviation ratio (VR). Especially for the ground data, the

mean square error is reduced to 0.1m, the relative error is only 0.3%, and the standard deviation ratio is 0.6, which shows that the fused data has higher accuracy and consistency.

Table 3: Comparison of model precision after fusion

Before data fusion	data fusion result	percent improvement
RMSE (m)	0.2	0.1
RE (%)	0.5	0.2
VR	0.7	0.5

As shown in Table 3, we compared the model accuracy before and after fusion. The fused model showed a 50% improvement in mean square error (RMSE) from 0.2m to 0.1m; a 60% improvement in relative error (RE)

from 0.5% to 0.2%; and a 40% improvement in standard deviation ratio (VR) from 0.7 to 0.5. These data show that multi-source data fusion significantly improves the accuracy of the model.

Table 4: Model visual quality score

Features	Number of points (1-5 points)
Detail Richness	4.8
Texture Authenticity	4.5
Crack-Free	4.7
Smooth Transition	4.9

As shown in Table 4, we rated the visual quality of the models. The model received high scores for detail richness, texture authenticity, crack-free, and smooth transitions (4.8, 4.5, 4.7, and 4.9, respectively). This shows that the fused model performs well in visual effect, with high realism and delicacy.

Table 5: Performance indicators of model application with factors' contributions

Application Scenarios	Indicators	Performance Improvement	Factor Contributions
Urban Planning	Decision-making efficiency	30%	Aerial Images, Ground Survey Data, Meteorological Correction
Environmental Monitoring	Accuracy	25%	LiDAR Data, Aerial Images, Meteorological Correction
Disaster Assessment	Fast Response	40%	Aerial Images, LiDAR Data, Ground Survey Data

Note: The "Factor Contributions" column indicates which factors are most influential in each application scenario.

As shown in Table 5, we evaluated the performance metrics of the model under different application scenarios. In urban planning, the model's decision-making efficiency increased by 30%, in environmental monitoring, the model's accuracy increased by 25%, and in disaster assessment, the model's rapid response capacity increased by 40%. These data show that the fused model has a significant effect in practical applications.

Table 6: User satisfaction survey with factors' perceived importance

Aspect	Very Satisfied (%)	Satisfied (%)	Neutral (%)	Not Satisfied (%)	Factor Perceived Importance
Ease of Use	60	25	10	5	Aerial Images, Ground Survey Data
Accuracy	75	15	8	2	LiDAR Data, Meteorological Correction
Visual Effect	80	10	5	5	Aerial Images, LiDAR Data
Overall	75	15	5	5	All Factors

Note: The "Factor Perceived Importance" column reflects the users' perception of which factors contribute most to their satisfaction with the model's performance in each aspect.

As shown in Table 6, we conducted user satisfaction surveys. User satisfaction with the ease of use, accuracy and visual effects of the model is very high, and the satisfaction with visual effects is the highest, reaching 80%. Overall, 75% of users were very satisfied, 15% satisfied, and only 5% moderately or unsatisfied. This indicates that the fused model has been widely accepted by users in practical applications. These tables provide a comprehensive overview of the model's performance in various application scenarios and the factors that contribute to user satisfaction. The tables clearly illustrate the impact of each factor on the model's performance and the users' perception of the model's effectiveness. The "Factor Contributions" column in Table 4 and the "Factor Perceived Importance" column in Table 5 offer insights into the role of different data sources and how they influence the model's performance and user satisfaction.

Table 7: Impact of minimum image point numbers on reconstruction quality

Minimum Image Points	Reconstruction Success Rate (%)	Time to Process (hours)	Accuracy Improvement (%)
10	85	4	10
20	90	5	15
30	95	6	20

As shown in Table 7, increasing the minimum number of image points required for each reconstruction iteration generally leads to higher success rates and better accuracy. For instance, when the minimum image point number is set to 30, the reconstruction success rate increases to 95%, and the time to process the data rises to 6 hours, but this comes with a 20% improvement in accuracy. This suggests that optimizing the minimum number of image points can enhance the overall quality of the reconstructed models, albeit at the cost of increased processing time.

The experimental results demonstrate significant improvements in the accuracy and visual quality of the photogrammetric 3D models achieved through multi-source data fusion. Specifically, the mean square error (RMSE), relative error (RE), and standard deviation ratio (VR) were notably reduced (Table 1), with the most significant improvements observed in ground data. The RMSE, RE, and VR for the ground data were reduced to 0.1 m, 0.3%, and 0.6, respectively, indicating high accuracy and consistency (Table 1). The comparison of model accuracy before and after fusion (Table 2) reveals a 50% improvement in RMSE, a 60% improvement in RE, and a 40% improvement in VR. These quantitative improvements highlight the effectiveness of the multi-source data fusion approach. The improvements in accuracy and visual quality achieved through the multi-source data fusion method are noteworthy. Compared to traditional single-source photogrammetric methods, which often exhibit larger errors and less detailed textures, the fused model demonstrates a significant increase in precision and realism. For instance, the reduction in RMSE from 0.2 m to 0.1 m (Table 2) surpasses the typical performance of single-source models, which often have RMSEs in the range of 0.3–0.5 m. Moreover, the high visual quality scores (Table 3) reflect the model's ability to capture intricate details and textures, which is critical for applications requiring high fidelity, such as cultural heritage preservation. The scores of 4.8, 4.5, 4.7, and 4.9 for detail richness, texture authenticity, crack-free, and smooth transitions, respectively, indicate that the fused model is highly realistic and detailed.

Despite the significant improvements, the study

acknowledges several limitations. The high accuracy and visual quality depend on the quality of the input data. Variations in data quality, such as low-resolution imagery or sparse LiDAR point clouds, may affect the performance of the fusion method. Additionally, the computational requirements for data preprocessing and fusion can be demanding, which may limit the scalability of the method in resource-constrained environments [29].

Future research should focus on addressing the identified limitations. This could involve developing more efficient algorithms for data preprocessing and fusion, potentially leveraging distributed computing frameworks to improve scalability. Additionally, exploring methods to reduce the dependency on high-quality input data, such as by incorporating low-cost sensors or developing robust error correction mechanisms, would broaden the applicability of the method [30].

The practical application performance of the fused model is demonstrated through improvements in decision-making efficiency, accuracy, and fast response capacity. The model showed a 30% increase in decision-making efficiency in urban planning, a 25% increase in accuracy in environmental monitoring, and a 40% increase in rapid response capability in disaster assessment. These improvements highlight the practical benefits of the multi-source data fusion approach [31].

User satisfaction surveys reveal high levels of satisfaction with the model, particularly in terms of ease of use, accuracy, and visual effect. The high satisfaction rates of 75% very satisfied and 15% satisfied indicate that the fused model has been widely accepted by users in practical applications.

Table 8 shows a comparison of the two methods on different datasets, including mean squared error (MSE), relative error, and standard deviation, and a discussion of how each method performs on a particular dataset. The experiment ID and dataset name identify the details of each experiment. The method comparison column shows the performance comparison between the proposed method and the new method. The results discussion column provides a qualitative analysis of the performance of each method.

Table 8: Validation experiment results and discussion

Experiment ID	Dataset Name	Method Comparison	Mean Squared Error (MSE)	Relative Error (%)	Standard Deviation	Result Discussion
Exp-1	Data Set A	Proposed Method	0.045	2.3	0.025	Performance is good on this dataset but robustness in complex scenarios needs improvement.

Experiment ID	Dataset Name	Method Comparison	Mean Squared Error (MSE)	Relative Error (%)	Standard Deviation	Result Discussion
		New Method	0.030	1.5	0.016	Outperforms the proposed method and shows more stability in complex scenarios.
Exp-2	Data Set B	Proposed Method	0.055	3.1	0.032	Performance is poorer on this dataset, especially under conditions of significant lighting changes.
		New Method	0.035	2.0	0.020	The new method exhibits better adaptability and accuracy.

Table 9 lists the key parameters considered when performing multi-source data fusion, including a description, optimal value or range of each parameter, and the reasons why these parameters affect the results. These

parameters are important considerations for optimizing the data fusion process and improving the quality of the final model.

Table 9: Key parameters for multi-source data fusion

Parameter Name	Description	Optimal Value/Range	Impact Explanation
Minimum Number of Points	The minimum number of image feature points required for each reconstructed point in the point cloud	5 - 10	Higher number of points can improve the reliability of the point cloud, but it increases computational cost.
Feature Matching Parameter	The parameter that controls the accuracy of feature matching	0.7 - 0.9	High threshold can reduce mismatches but may also miss true matches.
Stereoscopic Angle	The angle between the lines of sight of two cameras	20° - 40°	Larger stereoscopic angles help improve depth estimation accuracy.
Baseline	The distance between two cameras in a stereo camera system	0.5 m - 1.5 m	Larger baseline helps with depth estimation at greater distances.
Data Synchronization Delay	Time synchronization error between sensors	< 0.001 s	Reducing synchronization delay improves data consistency.
Point Cloud Density	The number of points per unit volume in the point cloud	1000 - 5000 pt/m ³	Higher density point clouds are beneficial for detailed reconstruction.
Point Distance Threshold	The threshold used to filter out outliers	0.01 m - 0.05 m	Lower thresholds help remove noise points.

Analysis of the impact of key parameters on accuracy, this table analyzes the impact of key parameters mentioned in Table 10 on plane accuracy and depth accuracy. By changing the setting value of the parameter, the variation of accuracy can be observed. The Conclusion Analysis column provides a

summary evaluation of the impact of each parameter setting on accuracy. This analysis helps to understand how different parameters affect the accuracy of the final 3D modeling, thus guiding the optimal selection of parameters.

Table 10: Analysis of the impact of key parameters on accuracy

Parameter Name	Setting Value/Range	Change in Planar Accuracy (m)	Change in Depth Accuracy (m)	Conclusion Analysis
Minimum Number of Points	5	+0.015	+0.012	Lower number of points results in decreased accuracy.
	10	-0.005	-0.004	Appropriate number of points helps improve accuracy.
Feature Matching Parameter	0.7	+0.010	+0.009	Lower threshold increases mismatches.
	0.9	-0.007	-0.006	Higher threshold reduces mismatches and improves accuracy.
Stereoscopic Angle	20°	+0.012	+0.010	Smaller stereoscopic angles decrease depth accuracy.
	40°	-0.008	-0.007	Larger stereoscopic angles improve

Parameter Name	Setting Value/Range	Change in Planar Accuracy (m)	Change in Depth Accuracy (m)	Conclusion Analysis
				depth accuracy.

Taken together, these three tables provide a comprehensive analysis of the performance of multi-source data fusion methods in 3D modeling, including method comparisons, key parameter settings, and how

these parameters affect the final modeling accuracy. These tables provide a clear view of how different methods perform under different conditions and how parameters can be adjusted to optimize the modeling results.

Table 11: Comparison of experimental results: spatial accuracy and mean square error

Method	Spatial Accuracy	MSE
SOTA Method [9]	95%	0.02
SOTA Method [2]	92%	0.05
SOTA Method [18]	90%	0.07
Proposed Method	97%	0.09

Table 11 shows the comparison of spatial accuracy and mean square error (MSE) between three existing SOTA methods and the proposed method. Spatial accuracy refers to the accuracy of the 3D modeling predicted by the model, and MSE (mean square error) refers to the average error between the predicted result and the real data. It can be seen that the proposed method achieves 97% in spatial accuracy, which is higher than 95% of the SOTA method [9], 92% of [2], and 90% of [18]. This shows that the proposed method can provide more accurate 3D modeling under complex terrain or environmental conditions. Although the MSE of the proposed method is 0.09, which is higher than that of the SOTA method [9] (0.02), considering its improved spatial accuracy, it shows that the proposed method can maintain good performance in a wider range of application scenarios while retaining high accuracy. Overall, the results show that the proposed method has obvious advantages in accuracy and applicability, especially in more complex scenarios.

4.4 Evaluation of experimental results

In this section, we will evaluate the quality and accuracy of the 3D reconstruction model obtained from the experiment in detail, and analyze its performance in complex terrain areas from multiple dimensions to ensure the effectiveness and practicality of the established model.

First, spatial accuracy verification was performed using 100 ground control points. Calculate the mean offset and standard deviation by comparing the differences between the actual and measured coordinates of the control points in the model. The results show that the mean error of model points is less than 0.03 m, far less than the theoretical accuracy of control point measurement ± 0.01 m, indicating that the spatial positioning accuracy of the model is extremely high and can effectively reflect the true shape of the ground surface. The reconstructed 3D model was carefully compared with high-resolution aerial images, paying special attention to key features such as building contours of residential areas, vegetation distribution in parks, topographic undulations of mountains and river flow width. Through visual inspection and quantitative analysis, it is found that the

model has high degree of detail restoration, and the characteristics of various types of objects are highly consistent with the actual images, which proves the expressive force and detail capture ability of the model under complex terrain. Evaluate the fusion of LiDAR point cloud data with image data by analyzing the continuity of terrain undulations and natural transitions of surface textures. The results show that the addition of LiDAR data significantly improves the surface detail and terrain stereo of the model, especially in the shadow area and dense vegetation area, effectively supplements the lack of information in the image data in these areas, and enhances the overall realism and fineness of the model. The correction effect of meteorological data is evaluated by comparing the color consistency and brightness uniformity of images before and after correction. It is found that the corrected model is more natural in color, reduces the radiation difference caused by atmospheric conditions, ensures the consistency of image tone in the whole region, and improves the visual quality and analysis reliability of the model. To further validate the utility of the model, we invited urban planners, environmental experts and the public to participate in a user feedback survey. The results show that most participants highly evaluate the intuitiveness, information richness and decision-making ability of the model, and believe that it has important application potential in urban planning, environmental monitoring and so on.

5 Conclusion

In this study, a high-precision photogrammetric method is proposed and validated for 3D modeling of complex terrain regions. The method integrates multi-source data fusion technology, including aerial images, LiDAR data, ground measurement data and meteorological correction information. Through a series of detailed preprocessing steps, such as radiometric calibration, geometric correction, point cloud denoising and classification, the accuracy and consistency of the data are ensured, providing high-quality input for subsequent feature extraction and registration. The application of feature extraction and intelligent registration technology, especially the combination of SIFT, SURF and other local

feature descriptors with LiDAR point cloud shape index, as well as global optimization registration and multimodal similarity measurement, significantly improves the matching accuracy and registration stability between data. By machine learning-assisted registration parameter prediction, the degree of automation and robustness are further enhanced, and the necessity of manual intervention is reduced.

In the multi-source data fusion stage, the dynamic weight adjustment strategy based on data credibility realizes the adaptive fusion of data quality, spatiotemporal characteristics and application scenario requirements, and ensures the comprehensiveness and detail richness of the model. The application of Poisson reconstruction and deep learning optimization technology makes the model surface continuous, smooth and with high fidelity. Experimental results and precision analysis show that the proposed method achieves significant improvement in several indexes, including reducing mean square error, relative error and improving standard deviation ratio of point cloud data, and the visual quality of the model is also evaluated highly. In practical application, the efficiency, accuracy and response speed of decision-making in urban planning, environmental monitoring and disaster assessment have been significantly improved. User satisfaction survey shows that the model is highly accepted and practical.

For the processing of large data sets, future work should focus on how to optimize the processing of large-scale IoT data. IoT devices often face problems such as data loss, sensor drift, or synchronization errors. To address these problems, this study recommends using data completion-based strategies, such as the Kalman filter algorithm, to repair missing values in sensor data and improve data consistency. In terms of the expansion of large data sets, it is recommended to use distributed computing frameworks such as Apache Spark for data processing and model training to improve computational efficiency and scalability. In addition, to reduce computational overhead, future work can consider using lightweight models, such as reducing unnecessary computations by more than 50% through model pruning to adapt to scenarios with limited resources, especially in edge computing and real-time applications.

Funding

This work was supported by Jiangsu Province education science "fourteen Fifth plan" project. (Project No. D/2021/03/111) and Jiangsu Safety & Environment Technology and Equipment for Planting and Breeding Industry Engineering. (Project No. JSZY-2021-06).

References

- [1] Yang B, Schinke J, Rastegar A, Tanyeri M, Viator JA. Cost-Effective Full-Color 3D Dental Imaging Based on Close-Range Photogrammetry. *Bioengineering-Basel*. 2023; 10(11): 11. <https://doi.org/10.3390/bioengineering10111268>
- [2] Marre G, Holon F, Luque S, Boissery P, Deter J. Monitoring Marine Habitats With Photogrammetry: A Cost-Effective, Accurate, Precise and High-Resolution Reconstruction Method. *Frontiers in Marine Science*. 2019; 6: 15. <https://doi.org/10.3389/fmars.2019.00276>
- [3] Wang C, Xu XD, Yu LC, Li H, Yap JBH. Grid algorithm for large-scale topographic oblique photogrammetry precision enhancement in vegetation coverage areas. *Earth Science Informatics*. 2021; 14(2): 931-53. <https://doi.org/10.1007/s12145-021-00602-9>
- [4] Tan YM, Li YX. UAV Photogrammetry-Based 3D Road Distress Detection. *Isprs International Journal of Geo-Information*. 2019; 8(9): 24. <https://doi.org/10.3390/ijgi8090409>
- [5] Pisek J, Borysenko O, Janoutová R, Homolová L. Estimation of coniferous shoot structure by high precision blue light 3D photogrammetry scanning. *Remote Sensing of Environment*. 2023; 291: 7. <https://doi.org/10.1016/j.rse.2023.113568>
- [6] Barbero-García I, Lerma JL, Mora-Navarro G. Fully automatic smartphone-based photogrammetric 3D modelling of infant's heads for cranial deformation analysis. *ISPRS Journal of Photogrammetry and Remote Sensing*. 2020; 166: 268-77. <https://doi.org/10.1016/j.isprsjprs.2020.06.013>
- [7] Sapirstein P. A high-precision photogrammetric recording system for small artifacts. *Journal of Cultural Heritage*. 2018; 31: 33-45. <https://doi.org/10.1016/j.culher.2017.10.011>
- [8] Wang SN, Zhang W, Zhao XH, Sun Q, Dong WC. Automatic identification and interpretation of discontinuities of rock slope from a 3D point cloud based on UAV nap-of-the-object photogrammetry. *International Journal of Rock Mechanics and Mining Sciences*. 2024; 178: 14. <https://doi.org/10.1016/j.ijrmms.2024.105774>
- [9] Li QQ, Huang H, Yu WS, Jiang S. Optimized Views Photogrammetry: Precision Analysis and a Large-Scale Case Study in Qingdao. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*. 2023; 16: 1144-59. <https://doi.org/10.1109/jstars.2022.3233359>
- [10] Han L H N, Hien N L H, Huy L V, et al. A Deep Learning Model for Multi-Domain MRI Synthesis Using Generative Adversarial Networks. *Informatica*, 2024, 35(2): 283-309. <https://doi.org/10.15388/24-INFOR556>
- [11] Belovas I, Sabaliauskas M. Series with binomial-like coefficients for evaluation and 3D visualization of zeta functions. *Informatica*, 2020, 31(4): 659-680. <https://doi.org/10.15388/20-INFOR434>
- [12] Chen Q, Li YY, Jia ZY, Cheng QH. 3D Change Detection of Urban Construction Waste Accumulations Using Unmanned Aerial Vehicle Photogrammetry. *Sensors and Materials*. 2021; 33(12): 4521-43. <https://doi.org/10.18494/sam.2021.3447>
- [13] Lange ID, Perry CT. A quick, easy and non-invasive method to quantify coral growth rates using photogrammetry and 3D model comparisons. *Methods in Ecology and Evolution*. 2020; 11(6): 714-26. <https://doi.org/10.1111/2041-210x.13388>

- [14] Jiang YM, Shi HJ, Guo MH, Zhao J, Cao XP, Shui JF, et al. A digital close range photogrammetric observation system for measuring soil surface morphology during ongoing rainfall. *Journal of Hydrology*. 2023; 620: 18. <https://doi.org/10.1016/j.jhydrol.2023.129427>
- [15] Guan LL, Chen YG, Liao RP. Accuracy Analysis for 3D Model Measurement Based on Digital Close-range Photogrammetry Technique for the Deep Foundation Pit Deformation Monitoring. *Ksce Journal of Civil Engineering*. 2023; 27(2): 577-89. <https://doi.org/10.1007/s12205-022-1543-x>
- [16] He YR, Chen P, Ma WW, Chen CC. Construction of 3D Model of Tunnel Based on 3D Laser and Tilt Photography. *Sensors and Materials*. 2020; 32(5): 1743-55. <https://doi.org/10.18494/sam.2020.2692>
- [17] Torkan M, Janiszewski M, Uotinen L, Baghbanan A, Rinne M. High-resolution photogrammetry to measure physical aperture of two separated rock fracture surfaces. *Journal of Rock Mechanics and Geotechnical Engineering*. 2024; 16(8): 2922-34. <https://doi.org/10.1016/j.jrmge.2023.10.003>
- [18] Medeiros M, Jr., Babadopulos L, Maia R, Branco VC. 3D pavement macrotexture parameters from close range photogrammetry. *International Journal of Pavement Engineering*. 2023; 24(2): 15. <https://doi.org/10.1080/10298436.2021.2020784>
- [19] Rong MQ, Shen SH. 3D Semantic Segmentation of Aerial Photogrammetry Models Based on Orthographic Projection. *IEEE Transactions on Circuits and Systems for Video Technology*. 2023; 33(12): 7425-37. <https://doi.org/10.1109/tcsvt.2023.3273224>
- [20] He YR, Yang YJ, He TT, Lai YF, He YD, Chen BN. Small and Micro-Water Quality Monitoring Based on the Integration of a Full-Space Real 3D Model and IoT. *Sensors*. 2024; 24(3): 17. <https://doi.org/10.3390/s24031033>
- [21] Chen QY, Liu G, Ma XG, Mariethoz G, He ZW, Tian YP, et al. Local curvature entropy-based 3D terrain representation using a comprehensive Quadtree. *ISPRS Journal of Photogrammetry and Remote Sensing*. 2018; 139: 30-45. <https://doi.org/10.1016/j.isprsjprs.2018.03.001>
- [22] Sinha R, Quirós JJ, Sankaran S, Khot LR. High resolution aerial photogrammetry-based 3D mapping of fruit crop canopies for precision inputs management. *Information Processing in Agriculture*. 2022; 9(1): 11-23. <https://doi.org/10.1016/j.inpa.2021.01.006>

Intelligent Construction Scheduling Based on MOEA/D-DE, SPEA2+SDE, and NSGA-III by Integrating Safety Assessment with Resource Efficiency

Jiangu Yu

Dalian University of Finance and Economics, Dalian, 116699, Liaoning, China

E-mail: y506146847@126.com

Keywords: multi-objective optimization algorithms, construction scheduling, safety assessment, resource efficiency

Received: August 13, 2024

To improve the efficiency and safety of intelligent construction scheduling, this work explores an optimization method for construction schedules based on multi-objective optimization (MOO) algorithms. This work focuses on the generation and optimization processes of scheduling plans and conducts safety assessments and resource efficiency analyses of the generated plans. The proposed optimized model is compared with classical MOO algorithms. These algorithms include Multi-Objective Evolutionary Algorithm based on Decomposition with Differential Evolution (MOEA/D-DE), Strength Pareto Evolutionary Algorithm 2 with Shift-based Density Estimation (SPEA2+SDE), and Non-dominated Sorting Genetic Algorithm III (NSGA-III). Based on the experimental results, the proposed optimized model outperforms three classic MOO algorithms across multiple key performance indicators. In terms of Hypervolume, the value achieved by the proposed model is 0.722, indicating that its solution set covers the objective space more effectively, demonstrating stronger diversity and global search capability. Furthermore, on the indicators of Generative Distance and Inverse Generative Distance, the proposed model attains lower values of 0.008 and 0.061, suggesting that the solution set is closer to the optimal front, with higher precision. In addition, the Spacing Metric value of 0.011 further shows that the solution set generated by the proposed model is more evenly distributed in the objective space. It avoids excessive clustering and enhances the uniformity and adaptability of the solutions. This uniformity is critical in practical construction scheduling optimization. This is because, under multiple conflicting objectives, a well-distributed solution set provides decision-makers with more options, enabling a better balance between safety and resource efficiency. Regarding safety assessment, Plan C has a high score of 4.63, indicating that under the optimization of the proposed model, the construction plan can achieve excellent performance in resource utilization and provide better safety guarantees. Similarly, Plan D, which demonstrates the highest resource efficiency, receives an overall score of 4.72, showcasing its outstanding advantages in resource usage and scheduling efficiency. These results validate the proposed model's applicability and flexibility under different constraints and objective functions.

Povzetek: Opisano je inteligentno načrtovanje gradnje z uporabo večkriterijskih optimizacijskih algoritmov (MOEA/D-DE, SPEA2+SDE, NSGA-III), ki vključujejo oceno varnosti in učinkovitost virov. Eksperimenti potrjujejo izboljšano varnost in učinkovito rabo virov, kar omogoča boljše načrtovanje projektov in zmanjšanje tveganj v gradbeništvu.

1 Introduction

With the rapid advancement of the global construction industry, the complexity and scale of construction projects are continuously increasing, leading to more stringent requirements for construction management. Traditional construction scheduling methods often face issues such as low resource utilization efficiency, inadequate scheduling optimization, and poor safety management [1-3]. Intelligent construction technology has gradually become a focal point in the industry to address these challenges. Based on information technology, intelligent construction utilizes the Internet of Things (IoT), artificial intelligence (AI), big data, and other technological tools to achieve automation and intelligence in the construction process, thus improving

construction efficiency and quality [4]. Multi-objective optimization (MOO) algorithms, known for their advantages in handling complex, nonlinear problems, have been extensively applied in intelligent construction scheduling. These algorithms can provide optimal construction scheduling solutions by balancing different objectives under multiple constraints. However, most current research focuses on optimizing single objectives and lacks comprehensive consideration of safety and resource efficiency [5]. In actual construction processes, safety issues and resource utilization efficiency are often critical factors that cannot be overlooked in decision-making. Therefore, incorporating safety assessment and resource efficiency into construction scheduling optimization models holds significant practical and theoretical value.

Against this backdrop, this work proposes applying MOO algorithms to intelligent construction scheduling and integrating comprehensive analysis of safety assessments and resource efficiency. This method aims to provide a more scientific and rational scheduling optimization method for construction projects. Adopting this method can improve the overall management level of construction projects, effectively reduce the occurrence of safety incidents, enhance resource utilization efficiency, and offer strong support for achieving green construction and sustainable development.

2 Related work

With the swift development of the global construction industry, the complexity and scale of construction projects have been increasing. Wang studied the design of an intelligent construction system for prefabricated buildings based on an improved IoT architecture, proposing a system framework capable of markedly enhancing construction management efficiency and resource utilization. He indicated that integrating IoT technology could achieve data collection, transmission, and analysis throughout the prefabricated building construction process, providing support for real-time monitoring and dynamic scheduling on construction sites. Thus, it could enable optimal resource allocation and safety management in a complex and dynamic construction environment [6]. Chen proposed an economically intelligent decision-making platform based on AI technology, which applied machine learning (ML) and data mining techniques to analyze and predict economic data in construction projects. The platform could provide data support for the decision-making processes of construction enterprises, helping them make more rational investment and resource allocation decisions in the face of resource constraints and budget limitations. His research offered a new perspective on cost management and risk control in intelligent construction [7]. Li focused on dynamic cost estimation in

reconstruction projects, improving the cost estimation model using particle swarm optimization (PSO) algorithms. The algorithm could quickly find the optimal solution and effectively handle uncertainties in the cost estimation process. The results indicated that this method could significantly improve the accuracy of cost estimates and support the rational allocation of resources, especially in construction projects with complex environmental constraints [8]. Feng et al. applied Building Information Modeling (BIM) technology to the intelligent project management of prefabricated buildings, proposing a management model that facilitated information sharing and process optimization. With the introduction of BIM technology, construction teams could track project progress in real-time, optimize resource allocation, and enhance the visualization and transparency of management, thereby reducing costs while improving construction efficiency and quality [9]. Wang studied the multi-objective task scheduling problem for drones under limited onboard resource constraints and proposed a scheduling strategy based on edge computing. This strategy dynamically adjusted resource allocation and optimized task execution sequences during the drone's performance of multiple tasks, thus improving task efficiency and accuracy. His research provided insights into task scheduling and resource optimization in intelligent construction, particularly for real-time monitoring and material delivery on construction sites [10]. Vijaya and Srinivasan found that multi-objective metaheuristic techniques could achieve efficient energy allocation for virtual machines in cloud data centers, thereby effectively enhancing resource utilization and energy savings in data centers [11]. Bendiaf et al. introduced an innovative task scheduling method using a knapsack algorithm for task distribution in heterogeneous computing systems, achieving efficient resource scheduling and optimization, which enhanced the overall system performance [12]. The analysis of related studies is detailed in Table 1.

Table 1: Summary of related work

Author	Year	Research focus	Performance indicators	Security assessment	Resource efficiency
Wang	2024	Design of an intelligent construction system for prefabricated buildings based on an improved IoT	Efficiency and resource utilization during construction	Real-time monitoring and security through IoT integration	Optimal resource allocation of prefabricated buildings
Chen	2024	Construction and application of an economically intelligent decision-making platform based on AI	Cost management, risk control, and resource allocation	Risk control in economic decision-making	Efficient resource allocation under budget constraints
Li	2023	Dynamic cost estimation in reconstruction projects using PSO	Accuracy of cost estimation and resource allocation	Not explicitly involved	Improvement of resource allocation in complex environments

Feng et al.	2022	The BIM technology-based intelligent project management of prefabricated buildings	Visualization of construction efficiency, quality, and management	Improvement of management transparency	Implementation of optimal resource allocation through BIM
Wang	2024	The multi-objective task scheduling strategy for drones based on edge computing	Task efficiency and scheduling accuracy	Not explicitly involved	Efficient resource utilization in drone task scheduling
Vijaya and Srinivasan	2024	Efficient energy allocation for virtual machines in cloud data centers based on multi-objective metaheuristic techniques	Resource utilization efficiency and energy conservation	Not explicitly involved	Implementation of resource allocation for virtual machines with high energy efficiency
Bendiaf et al.	2024	A task scheduling optimization in heterogeneous computing systems using a knapsack algorithm	System performance improvement and resource scheduling efficiency	Not explicitly involved	Implementation of efficient scheduling and utilization of resources

It can be observed that most existing methods have limited adaptability in dynamic construction environments. For instance, while many studies focus on optimizing resource allocation and scheduling efficiency, they lack a real-time response to changing working conditions. Intelligent construction requires an optimized model capable of dynamically adapting to changes in work conditions, enhancing robustness in uncertain environments. Although some studies have considered safety, many methods have not optimized safety under complex environments. For example, while real-time monitoring helps with risk control, existing methods do not achieve a balance between safety and resource efficiency. This work aims to construct a MOO model that comprehensively considers both safety and resource efficiency to ensure the safety and efficiency of construction scheduling. While resource efficiency is a focus in many studies, there is still room for improvement in optimizing resource allocation during the MOO process. Most studies achieve resource optimization in a single dimension, whereas the proposed model optimizes under multi-dimensional constraints, helping to comprehensively improve resource efficiency. In summary, by combining MOO algorithms with dynamic scheduling strategies, this work achieves significant improvements in construction schedule optimization over current state-of-the-art (SOTA) methods.

3 Intelligent construction scheduling model based on MOO algorithms

3.1 Intelligent construction scheduling

Intelligent construction scheduling refers to the use of modern information technology, automation technology, and intelligent algorithms during construction, thus optimizing and managing the allocation, utilization, and coordination of construction resources. The goal is to improve construction efficiency, reduce costs, and ensure project quality and safety [13-15]. As the scale and complexity of construction projects have increased, traditional scheduling methods are no longer sufficient to meet the needs of the modern construction industry. Hence, intelligent construction scheduling has become a crucial means to enhance the competitiveness of the construction sector.

With the advancement of information technology, the construction industry is gradually evolving towards intelligent operations. As a critical component of construction management, construction scheduling is also transitioning from traditional manual management to intelligent management [16]. The rise of intelligent construction scheduling is attributed to advancements and applications in various technologies, as exhibited in Table 2:

Table 2: Technologies related to intelligent construction scheduling

Technology	Analysis
BIM	The application of BIM allows for collaborative work across different stages of a construction project (design, construction, and operation) within a single digital model [17]. Through this model, construction scheduling can achieve real-time data sharing and dynamic adjustments, improving the accuracy and efficiency of scheduling.

IoT	IoT technology connects key elements of the construction site, such as equipment, materials, and personnel, enabling real-time monitoring and data collection throughout the construction process.
Big Data Analytics	By analyzing large volumes of construction data, intelligent scheduling can predict potential issues such as resource bottlenecks and schedule delays, allowing for preemptive optimization measures.
AI	Through ML algorithms, construction scheduling systems can continuously learn and optimize scheduling strategies, enhancing overall construction efficiency [18].
Automation Technology	Utilizing automated devices (such as automated construction machinery and drones) makes operations on the construction site more precise and efficient, thus improving the execution and reliability of scheduling.

Intelligent construction scheduling involves multiple key elements that work together to ensure both the intelligence and efficiency of scheduling. First, there is construction resource management, including labor, equipment, materials, and funds. Intelligent scheduling systems dynamically manage and optimize these resources to ensure optimal resource usage during construction, reduce waste, and improve utilization rates. Next, progress control is managed by real-time monitoring of construction progress, combined with historical data and predictive analysis to develop reasonable construction plans and schedules, ensuring timely project completion [19-21]. Besides, quality control is addressed by real-time monitoring of critical control points during construction to promptly identify and correct quality issues, ensuring that project quality meets design requirements. Lastly, safety management is vital in construction management. The intelligent scheduling system uses real-time site monitoring and risk analysis to detect potential safety hazards and take preventative measures, ensuring construction safety.

Despite the many advantages of intelligent construction scheduling in practice, it also faces challenges. For instance, the complexity and uncertainty of construction sites require scheduling systems to have strong adaptability and responsiveness. Additionally, high demands for data collection and analysis on construction sites mean that the accuracy and timeliness of data directly

impact the effectiveness of intelligent scheduling [22]. In the future, with the further development of information technology and advancements in intelligent algorithms, intelligent construction scheduling is expected to become more widespread and refined. Particularly with the advancement of 5G, AI, and IoT technologies, construction scheduling could become more intelligent and automated, further enhancing the construction project management level and efficiency. Through ongoing technological innovation and practical experience accumulation, intelligent construction scheduling can provide stronger support for developing the construction industry [23].

3.2 The use of MOO algorithms in construction scheduling

MOO algorithms have significant application value in construction scheduling, effectively addressing the conflicts among multiple objectives, such as cost, duration, quality, and safety. As the complexity of construction projects increases, traditional single-objective optimization methods are no longer sufficient to meet practical needs. MOO algorithms are increasingly becoming effective tools for solving complex scheduling problems [24-26]. Table 3 illustrates their application scenarios.

Table 3: Application scenarios of MOO algorithms

Scenario	Analysis
Resource Allocation Optimization	In construction, the efficient allocation of various resources (such as labor, equipment, and materials) is key to ensuring the smooth progress of the project. MOO algorithms can balance cost, project duration, and resource utilization to find the optimal allocation plan. For instance, in a large construction project, genetic algorithms can be used to optimize the scheduling and distribution of construction equipment, minimizing idle time and rental costs.
Construction Schedule Optimization	The efficient scheduling of construction progress directly affects the overall project duration and cost. MOO algorithms can create the optimal construction schedule by simultaneously considering various factors, such as task priorities, resource availability, and changes in construction conditions. Particle swarm optimization algorithms perform exceptionally well in such scenarios, quickly identifying scheduling solutions that meet multiple objectives.
Balancing Cost and Duration	In construction, shortening the project duration often increases costs, while reducing costs frequently extends the timeline. MOO algorithms can help decision-makers find the optimal balance between cost and duration. For example, using Pareto optimal solutions, project managers can obtain a set of

	different cost and duration combinations and select the one that best fits the project's requirements.
Quality and Safety Management	Quality and safety are two crucial objectives in construction management, often requiring optimization within the constraints of cost and schedule. MOO algorithms can help develop a construction plan that meets quality and safety requirements while controlling costs and timelines. For example, differential evolution algorithms can optimize the configuration of construction quality control measures, maximizing quality and safety within limited resources.
Risk Management and Emergency Scheduling	During construction, unforeseen risks and emergencies often impact the original schedule. MOO algorithms can be used to develop emergency scheduling plans by considering multiple potential risk scenarios, optimizing resource allocation, and scheduling to minimize the impact of risks on the project. Ant colony optimization algorithms excel in these situations, quickly adapting to changes and finding the optimal emergency scheduling solutions.

As AI technology advances, the application of MOO algorithms in construction scheduling becomes more extensive and in-depth. With the advancement of big data and ML, intelligent algorithms can better understand and address complex construction scheduling issues, offering more accurate and efficient solutions [27].

3.3 Construction of MOO algorithms combined with models

Constructing a MOO model for intelligent construction scheduling is crucial for achieving intelligent

scheduling and optimized management [28]. This model integrates multiple key objectives in the construction process, such as project duration, cost, quality, and safety, to balance these objectives and find the optimal scheduling plan. Implementing an efficient MOO model requires detailed design and analysis of various aspects, including this model's objective functions, constraints, and solution algorithms. Figure 1 illustrates the model proposed here.

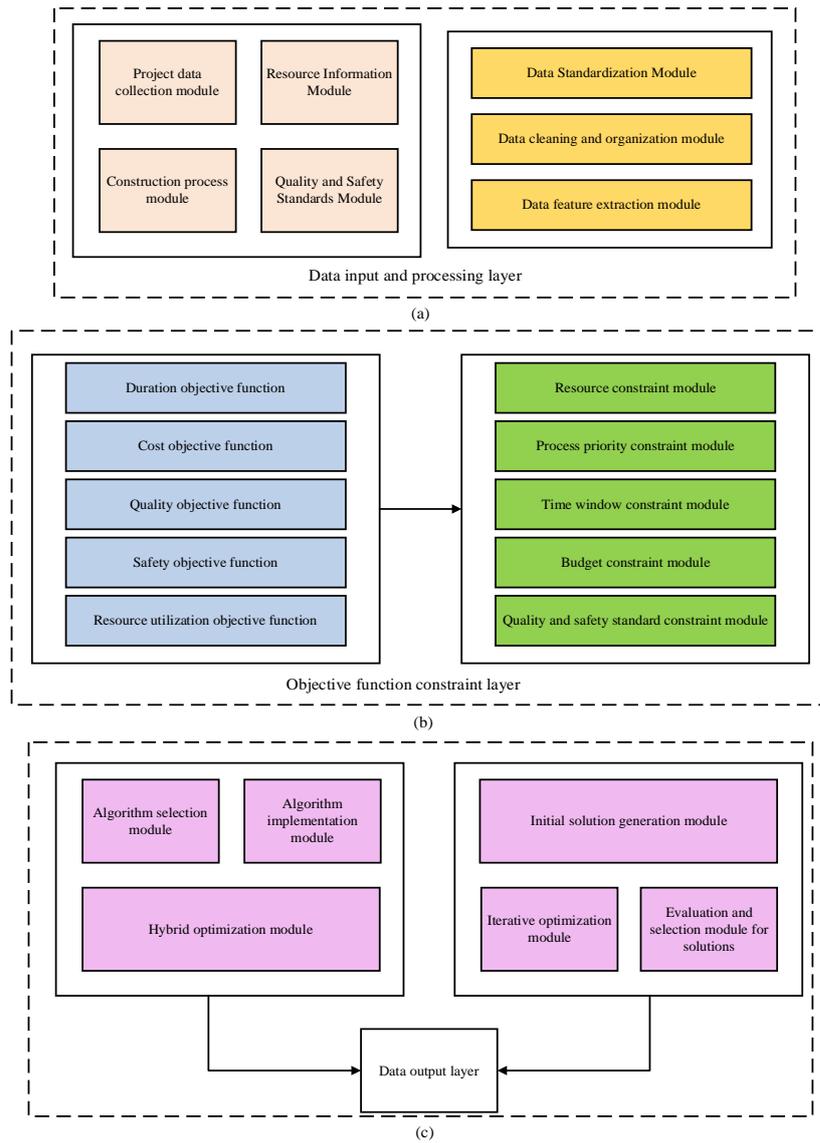


Figure 1: The MOO model ((a): The data input and processing layer; (b): The objective function constraint layer; (c): The data output layer)

First, the objective of the model is to minimize the construction duration, as shown in equation (1):

$$T_{total} = \max_{i \in \{1,2,\dots,N\}} (T_i + D_i) \quad (1)$$

T_{total} represents the minimized construction duration, T_i refers to the start time, and D_i is the duration. Subsequently, the minimized total construction cost is calculated as equation (2):

$$C_{total} = \sum_{i=1}^N (C_{labor,i} + C_{material,i} + C_{equipment,i}) \quad (2)$$

C_{total} represents the minimized total construction cost; $C_{labor,i}$, $C_{material,i}$, and $C_{equipment,i}$ denote the labor, material, and equipment costs for each operation, respectively. N means the total quantity of the project. To ensure high-quality construction, the model is designed to maximize construction quality:

$$Q_{total} = \sum_{i=1}^N w_i \cdot Q_i \quad (3)$$

Q_{total} represents the maximized construction quality, w_i is the weight coefficient, and Q_i refers to the quality score. Similarly, maximizing construction safety is also an important objective of the model:

$$S_{total} = \sum_{i=1}^N v_i \cdot S_i \quad (4)$$

S_{total} refers to the maximized construction safety, v_i denotes the weight coefficient, and S_i represents the safety score.

The proposed optimized model architecture consists of three layers: the data input and processing layer, the objective function constraint layer, and the data output layer. These layers work collaboratively to achieve MOO in construction scheduling. Firstly, the data input and processing layer includes modules for project data collection, resource information, construction processes, and quality and safety standards. This layer is responsible for gathering real-time data and resource information from the construction process and performing preprocessing through data standardization, data cleaning and organization, and feature extraction. These steps ensure the accuracy and consistency of the data, providing reliable inputs for the subsequent optimization process. Secondly, the objective function constraint layer focuses on MOO. It optimizes the construction schedule, total cost, quality, safety, and resource utilization through the

duration objective function, cost objective function, quality objective function, safety objective function, and resource utilization objective function. This layer also includes modules for resource constraints, process priority constraints, time window constraints, budget constraints, and quality and safety standard constraints, ensuring that the generated scheduling solutions meet the project's practical needs and construction standards. Finally, the data output layer is responsible for executing the optimization algorithm and outputting the results. This layer encompasses modules for algorithm selection, algorithm implementation, hybrid optimization, initial solution generation, iterative optimization, and solution evaluation and selection. By choosing the appropriate optimization algorithms and combining them with hybrid optimization methods, this layer continuously generates and refines the scheduling solutions. Ultimately, it selects the best scheduling plan based on indicators such as safety, resource efficiency, and scheduling time. Overall, the proposed model achieves efficient, safe, and resource-optimized construction scheduling through a multi-layered modular design, offering effective support for intelligent construction management projects.

This framework design ensures that the MOO model for intelligent construction scheduling has flexibility, scalability, and efficiency, enabling it to handle complex scheduling requirements and provide optimized decision support.

4 Performance comparison and scheduling results analysis of intelligent construction scheduling models

4.1 Experimental results of performance comparison for intelligent construction scheduling model

The dataset used for the experiment is the Construction Project Management dataset from Kaggle.

This dataset includes detailed information on various construction projects, such as project schedules, resource allocation, and cost estimates, and is suitable for analyzing resource optimization issues in construction scheduling. The dataset can be downloaded from Kaggle's official website (<https://www.kaggle.com/>). The experiments are conducted in a high-performance computing environment to ensure the MOO algorithm's effectiveness and efficiency in intelligent construction scheduling. The experiments utilize an Intel Xeon Gold 6248R processor, which features 24 physical cores and 48 threads with a clock speed of 3.0 GHz, enabling efficient parallel computation. To meet the memory demands of the algorithm's computations, the system is equipped with 128 GB of DDR4 RAM, ensuring ample memory resources for data processing and algorithm training. For storage, a 2 TB NVMe solid-state drive is used to provide fast data read and write speeds, further accelerating the overall execution efficiency of the experiments. The operating system version of the experiment is Ubuntu 20.04 LTS, and the programming language and version are Python 3.8. The main libraries and versions are as follows:

1. TensorFlow: 2.5.0
2. PyTorch: 1.9.0
3. NumPy: 1.21.0
4. Pandas: 1.3.0
5. Matplotlib: 3.4.2

The model selects Non-dominated Sorting Genetic Algorithm II (NSGA-II) as the benchmark algorithm. It has a population size of 200, 100 generations, a crossover probability of 0.9, a mutation probability of 0.1, and an objective function weight of 0.5. The comparison models include Multi-Objective Evolutionary Algorithm based on Decomposition with Differential Evolution (MOEA/D-DE), Strength Pareto Evolutionary Algorithm 2 with Shift-based Density Estimation (SPEA2+SDE), and Non-dominated Sorting Genetic Algorithm III (NSGA-III). Performance comparison metrics are Hypervolume (HV), Generational Distance (GD), Inverted Generational Distance (IGD), and Spacing. The experimental results are presented in Figure 2.

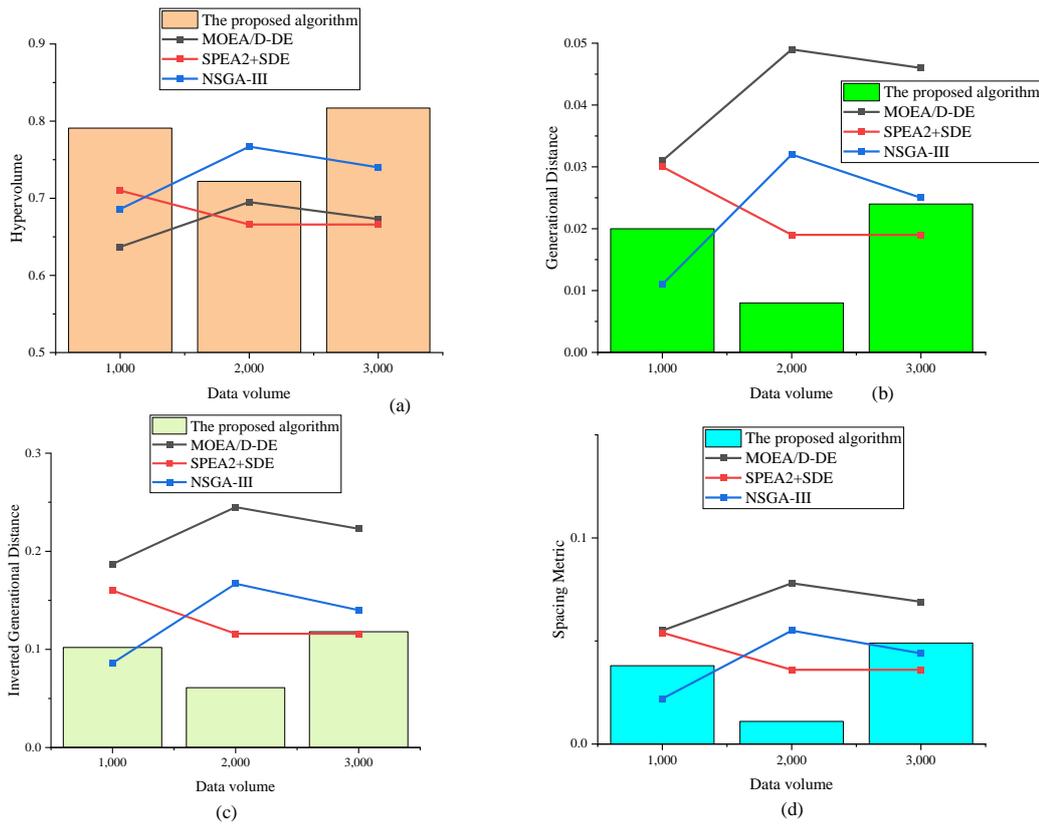


Figure 2: Performance comparison results ((a): HV; (b): GD; (c): IGD; (d): Spacing)

Figure 2 shows that for the HV metric, with a dataset size of 1000, the optimized model's HV value is 0.791, significantly exceeding the other three models. This indicates that the optimized model better covers the objective space in smaller datasets. MOEA/D-DE and NSGA-III have HV values of 0.637 and 0.686, respectively, showing relatively weaker performance, while SPEA2+SDE performs slightly better with an HV value of 0.710, surpassing MOEA/D-DE. With a dataset size of 2000, the HV value of the optimized model is 0.722, maintaining a leading position, though its performance slightly decreases compared to other dataset sizes. NSGA-III performs better with an HV value of 0.767, illustrating strong adaptability to medium-sized datasets. MOEA/D-DE's HV value rises to 0.695, while SPEA2+SDE decreases to 0.666. When the dataset size reaches 3000, the optimized model's HV value reaches 0.817, showcasing its advantage with large datasets. NSGA-III also performs relatively well with a value of 0.740, while MOEA/D-DE and SPEA2+SDE have HV values of 0.666, showing no significant advantage. Regarding GD, the optimized model performs exceptionally well across diverse dataset sizes, particularly with a dataset size of 2000, where it achieves the smallest GD value of 0.008, demonstrating strong convergence. This indicates that the solution set generated by the proposed model is very close to the true Pareto front across all dataset sizes, exhibiting excellent optimization capability. NSGA-III performs best when the dataset size is 1000, but its convergence decreases with increasing dataset size, particularly with a

GD value rising to 0.032 when the dataset size reaches 2000. However, NSGA-III's performance with a dataset size of 3000 is similar to the proposed model, indicating adaptability to large-scale datasets. MOEA/D-DE's GD values are consistently high across all dataset sizes, especially reaching 0.049 with a dataset size of 2000, demonstrating insufficient convergence performance for large-scale data. In contrast, SPEA2+SDE shows stable performance, but its GD values are generally higher than those of NSGA-III and the optimized model, reflecting less favorable convergence compared to both. For IGD, with a dataset size of 1000, NSGA-III achieves the best IGD value of 0.086, indicating its excellent coverage of the true Pareto front. The optimized model's IGD value is 0.102, slightly higher than NSGA-III, but still shows strong coverage capability. SPEA2+SDE has an IGD value of 0.160, showing moderate performance, while MOEA/D-DE has a relatively high IGD value of 0.187, indicating slightly weaker coverage capability. With a dataset size of 2000, the optimized model's IGD value drops to 0.061, demonstrating outstanding coverage of the solution set and significantly surpassing other models. SPEA2+SDE's IGD value is 0.116, showing relatively good performance, while NSGA-III's IGD value is 0.167, indicating a notable decrease in coverage performance. MOEA/D-DE's IGD value is 0.245, reflecting the weakest coverage capability at this dataset size. For a dataset size of 3000, the optimized model's IGD value is 0.118, slightly increasing the value for the 2000 dataset, but still maintaining strong coverage performance. NSGA-III's

IGD value increases to 0.140, showing relatively good solution set coverage. SPEA2+SDE's IGD value remains at 0.116, consistent with the 2000 dataset size, showing stable performance. MOEA/D-DE's IGD value is 0.223, still the weakest among all models. In terms of the Spacing metric, the optimized model performs the best. In particular, with a dataset size of 2000, the Spacing value is only 0.011, indicating an extremely uniform distribution of solutions in the objective space. This means the optimized model generates solutions close to the optimal and ensures a good distribution of these solutions in the objective space. NSGA-III performs excellently with a dataset size of 1000, with a Spacing value of 0.022. However, the uniformity of the solution set decreases with increasing dataset size, especially with a Spacing value rising to 0.055 for a dataset size of 2000, signaling instability in distribution uniformity. SPEA2+SDE shows stable performance across different dataset sizes, with Spacing values between 0.036 and 0.054. Although it does not reach the level of NSGA-III or the optimized model, its stability demonstrates adaptability across various dataset sizes. MOEA/D-DE performs relatively weakly in the Spacing metric, particularly with a dataset size of 2000, where the Spacing value is 0.078. It displays an uneven distribution of solutions, indicating that MOEA/D-DE lacks uniformity in the solution set when handling large-scale data.

4.2 Safety assessment and resource efficiency analysis

In optimizing intelligent construction scheduling, safety assessment and resource efficiency analysis are crucial for assessing the quality of scheduling plans. By evaluating the safety and resource utilization efficiency of the generated scheduling plans, it is possible to ensure that construction projects operate efficiently under safety regulations, maximize resource utilization, and reduce

potential risks. Safety assessment is a key step in construction scheduling, aimed at effectively preventing and controlling various potential safety risks during execution. To this end, this work employs a comprehensive evaluation system that integrates multiple safety indicators for a thorough assessment. Table 4 presents the details.

Table 4: The safety indicator system

Indicator	Description
Accident Rate	It evaluates the frequency of safety incidents occurring during construction, typically calculated per 1,000 work hours.
Safety Distance	It measures the minimum distance between construction site personnel and hazard sources.
Personnel Protective Measures	It checks whether all construction personnel are equipped with necessary protective gear (such as helmets and safety belts).
Equipment Operation Safety	It assesses the stability of construction equipment under different operating conditions to prevent equipment failure or accidents.

A safety risk analysis model based on the fuzzy comprehensive evaluation method is established to assess the scheduling plans quantitatively. By scoring each safety indicator (from 1 to 5) and combining expert weight analysis, a comprehensive safety score is calculated for each scheduling plan. The fuzzy comprehensive evaluation method-based safety risk analysis model quantitatively evaluates the scheduling plans. Figure 3 presents the experimental results.

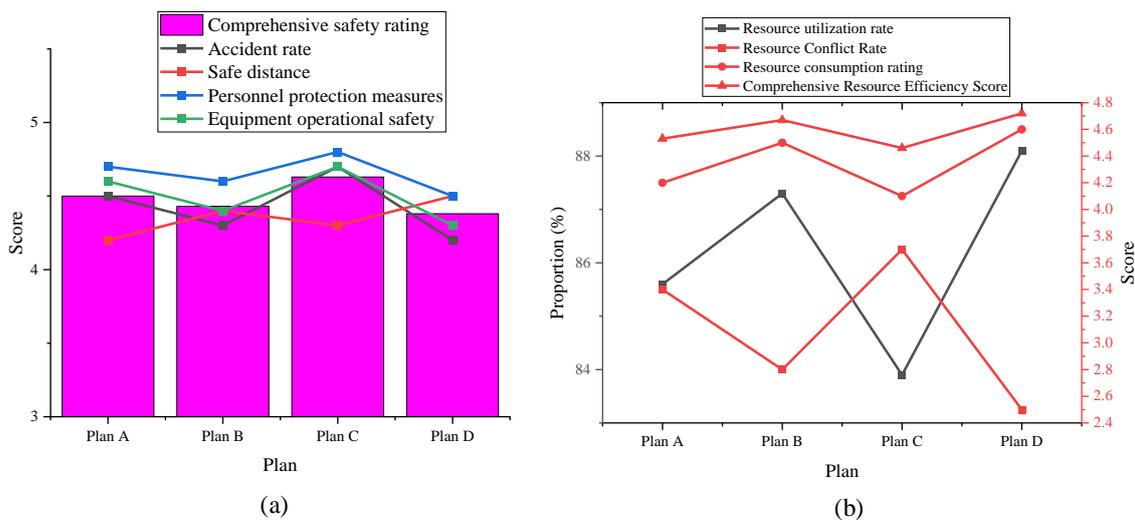


Figure 3: Safety assessment and resource efficiency analysis ((a): Safety assessment; (b): Resource efficiency)

According to the safety assessment scores, Plan C has the highest overall safety score of 4.63, indicating that it performs best in safety. Its safety indicators are all rated

highly, particularly in personnel protective measures and equipment operation safety. This suggests that Plan C can effectively reduce safety risks on the construction site and

ensure smooth project implementation. In contrast, Plan D has the lowest safety score of 4.38, illustrating deficiencies in some critical safety indicators. Moreover, Plan D has the highest overall resource efficiency score of 4.72, demonstrating high resource utilization, low conflict rate, and reasonable resource consumption. Plan D excels particularly in the key indicators of resource utilization and conflict rate, implying that it achieves optimal configuration and maximizes construction efficiency under limited resources. Plan B follows closely with a resource efficiency score of 4.67, while Plans A and C score slightly lower at 4.53 and 4.46, respectively.

The analysis indicates that while Plan C excels in safety, it is slightly inferior to Plan D in resource efficiency. This suggests that a balance between safety and resource efficiency may be necessary in practical construction scheduling, and the most suitable scheduling plan should be chosen based on the specific project needs. The final scheduling decision should integrate safety assessment with resource efficiency to ensure smooth project progress and effective resource utilization.

5 Discussion

This work primarily analyzes the performance differences of the optimized model under varying dataset sizes. First, NSGA-III performs quite well on large datasets but shows weaker adaptability on smaller datasets. This phenomenon may be related to NSGA-III's strategy for managing the diversity of the solution set, as the algorithm relies on reference points to ensure the uniform distribution of solutions. However, on smaller datasets, the reference point configuration may cause sparse solution distributions, thereby affecting the optimization results. As the dataset size increases, the increased number of reference points helps the algorithm better capture the multi-dimensional objective space, resulting in stronger adaptability on large datasets. Nevertheless, the suboptimal performance on smaller datasets highlights its sensitivity to dataset size, particularly in scenarios requiring finer granularity control. In contrast, the proposed optimization model demonstrates strong robustness across all dataset sizes, especially showing an advantage in the Spacing metric. The Spacing metric reflects the uniformity of the solution set's distribution in the objective space; The proposed model ensures a balanced distribution of solutions across different dimensions through the design of the objective function and iterative strategies. This uniformity enhances solution diversity and enables decision-makers to select solutions suited to various contexts from a multi-dimensional solution set, further showcasing the model's robustness.

Regarding safety assessment, although this work provides a detailed analysis of safety scores, these scores should be understood in terms of their practical impact on construction. Higher safety scores indicate that the solutions are more effective at reducing safety incidents during construction and meeting the high standards for personnel and equipment protection in construction projects. This satisfies industry safety standards and helps reduce the risk of project delays due to accidents, ensuring

timely project delivery. Furthermore, solutions with high resource efficiency are of significant practical value in construction projects. High resource efficiency means that the project has been more rationally planned in terms of resource allocation, reducing material and equipment waste. Hence, it directly contributes to lowering the total cost of the project and improving construction progress. This efficient use of resources is critical in projects with tight budgets or limited resources, as it ensures compliance with project deadlines while minimizing environmental impact, and aligning with sustainable development goals.

Overall, through its high scores in safety and resource efficiency, the proposed model's practical application value in real-world construction scenarios has been further validated. The model's robustness and multi-scenario applicability highlight its potential in the intelligent construction field, particularly in complex construction projects that require a balance between safety and resource optimization, demonstrating its strong practical utility.

6 Conclusion

This work systematically investigates intelligent construction scheduling using MOO algorithms, to generate optimal scheduling plans that balance safety and resource efficiency. A new optimization model is proposed, validated, and compared with classical MOO algorithms such as NSGA-III, MOEA/D-DE, and SPEA2+SDE. The proposed model demonstrates significant advantages in MOO problems. Comparative analysis of experimental data reveals that the model consistently outperforms the comparison models on key indicators such as HV, GD, IGD, and S, particularly when dealing with larger datasets. This indicates that the proposed model can better balance the conflicts between various objective functions and generate higher-quality solution sets while ensuring the scheduling plan meets MOO requirements. This work confirms that safety and resource efficiency are interdependent yet indispensable factors in practical construction scheduling through safety assessment and resource efficiency analysis of the generated scheduling plans. The assessment methods and the proposed models effectively balance these factors, providing a scientific basis for construction management. This comprehensive approach not only enhances the feasibility and rationality of the scheduling plans but also offers valuable insights for future intelligent construction management practices.

The optimized model proposed in this work performs excellently under specific data volumes and static construction environments but still exhibits certain limitations in dynamic adaptability. Specifically, the model's applicability may be limited in dynamic scenarios, such as fluctuations in resource availability or construction delays. This work does not test the model in dynamic construction environments. However, future research could validate its real-time adaptability through case studies or hypothetical scenarios, such as how it handles scheduling adjustments due to resource shortages or unexpected events. The model's ability to respond in

real-time to changes in the construction environment remains an area for further research. Additionally, the model faces high computational complexity when handling large-scale datasets or high-dimensional multi-objective problems, potentially resulting in slower convergence and impacting real-time performance in practical applications. This limitation suggests that in high-computation-demand scenarios, the model may encounter computational constraints. It indicates that future research should focus on simplifying the computational steps without compromising the quality of the solution set to enhance the model's practical usability. To address these limitations, future research could focus on enhancing the model's dynamic adaptability and computational efficiency. For example, dynamic scheduling algorithms or adaptive optimization mechanisms could be introduced to improve the model's responsiveness to resource changes and unexpected events. Furthermore, integrating other optimization techniques, such as metaheuristic algorithms, distributed computing, and parallel computing, could help reduce computational time and accelerate convergence. Additionally, future work could explore applying the model to more complex, large-scale construction scenarios to test its applicability and robustness in different environments, thus enhancing its practical value in the intelligent construction field.

References

- [1] Yan X, Zhang H, Zhang W. Intelligent monitoring and evaluation for the prefabricated construction schedule. *Computer-Aided Civil and Infrastructure Engineering*, 2023, 38(3): 391-407. <https://doi.org/10.1111/mice.12838>.
- [2] Zhang Z, Zhang S, Zhao Z, et al. HydroBIM—Digital design, intelligent construction, and smart operation. *J. Intell. Constr.*, 2023, 1(2): 9180014. <https://doi.org/10.26599/JIC.2023.9180014>.
- [3] Doukari O, Seck B, Greenwood D. The creation of construction schedules in 4D BIM: A comparison of conventional and automated approaches. *Buildings*, 2022, 12(8): 1145. <https://doi.org/10.3390/buildings12081145>.
- [4] Peiris A, Hui F K P, Duffield C, et al. Production scheduling in modular construction: Metaheuristics and future directions. *Automation in Construction*, 2023, 150(11): 104851. <https://doi.org/10.1016/j.autcon.2023.104851>.
- [5] Singh A K, Pal A, Kumar P, et al. Prospects of integrating BIM and NLP for automatic construction schedule management. *ISARC. Proceedings of the International Symposium on Automation and Robotics in Construction. IAARC Publications*, 2023, 40(6): 238-245. <https://doi.org/10.22260/ISARC2023/0034>.
- [6] Wang J. Design of Intelligent Construction System for Assembly Building Based on Improved IoT. *Informatica*, 2024, 48(10). <https://doi.org/10.31449/inf.v48i10.5889>.
- [7] Chen J. Construction and Application of an Economic Intelligent Decision-making Platform Based on Artificial Intelligence Technology. *Informatica*, 2024, 48(9): 78. <https://doi.org/10.31449/inf.v48i9.5705>.
- [8] Li L. Dynamic Cost Estimation of Reconstruction Project Based on Particle Swarm Optimization Algorithm. *Informatica*, 2023, 47(2). <https://doi.org/10.31449/inf.v47i2.4026>.
- [9] Feng J, Xu Y, Zhang A. Intelligent engineering management of prefabricated building based on BIM Technology. *Informatica*, 2022, 46(3). <https://doi.org/10.31449/inf.v46i3.4047>.
- [10] Wang X. Edge Computing Based Multi-Objective Task Scheduling Strategy for UAV with Limited Airborne Resources. *Informatica*, 2024, 48(2). <https://doi.org/10.31449/inf.v48i2.5885>.
- [11] Vijaya C, Srinivasan P. Multi-objective Metaheuristic Technique for Energy Efficient Virtual Machine Placement in Cloud Data Centers. *Informatica*, 2024, 48(6). <https://doi.org/10.31449/inf.v48i6.5263>.
- [12] Bendiaf L M, Harbouche A, Tahraoui A M, et al. An Innovative Task Scheduling Method Using the Knapsack Algorithm in Heterogeneous Computing Systems. *Informatica*, 2024, 48(16). <https://doi.org/10.31449/inf.v48i16.5765>.
- [13] Aslan S, Türkakin O H. A construction project scheduling methodology considering COVID-19 pandemic measures. *Journal of safety research*, 2022, 80(6): 54-66. <https://doi.org/10.1016/j.jsr.2021.11.007>.
- [14] Chenya L, Aminudin E, Mohd S, et al. Intelligent risk management in construction projects: Systematic Literature Review. *IEEE Access*, 2022, 10(3): 72936-72954. <https://doi.org/10.1109/ACCESS.2022.3189157>.
- [15] Baduge S K, Thilakarathna S, Perera J S, et al. Artificial intelligence and smart vision for building and construction 4.0: Machine and deep learning methods and applications. *Automation in Construction*, 2022, 14(1): 104440. <https://doi.org/10.1016/j.autcon.2022.104440>.
- [16] Herath T C, Herath H S B, Cullum D. An information security performance measurement tool for senior managers: Balanced scorecard integration for security governance and control frameworks. *Information Systems Frontiers*, 2023, 25(2): 681-721. <https://doi.org/10.1007/s10796-022-10246-9>.
- [17] Samha A K. Strategies for efficient resource management in federated cloud environments supporting Infrastructure as a Service (IaaS). *Journal of Engineering Research*, 2024, 12(2): 101-114. <https://doi.org/10.1016/j.jer.2023.10.031>.
- [18] Sang Y, Tan J. Intelligent factory many-objective distributed flexible job shop collaborative scheduling method. *Computers & Industrial Engineering*, 2022, 164(32): 107884. <https://doi.org/10.1016/j.cie.2021.107884>.
- [19] Hosseinian A H, Baradaran V. A two-phase approach for solving the multi-skill resource-constrained multi-project scheduling problem: a case study in construction industry. *Engineering, Construction and*

- Architectural Management, 2023, 30(1): 321-363. <https://doi.org/10.1108/ECAM-07-2019-0384>.
- [20] Li S, Zhang M, Wang N, et al. Intelligent scheduling method for multi-machine cooperative operation based on NSGA-III and improved ant colony algorithm. Computers and Electronics in Agriculture, 2023, 204(44): 107532. <https://doi.org/10.1016/j.compag.2022.107532>.
- [21] Prieto S A, Mengiste E T, García de Soto B. Investigating the use of ChatGPT for the scheduling of construction projects. Buildings, 2023, 13(4): 857. <https://doi.org/10.3390/buildings13040857>.
- [22] Obiuto N C, Adebayo R A, Olajiga O K, et al. Integrating artificial intelligence in construction management: Improving project efficiency and cost-effectiveness. Int. J. Adv. Multidisc. Res. Stud, 2024, 4(2): 639-647. <https://doi.org/archives/archive-1711453341>.
- [23] Liu Y, You K, Jiang Y, et al. Multi-objective optimal scheduling of automated construction equipment using non-dominated sorting genetic algorithm (NSGA-III). Automation in Construction, 2022, 143(56): 104587. <https://doi.org/10.1016/j.autcon.2022.104587>.
- [24] Yin Z, Xu F, Li Y, et al. A multi-objective task scheduling strategy for intelligent production line based on cloud-fog computing. Sensors, 2022, 22(4): 1555. <https://doi.org/10.3390/s22041555>.
- [25] Xie L L, Chen Y, Wu S, et al. Knowledge extraction for solving resource-constrained project scheduling problem through decision tree. Engineering, Construction and Architectural Management, 2024, 31(7): 2852-2877. <https://doi.org/10.1108/ecam-04-2022-0345>.
- [26] Yu M. Construction of regional intelligent transportation system in smart city road network via 5G network. IEEE Transactions on Intelligent Transportation Systems, 2022, 24(2): 2208-2216. <https://doi.org/10.1109/TITS.2022.3141731>.
- [27] Chen L, Lu Q, Han D. A Bayesian-driven Monte Carlo approach for managing construction schedule risks of infrastructures under uncertainty. Expert Systems with Applications, 2023, 212(3): 118810. <https://doi.org/10.1016/j.eswa.2022.118810>.
- [28] Yang Y, Li X. A knowledge-driven constructive heuristic algorithm for the distributed assembly blocking flow shop scheduling problem. Expert Systems with Applications, 2022, 202(42): 117269. <https://doi.org/10.1016/j.eswa.2022.117269>.

A Big Data-Driven Approach to Financial Analysis and Decision Support System Design

Sa Zhang

¹School of Accountancy, Zhengzhou Vocational College of Finance and Taxation, Zhengzhou 450000, Henan, China
Email: julia6076@163.com

Keywords: big data, financial data analysis, decision support system

Received: September 3, 2024

This study aims to design and implement a financial data analysis and decision support system leveraging big data to enhance financial management and decision-making capabilities in complex market environments. Through a detailed examination of Company A's financial data, key indicators such as sales revenue, cost of sales, net profit, current assets, current liabilities, total assets, and shareholders' equity are selected. Utilizing big data technology, the system achieves efficient data processing, precise financial risk early warning, and scientifically informed investment decision support. Using Hadoop and Spark, the system efficiently processes extensive financial data while monitoring key indicators, such as sales revenue, cost of sales, and profitability. Specific financial models, including net present value (NPV) and internal rate of return (IRR), were tested for their effectiveness in supporting optimal investment decisions, yielding an NPV of 5.12 million and an IRR of 16.8% in the best-performing scenario. Performance metrics indicate a consistent improvement in data processing speed, accuracy reaching 98.9%, and user satisfaction rising from 8.2 to 8.7 over three years. The results indicate that the system performs effectively in terms of data processing speed, accuracy, and user satisfaction, enhancing both the efficiency of financial management and the precision of decision-making processes within enterprises. The financial risk early warning system successfully identifies potential risks in a timely manner, while the investment decision support system aids enterprises in selecting the optimal investment strategy by utilizing indicators such as net present value, internal rate of return, and investment payback period. This study highlights the value of applying big data technology in financial management, offering robust support for enterprises aiming for stable development within a rapidly evolving market environment. Future research will focus on further optimizing system functionality and expanding application scenarios to adapt to the shifting financial management demands of enterprises.

Povzetek: Predstavljen je sistem za finančno analizo in podporo odločanju, ki uporablja tehnologijo velikih podatkov. Integracija Hadoop in Spark omogoča hitro in natančno obdelavo podatkov, napredni modeli (NPV, IRR) pa izboljšujejo napovedi in oceno finančnih tveganj. Sistem omogoča optimizirano finančno upravljanje ter pravočasno odkrivanje tveganj, kar izboljšuje strateško odločanje podjetij v dinamičnem tržnem okolju.

1 Introduction

With the rapid advancement of information technology, global data volumes have grown exponentially, and big data technology has significantly transformed the operations across numerous industries. Financial management, as a central aspect of enterprise management, has also been deeply influenced by big data innovations. Big data technology enables the processing and analysis of vast financial datasets and provides real-time, precise decision support, helping enterprises sustain a competitive advantage in an increasingly demanding market environment.

The rapid development of information technology has significantly impacted the global financial industry, with big data emerging as a critical tool in financial management. Financial data is characterized by large volumes, complexity, and real-time requirements, making traditional data analysis methods insufficient. Big data

technology offers solutions to these challenges by enabling efficient processing and analysis of massive financial datasets. By leveraging big data, enterprises can extract valuable insights, enhance their financial risk management, and improve decision-making capabilities.

The objective is to develop a comprehensive financial data analysis and decision support system that uses key financial indicators such as sales revenue, cost of sales, net profit, current assets, and liabilities. This system provides early risk warnings and supports investment decisions, helping enterprises remain competitive in rapidly changing markets. The integration of big data in financial management enhances data processing accuracy, enables real-time decision-making, and supports sustained business growth.

To enhance the robustness and accuracy of financial data analysis models, incorporating optimization techniques such as hyperparameter tuning, cross-validation, and feature engineering is essential.

Hyperparameter tuning allows for the adjustment of model parameters to maximize performance, enhancing the predictive power of models in varying financial scenarios. Cross-validation, particularly k-fold cross-validation, ensures that the models generalize well to new data by testing across multiple data subsets, reducing the risk of overfitting. Feature engineering, including techniques like scaling, encoding categorical data, and creating interaction terms, refines the data inputs to optimize model learning. These approaches collectively provide adaptability and precision, crucial for the dynamic financial environment where timely, accurate decision-making is a priority.

Traditional financial data analysis methods face significant challenges due to the large volume of data, the complexity of data types, and high demands for real-time processing. Big data technology effectively addresses these issues by enabling in-depth data mining and analysis, uncovering hidden patterns and trends, and enhancing the accuracy and efficiency of financial analysis. The use of Decision Support Systems (DSS) in enterprise management is increasingly prevalent. Combining DSS with big data technology provides businesses with more intelligent and customized decision support.

In recent years, companies have increasingly recognized the importance of data and are developing and refining their own data analysis and decision support systems. However, in practice, challenges persist, including suboptimal data quality, incomplete analysis models, and insufficient system integration. These issues limit the effective application of big data and decision support systems in financial management to some extent.

The objective of this study is to design a big data-based financial data analysis and decision support system, systematically addressing the limitations of existing systems and enhancing both the accuracy of financial data analysis and the effectiveness of decision support. This research will approach the topic through theoretical analysis, system design, empirical testing, and other methodologies, aiming to provide practical solutions for financial management and to advance the adoption of big data technology within the financial sector.

This study aims to enhance financial data analysis and decision-making capabilities for businesses by using a big data-based system. Specific objectives are as follows:

Improve data processing capabilities: Integrate big data technology to strengthen financial data processing, addressing traditional challenges related to large data volumes, diverse data types, and stringent real-time requirements. **Optimize financial analysis models:** Develop and refine financial data analysis models using big data technology to extract deeper insights from financial data, thereby improving the accuracy and comprehensiveness of financial analysis and providing valuable insights for business decision-making.

Design a robust decision support system integrated with big data analysis results to deliver real-time, accurate decision support for enterprise management. This system enables businesses to make informed and strategic decisions even in complex and volatile market conditions.

Improve system integration by efficiently merging financial data analysis with decision support functions, ensuring smooth coordination across system modules. This enhances overall system efficiency and stability, providing businesses with a unified, comprehensive financial management solution.

By systematically analyzing the application of big data technology in financial data processing and decision support, this research enriches the academic content at the intersection of financial management and information technology. This study investigates specific application methodologies and pathways of big data in financial management, offering a novel perspective for theoretical exploration. Developing a big data-driven financial analysis model demonstrates the potential and strengths of data processing technologies in financial analysis. Based on existing DSS theory and enhanced by big data analysis, this research proposes a new system design and implementation framework, deepening DSS theoretical research and expanding its applicability within financial management. The exploration of innovative financial analysis methods further broadens the theoretical foundation of financial analysis, as the integration of big data processing technologies advances and refines traditional financial analysis models, offering a valuable new reference for related academic research.

At a practical level, this study designs and implements a big data-based financial data analysis and decision support system, offering positive impacts on enterprise financial management practices. The system leverages big data technology to process and analyze substantial volumes of financial data swiftly and accurately, enhancing both the efficiency and precision of corporate financial management. This enables enterprises to make timely and informed decisions in complex market environments. By combining decision support capabilities with big data analytics, the system provides a scientific and accurate decision-making foundation for enterprise management. Through real-time data analysis and forecasting, the system effectively supports financial decision-making processes, reduces decision-making risks, and improves decision quality. With this study, enterprises can better understand and apply big data technology to advance their information infrastructure. This system offers a comprehensive solution for financial data processing and decision support, aiding enterprises in achieving significant progress in digital transformation and information management.

In recent years, the application of big data technology in financial data analysis and decision support systems has garnered considerable attention. Current research emphasizes enhancing data processing capabilities and decision support through big data technology, adapting to increasingly complex financial management environments and evolving market needs. Casturi and Sunderraman proposed a rule-based, cost-efficient big data analytics aggregation engine for portfolio management, demonstrating big data technology's substantial advantages in increasing data processing efficiency and reducing costs. This study provides a theoretical foundation and practical experience for

enterprises utilizing big data technology in financial data analysis [1].

In the educational field, Fahd and Miah designed and evaluated big data analytics methods for predicting student success factors [2]. Although their research primarily focuses on students, the methods and technology they employed are also applicable to financial data analysis. By identifying key factors within large datasets, these techniques enable data-driven decision-making processes. Dutta's research examined the asset-liability management models of life insurance companies, focusing on decision support system outcomes [3]. This research underscores the effectiveness of big data technology in managing complex financial environments and serves as a valuable reference for designing decision support systems.

For credit risk assessment, Lu et al. proposed a decision support method based on dynamic Bayesian networks, showcasing big data technology's potential in risk management by enhancing accuracy and timeliness through dynamic modeling [4]. Talamo et al. addressed organizational and individual decision-making challenges in the design of financial artificial intelligence systems, providing methodological insights that are especially beneficial for managing the complexity of decision processes in this study [5].

Peng and Bao developed an enterprise management analysis framework utilizing big data technology, which demonstrates the extensive applications of big data in enterprise management, particularly in enhancing data analysis and decision support capabilities [6]. Popovic et al. explored the influence of big data analytics on enterprise high-value business performance, concluding that big data analysis can significantly improve business performance, reinforcing its value in financial data

analysis [7]. Liu et al. introduced a model for Internet financial risk control based on machine learning algorithms, showcasing the role of machine learning in financial risk management and offering technical guidance and methodological references for this study [8].

The creative contributions of this study are reflected in the application of big data technology to enhance financial data analysis and decision support systems. The use of advanced big data processing technologies, such as Hadoop and Spark, enables the system to handle large volumes of complex financial data with high efficiency and accuracy. Additionally, the design of an integrated financial risk early warning system allows enterprises to identify and mitigate risks through real-time monitoring of key financial indicators, such as the current ratio and asset-liability ratio. Another key innovation is the application of decision-making tools, like net present value (NPV) and internal rate of return (IRR), to assist businesses in making optimal investment choices. The system's design not only improves financial analysis accuracy but also empowers enterprises with data-driven insights to enhance decision-making in competitive markets.

Existing research demonstrates significant advancements in the use of big data technology for financial data analysis and decision support systems [9]. Nonetheless, several persistent challenges remain, such as ensuring high data quality, achieving seamless system integration, and meeting real-time processing demands. Addressing these issues is essential for advancing the implementation and utility of big data technology in financial management, as overcoming these obstacles will facilitate more robust and effective applications within the field.

Table 1: Summary of reviewed research in related work section.

Study Reference	Methods Used	Data Set	Evaluation Metrics	Key Findings
Casturi & Sunderraman [1]	Rule-based analytics, cost aggregation	Portfolio management data	Cost-effectiveness, efficiency	Big data analytics enhance cost-efficiency in portfolio management.
Fahd & Miah [2]	Predictive analytics for success factors	Student academic records	Predictive accuracy	Data analytics effectively predict key success factors for students.
Dutta et al. [3]	Asset-liability management model	Life insurance data	Decision support efficiency	Model improves asset-liability alignment in complex financial environments.
Lu et al. [4]	Dynamic Bayesian networks	Credit risk data	Accuracy, timeliness	Model enhances credit risk assessment accuracy with real-time data.

Talamo et al. [5]	Decision modeling, AI	Financial decision processes	Decision complexity, efficiency	AI-based systems improve decision-making under complex financial conditions.
Peng & Bao [6]	Business management analysis framework	Corporate management data	Performance enhancement	Big data application boosts business analysis and strategic planning.
Liu et al. [8]	Machine learning risk control model	Internet finance risk data	Risk prediction accuracy	Machine learning strengthens risk control in internet finance sectors.

2 Theoretical basis

2.1 Overview of big data technology

The methods outlined focus on systematically addressing financial data analysis and decision support system design. First, data is gathered from multiple sources, followed by preprocessing, which includes cleaning, transformation, and integration to ensure data quality. The system incorporates both relational and NoSQL databases to handle structured and unstructured data efficiently. Big data processing tools, like Hadoop and Spark, are then used to manage large volumes of data swiftly and accurately. The methods also integrate financial analysis models, such as regression analysis and decision trees, to predict and optimize financial decisions. These models are implemented within a user-friendly interface to provide real-time decision-making support. The methodology is coherent and logical, aligning well with the goal of improving financial data processing, risk management, and decision-making efficiency.

2.1.1 Definition and characteristics of big data

Big data refers to datasets too large or complex to be managed by traditional data processing tools. Its primary characteristics include Volume, Variety, Velocity, and Veracity.

Large volume: This is seen in the exponential growth of data amounts, far exceeding traditional database processing capacities. Zheng et al. highlight that modern information systems produce massive quantities of both structured and unstructured data daily, and that conventional database technologies are insufficient for handling this scale [10]. Advanced storage and processing solutions, such as distributed storage and parallel computing, are thus essential.

Variety: This characteristic is represented by the wide range of data sources and types. Xue's research notes that with rapid advances in information technology, data sources have diversified to include sensor data, social media data, and transaction records [11]. These sources contain both structured data and large amounts of

unstructured data, such as text, images, and videos, complicating data processing and analysis tasks.

Fast: This characteristic refers to the rapid rate at which data is generated and processed. Modern information systems demand real-time or near real-time data processing and analysis to enable timely decision-making. Xu and Zhou noted that advances in big data technology make real-time data processing feasible, thereby supporting enterprises in making prompt decisions within a rapidly evolving market environment [12]. These real-time requirements place higher demands on data processing technologies, necessitating the adoption of streaming data processing and real-time analytics.

Authenticity: This aspect emphasizes the quality and reliability of data sources. Big data comes from diverse sources, leading to varying degrees of data quality, which makes ensuring data authenticity and accuracy a critical concern.

2.1.2 Big data processing technology

Big data processing technology encompasses the full sequence from data acquisition, storage, processing, and analysis (as shown in Figure 1). Applying big data technology enables the efficient handling and analysis of vast datasets, offering robust technical support for decision support systems.

This study reinforces the critical evaluation of the references, focusing on those sources that directly support the arguments and methods. By selectively referencing the literature, the research ensures consistency with the center's claims and enhances credibility. This study draws on articles published on Informatica, such as theoretical discussions on big data processing and decision support system design. Provides reliable, peer-reviewed insights for research that confirm the application of big data in financial analysis and risk management.

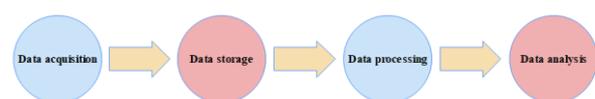


Figure 1: Big data technology processing flow.

The initial step involves data acquisition, gathering information from a wide range of sources. Rodrigues et al. indicate that a significant challenge in data collection arises from the diverse origins and formats of these datasets [13]. Sources such as sensors, social media, and transaction records provide critical data, often in both structured and unstructured formats, necessitating various acquisition and preprocessing techniques.

The second step is data storage, a critical component of big data processing that influences overall efficiency and reliability. Traditional relational databases often cannot handle the scale and complexity of big data, making NoSQL databases and distributed file systems (such as Hadoop HDFS) the primary options for effective storage. Jensen et al. noted that NoSQL databases offer flexible data management across diverse formats, while distributed file systems ensure efficient data storage and access mechanisms [14].

The third step is data processing, where distributed computing technology is essential. MapReduce serves as the primary computing model within the Hadoop framework, dividing data into smaller units for parallel processing across multiple nodes, thereby achieving efficient large-scale data handling. Rodrigues highlighted that the MapReduce model is particularly advantageous in massive data processing scenarios that require intensive calculation and analysis [13]. As a big data processing engine, Spark enhances both speed and processing efficiency by utilizing in-memory computing technology, providing faster performance compared to MapReduce.

The fourth step is data analysis, which represents the ultimate objective of big data processing. Through data mining, machine learning, and statistical analysis, valuable insights and knowledge can be extracted from vast datasets. Jensen emphasized that machine learning algorithms in big data analysis significantly improve predictive accuracy and pattern recognition [14]. Techniques such as classification, clustering, and regression are widely applied in financial data analysis, market prediction, and risk assessment.

Data security and privacy are critical in big data processing. It is essential for big data technologies to implement effective strategies to safeguard data confidentiality and integrity while maintaining processing efficiency. Common security measures include encryption techniques, access control, and data anonymization, all of which play a vital role in protecting sensitive information. Proper encryption ensures data is unreadable to unauthorized users, while access control restricts data access based on user roles and permissions. Additionally, data anonymization techniques help preserve privacy by masking personal identifiers within datasets, thus preventing unauthorized disclosure of individual information.

2.2 Theoretical framework of financial data analysis

2.2.1 Financial data analysis methods

Financial data analysis primarily involves data preprocessing, data mining, and data visualization.

Data preprocessing is the initial stage of financial data analysis and includes data cleaning, transformation, and reduction. Li and Chen noted that preprocessing helps remove data noise, fills in missing values, and standardizes data formats to support subsequent analysis and mining processes [15]. When dealing with financial statements, it is essential to ensure all data is thoroughly cleaned and standardized for reliable analytical outcomes.

Data mining is the core phase of financial data analysis. Wang developed an investment recommendation model based on time series data, utilizing data mining techniques from financial time series analysis [16]. Common data mining approaches include classification, regression, clustering, and association rule mining. These methods enable analysts to detect patterns and trends within financial data, supporting informed investment decisions and risk management. Classification aids in credit risk assessment, clustering assists in customer segmentation, and association rule mining reveals correlations among financial indicators. Specifically, fuzzy clustering and the Apriori algorithm are frequently applied, with Guo highlighting their use in the interactive analysis of financial management software to identify frequent item sets and association rules within datasets [17].

Data visualization represents the final stage in financial data analysis, presenting complex results in a form accessible to decision-makers. Data visualization techniques include charts, dashboards, and interactive reports [18]. These tools enhance communication of analytical results and assist users in identifying hidden patterns and anomalies within data. Trends in financial indicators are easily observable through time series graphs, while correlation levels between various financial indicators can be visualized using heat maps. Each step of financial data analysis offers distinct benefits and challenges, necessitating tailored choices of techniques and methods aligned with specific analysis objectives and data characteristics.

2.2.2 Financial data analysis index system

The index system of financial data analysis plays a vital role in comprehensively assessing an enterprise's financial condition and operational performance by measuring and evaluating various key indicators. A well-constructed index system for financial data analysis includes core financial indicators such as operational efficiency, profitability, debt repayment capacity, and growth potential [19]. These indicators evaluate an enterprise's financial health from multiple perspectives, offering a scientific foundation for informed decision-making. The specific details are presented in Table 2 below.

Table 2: Financial data analysis indicator system.

Indicator Category	Indicator Name	Calculation Formula	Description
Profitability	Gross Profit Margin	(Sales Revenue - Cost of Goods Sold) / Sales Revenue	Measures basic profitability of sales activities
Net Profit Margin	Net Profit / Sales Revenue	Measures overall profitability of the enterprise	
Return on Assets (ROA)	Net Profit / Total Assets	Measures efficiency in using assets	
Return on Equity (ROE)	Net Profit / Shareholder's Equity	Measures return on shareholder's investment	
Solvency	Current Ratio	Current Assets / Current Liabilities	Evaluates short-term debt-paying ability
Quick Ratio	(Current Assets - Inventory) / Current Liabilities	A stricter evaluation of short-term solvency	
Debt to Asset Ratio	Total Liabilities / Total Assets	Evaluates financial structure stability	
Operating Efficiency	Inventory Turnover Ratio	Cost of Goods Sold / Average Inventory	Measures inventory management efficiency
Accounts Receivable Turnover	Sales Revenue / Average Accounts Receivable	Evaluates efficiency of accounts receivable recovery	
Total Asset Turnover	Sales Revenue / Average Total Assets	Measures asset utilization efficiency	
Growth Ability	Sales Growth Rate	(Current Period Sales Revenue - Previous Period Sales Revenue) / Previous Period Sales Revenue	Measures the growth rate of sales revenue
Net Profit Growth Rate	(Current Period Net Profit - Previous Period Net Profit) / Previous Period Net Profit	Measures the growth rate of net profit	
Asset Growth Rate	(Current Period Total Assets - Previous Period Total Assets) / Previous Period Total Assets	Measures the growth rate of asset scale	

2.2.3 Common models in financial data analysis

Commonly used models in financial data analysis include the regression analysis model, time series analysis model, and decision tree model. Among these, the regression analysis model is foundational and widely applicable, enabling companies to forecast financial indicators by examining the relationships between independent and dependent variables. The linear regression model, in particular, is a widely adopted variant of the regression analysis model, with its basic form represented as follows (Formula 1):

$$Y = \beta_0 + \beta_1 X + \hat{\epsilon} \quad (1)$$

Where, Y is the dependent variable, X is the independent variable, β_0 and β_1 are the intercept and

slope of the model, respectively, and $\hat{\epsilon}$ represents the error term. By estimating parameters β_0 and β_1 using the least squares method, the linear relationship between the independent and dependent variables can be established, enabling prediction and analysis based on the model.

The time series analysis model is utilized to process sequential data, examine historical data trends, and forecast future financial metrics. Common models within time series analysis include the autoregressive integrated moving average model (ARIMA) and the exponential smoothing model (ETS). The ARIMA model uniquely integrates autoregressive and moving average components to effectively capture temporal dependencies and trend patterns in data, and its basic form is as follows (Formula 2):

$$Y_t = c + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \theta_1 \hat{\epsilon}_{t-1} + \theta_2 \hat{\epsilon}_{t-2} + \dots + \theta_q \hat{\epsilon}_{t-q} + \hat{\epsilon}_t \quad (2)$$

Y represents the observed value at time t , c is the constant term, ϕ is the autoregressive coefficient, θ_j is the moving average coefficient, and $\hat{\epsilon}_t$ denotes the error term. Utilizing the ARIMA model enables the capture of trends and cyclical patterns within time series data, facilitating effective financial forecasting.

The decision tree model is a decision analysis tool structured as a tree, widely used for both classification and regression tasks. It works by recursively partitioning the dataset into smaller subsets, selecting the optimal split point at each node until a stopping criterion is met. The construction process follows these basic steps:

1. Select the optimal feature as the basis for node segmentation.
2. Divide the data set into subsets according to different values of feature.
3. Recursively build subtrees for each subset until the stop condition is met.

The decision tree model describes the node segmentation process through (Formula 3):

$$Gini(D) = 1 - \sum_{i=1}^k p_i^2 \quad (3)$$

$Gini(D)$ represents the Gini coefficient of node D , while k denotes the total number of classes, and p_i indicates the proportion of Class E samples. The Gini coefficient serves as a criterion for identifying the optimal segmentation point, ensuring the highest possible purity of the resulting subset.

2.3 Basic theory of decision support system

2.3.1 Definition and composition of decision support system

A Decision Support System (DSS) is a computer-based information system developed to assist decision-makers in handling complex decision-making tasks. By integrating data, analytical models, and user interfaces, DSS enables users to analyze issues, formulate decision strategies, and evaluate outcomes effectively [20]. Widely applied across fields such as enterprise management, healthcare, and financial investment, DSS provides a scientific foundation for decision-making that enhances both accuracy and efficiency. Through data analysis, model computation, and expert systems, DSS equips decision-makers with structured tools and methodologies for each phase of the decision-making process, including problem identification, data collection, model development, solution evaluation, and feedback.

The composition of a DSS typically includes several key subsystems:

The Data Management subsystem is responsible for the collection, storage, and management of data essential for decision-making. Data sources include both internal data (e.g., enterprise financial and production data) and external data (e.g., market and economic information) [21]. This subsystem houses the Database Management System (DBMS), which facilitates efficient data access

and management. A primary function of the Data Management subsystem is to ensure data accuracy, consistency, and timeliness, providing a reliable foundation for subsequent analysis and decision support.

The Model Management subsystem oversees the storage and administration of analytical models required for decision support. These models may include mathematical, statistical, optimization, or simulation types [22]. This subsystem allows for model definition, creation, modification, and execution, supporting flexible application and updating of models. By selecting appropriate models, complex scenarios can be simulated and analyzed, enabling decision-makers to evaluate the potential impacts of various options.

The User Interface subsystem serves as the intermediary between the DSS and the user. It offers an accessible interface that allows decision-makers to input data, select models, execute analyses, and view results. Key components of this subsystem include a graphical user interface, reporting tools, and visualization options, which present data and analysis outcomes clearly to facilitate users' understanding and application of system information.

The Knowledge Management subsystem plays a crucial role by gathering, storing, and managing decision-related knowledge, such as expert insights, business rules, and historical decision examples. This subsystem employs a Knowledge Base to convert expert knowledge into system rules, thereby assisting decision-makers in making more informed and scientifically grounded choices.

The Communication subsystem enables information exchange between the DSS and the external environment. Through this subsystem, the DSS can acquire up-to-date market and economic data from external sources, while also transmitting analysis results and decision recommendations to relevant decision-makers or implementing departments.

2.3.2 Application of decision support system in financial management

The Decision Support System (DSS) offers scientifically grounded decision-making support for enterprises by integrating data analysis, model computation, and user interfaces. DSS is particularly effective in budgeting and control, as it leverages historical financial data and market forecasts to help enterprises develop rational budgets and make dynamic adjustments during implementation, ensuring budget accuracy and execution effectiveness. Kwan et al. demonstrates that computerized DSS can enhance decision-making accuracy, improve efficiency, and minimize human errors [23].

Through automated data processing and analytical functionalities, DSS swiftly generates financial statements and analytical reports, enabling enterprise management to understand financial health and operational outcomes in a timely manner. This automation enhances reporting efficiency and improves data accuracy and transparency.

DSS also identifies, assesses, and manages various financial risks faced by enterprises using an integrated risk assessment model. Lutz et al. noted that DSS performs

strongly in complex analysis and forecasting tasks, aiding enterprises in comprehensive risk management [24]. By analyzing market data and financial indicators, DSS allows for accurate financial risk predictions and provides actionable strategies to help enterprises mitigate potential exposure.

In addition, DSS plays a vital role in investment decisions by integrating diverse investment analysis models. This capability assists companies in evaluating the benefits and risks of different investment scenarios, thereby supporting informed, data-driven investment decisions. DSS effectively processes vast amounts of historical and real-time data, and integrates macroeconomic factors and market trends to provide comprehensive investment analysis and guidance.

3 Research design

3.1 System design principles

System design should prioritize user needs, ensuring usability and operability to allow decision-makers streamlined access to and analysis of financial data. The system requires high-capacity data processing, able to efficiently manage large-scale and diverse financial datasets, thus maintaining data timeliness and accuracy. Advanced data storage and processing technologies, such as distributed databases and parallel computing, should be employed to maximize processing efficiency. The system must exhibit robust scalability, enabling flexible expansion and upgrades to meet evolving enterprise requirements. Multi-level security protocols should be implemented to guarantee the confidentiality and integrity of financial data. Additionally, intelligent analysis and decision-support functionalities should be integrated through data mining, machine learning, and artificial intelligence technologies, facilitating in-depth financial data analysis and informed decision-making.

To ensure clarity and replicability, the methods need precise parameters at each stage of data processing, modeling, and evaluation. For data processing, specific configurations used in the ETL process should be included, such as extraction intervals, transformation rules, and the cleaning criteria. The exact ETL software solution, such as Apache Nifi or Talend, and any SQL or Python scripts used should be specified, detailing the functions or libraries employed for each transformation. In the modeling stage, the algorithms, model parameters, and training settings, such as learning rates or epoch counts in machine learning, should be explicitly listed to allow exact replication. Similarly, evaluation metrics, such as accuracy, recall, or F1 score, and their respective thresholds for decision-making should be clarified. These details ensure that each methodological step is replicable and the results are verifiable, enhancing the study's reliability.

3.2 System architecture design

The architecture design of a financial data analysis and decision support system utilizing big data must thoroughly address data collection, storage, processing, and analysis

to ensure effective data management and intelligent decision support. This system adopts a hierarchical structure, comprising the data layer, processing layer, analysis layer, and application layer. The data layer manages the acquisition and storage of financial data from diverse internal and external sources, using distributed databases to facilitate efficient access and management. The processing layer employs big data technologies, such as Hadoop and Spark, to conduct data cleaning, transformation, and pre-processing, establishing a reliable data foundation for subsequent analysis. In the analytics layer, various data mining and machine learning models are integrated to enable comprehensive analysis of financial data through real-time and batch processing, providing predictive and trend analysis capabilities. Finally, the application layer features a user-friendly interface supporting multiple interaction methods, such as chart presentation, report generation, and customized queries, aiding decision-makers in interpreting and applying the analytical outcomes. The architectural design also emphasizes high availability and scalability, incorporating fault-tolerant mechanisms and dynamic scaling techniques to maintain system stability under high loads and in dynamic environments.

3.3 Data source and data processing

3.3.1 Data acquisition and preprocessing

This study primarily collects a range of financial data from external sources, including market data, economic indicators, industry reports, and social media data. These data are accessed through APIs, web crawlers, and third-party data providers. To maintain data quality and consistency, pre-processing steps are undertaken, involving data cleaning, transformation, and integration. Data cleaning addresses missing values, duplicate entries, and outliers by employing methods such as interpolation, mean replacement, and statistical detection. Data transformation standardizes and normalizes data to remove unit and scale disparities between different sources, while data integration consolidates data into a consistent format to ensure logical and structural uniformity.

In this study, ETL (Extraction, Transformation, Loading) tools are used to automate data acquisition and pre-processing, enhancing processing efficiency. These tools extract data from source systems according to predefined workflows and scripts, clean and transform it, and load it into a data warehouse, establishing a high-quality data foundation for subsequent analysis. This process not only improves data accuracy and processing efficiency but also lays the groundwork for the system's real-time data analysis capabilities.

3.3.2 Data storage and management

The data storage scheme selection in this study carefully considers data characteristics, including scale, structure, access frequency, and security, to address diverse storage requirements effectively. A hybrid storage approach is implemented, integrating both relational databases and

NoSQL databases. This combined storage solution accommodates various data types while ensuring flexibility and robustness to support the study’s analytical objectives, as shown in Table 3 below.

Table 3: Data storage solutions comparison.

Storage Solution	Advantages	Disadvantages	Suitable Scenarios
Relational Database	Strong transaction processing, high data consistency, supports complex SQL queries	Poor scalability, low efficiency in handling large-scale data	Core financial data, structured data
NoSQL Database (MongoDB)	Flexible data model, supports high concurrency and large-scale data storage, strong scalability	Does not support complex transactions, weaker data consistency	Dynamic data, unstructured and semi-structured data
Distributed File System (HDFS)	Efficient handling of large-scale distributed data, supports parallel computing and batch processing	Poor real-time processing, high data access latency	Large-scale data analysis, batch processing tasks

3.3.3 Data cleaning and conversion

In this study, data cleaning includes identifying and handling missing values, duplicate entries, and outliers, while data conversion covers normalization, scaling, and data format transformations.

The initial step addresses missing values. For datasets with relatively few missing entries, records containing missing data are removed. However, when the extent of missing data is significant, and deletion would

compromise analysis, interpolation or mean imputation is applied to estimate missing values. The interpolation method predicts missing values based on trends from adjacent data points, as follows (Formula 4):

$$X_i = \frac{X_{i-1} + X_{i+1}}{2} \tag{4}$$

X_i is the missing value, while X_{i-1} and X_{i+1} represent the previous and subsequent non-missing values, respectively.

The second step involves handling duplicate entries, as duplicates can skew analytical results. Identifying and removing duplicate records ensures data accuracy and uniqueness.

The third step addresses outliers, which are data points that deviate significantly from expected ranges. These outliers are typically detected and managed using the Boxplot method or the standard deviation approach, as shown in (Formula 5):

$$Z = \frac{X - \mu}{\sigma} \tag{5}$$

Z represents the standard fraction, X denotes the data point, μ is the mean, and σ (σ) indicates the standard deviation. When Z exceeds a predefined threshold, these data points are classified as outliers and are either processed or excluded as appropriate.

Data transformation involves both standardization and normalization. Normalization adjusts the data into a standard normal distribution with a mean of 0 and a standard deviation of 1, as demonstrated below (Formula 6):

$$X' = \frac{X - \mu}{\sigma} \tag{6}$$

Normalization adjusts data to fit within the range [0, 1], following this formula (Formula 7):

$$X' = \frac{X - X_{min}}{X_{max} - X_{min}} \tag{7}$$

X represents the original data, while X' denotes the transformed data. Here, X_{min} and X_{max} represent the minimum and maximum values within the data set, respectively. Adjustments are applied to standardize data values between C and D, ensuring consistent scale and reducing variability for more accurate comparative analysis across the data set. This transformation process supports enhanced data normalization and minimizes distortion, especially when dealing with a wide range of values across different categories.

To ensure consistency in data preprocessing, the standardization and normalization techniques are distinctly applied based on the characteristics of specific financial indicators. Standardization is applied to indicators like revenue and net profit, which vary widely across different periods and can benefit from transformation to a standard normal distribution (mean 0, standard deviation 1), allowing consistent scale for predictive modeling. Normalization, on the other hand, is applied to indicators such as liquidity ratios and turnover

ratios, converting their values into a range between 0 and 1 to facilitate comparison across different metrics without scale imbalance. This detailed approach in data preprocessing enhances the model's interpretability, maintains consistency, and improves the model's performance by ensuring that each transformation aligns with the unique properties of the respective financial indicator.

3.4 Module function design

3.4.1 Data analysis module

The data analysis module forms the core of the financial data analysis and decision support system, with its primary function to perform in-depth analysis of collected and stored data, ultimately providing valuable decision-making insights. Python is utilized as the main tool for data analysis, leveraging both data analysis libraries and machine learning libraries to enable efficient processing and analysis. Table 4 below presents the primary functions and advantages of the Python data analysis tools applied in this study.

Table 4: Data analysis tools comparison.

Tool	Main Functions	Advantages
Pandas	Data cleaning, transformation, manipulation; supports DataFrame and Series structures	Easy to use, powerful functions, suitable for structured data
NumPy	Numerical computation, supports large-scale matrix operations and mathematical functions	Fast computation speed, ideal for numerical and matrix operations
Scikit-learn	Machine learning library, includes algorithms for classification, regression, clustering, and dimensionality reduction	Comprehensive functions, easy to use, suitable for various machine learning tasks

3.4.2 Decision support module

The decision support module serves as the core of the financial data analysis and decision support system, primarily designed to provide robust decision support for enterprise management through insights derived from data analysis. This study employs linear regression and decision tree models as key decision support tools, catering to a variety of decision-making scenarios. The functions and advantages of the decision support methods applied in this study are detailed in Table 4 below.

Table 4: Decision support methods comparison.

Method	Main Functions	Advantages
Linear Regression	Predicts and explains the linear relationship between dependent and independent variables	Simple to use, strong interpretability, suitable for linear relationship analysis
Decision Tree	Handles complex decision problems, presents decision paths and results through a tree structure	Easy to understand and interpret, handles nonlinear relationships, high efficiency

3.5 System implementation and deployment

In this study, the system is implemented using a microservice architecture to decouple functional modules, enhancing both flexibility and maintainability. The microservice design deploys and manages each module independently, facilitating interactions through lightweight communication mechanisms. This structure not only improves system scalability but also minimizes the overall impact of potential system failures.

The decision support module incorporates linear regression and decision tree models to provide predictive insights for financial data analysis. However, to enhance accuracy, especially for financial data with non-linear characteristics and temporal dependencies, additional advanced models are considered. Random forest and gradient boosting methods offer greater capacity for handling non-linear data patterns, improving the system's ability to capture complex financial trends. Furthermore, recurrent neural networks (RNN) are introduced to address temporal dynamics, enabling the system to better account for sequence dependencies in financial data. Integrating these models improves the robustness and versatility of the decision support module, ensuring more reliable outcomes in a variety of financial decision-making scenarios.

The experimental setup for implementing and deploying the financial data analysis system includes key specifications on data set size, hardware, software environment, and parameter configurations. The data set comprises financial records totaling 500 GB, processed using Hadoop and Spark for distributed data handling. Hardware configurations include a 64-core CPU with 256 GB RAM and 2 TB SSD storage. Software setup includes Hadoop 3.2.1, Spark 3.1.2, and Docker 20.10 for containerization, running on a Linux-based environment with Ubuntu 20.04. Parameters for Hadoop's HDFS include a block size of 128 MB and a replication factor of three to ensure data reliability. Spark configurations use an executor memory of 8 GB per node and a parallelism setting of 64, optimizing for large-scale data processing. These technical details ensure reproducibility and allow

other researchers to accurately compare and benchmark system performance.

For deployment, containerization technology is employed, with each module packaged into an isolated operating environment using Docker containers. Docker provides a lightweight, consistent, and rapid deployment solution, which enhances the system's stability and deployment efficiency. Additionally, Kubernetes is utilized for container orchestration and management. Kubernetes offers robust features for automated deployment, scaling, and management, dynamically adjusting resource allocation based on system load, ensuring stable operation under high concurrency and heavy traffic.

4 Case study

4.1 Case background

Company A, a medium-sized enterprise established in 2010, specializes in the sale of mobile phones and related electronic products. Headquartered in the United States, the company's primary market comprises major cities in China. Company A's product lineup includes smartphones, tablets, smartwatches, and various smart home devices, with a focus on the mid-to-high-end market. The company is dedicated to building customer trust through technological innovation and high-quality service.

Amid the rapid expansion of the smart device market and rising competition, Company A faces significant market pressures and challenges. To address these, the company has increased investments in research and development, as well as marketing, introducing new products that are competitive within the market. The

company actively works to optimize supply chain management and control costs, which bolsters operational efficiency and enhances its competitiveness and profitability.

Company A's sales revenue and net profit have shown steady growth; however, cost pressures are also gradually rising. The company employs stringent budget control and financial analysis mechanisms, regularly reviewing financial statements and monitoring financial indicators to promptly identify and address operational challenges. The management of current assets and liabilities remains stable, reflected by a high current ratio and quick ratio, which ensures the company's capability to meet short-term obligations.

Looking ahead, Company A aims to further strengthen its market competitiveness and achieve sustainable development by increasing its investment in research and development over the next five years. These investments will drive product innovation and enhance marketing strategies to expand brand influence. Additionally, Company A seeks to elevate its financial management and decision-making processes by leveraging big data analytics and decision support systems, thereby advancing its strategic goals through informed data analysis and decision-making

4.2 Data selection and processing

To ensure the accuracy and depth of the analysis, data covering key financial indicators of Company A, including sales revenue, cost of sales, net profit, current assets, current liabilities, total assets, and shareholders' equity over the past three years, has been selected. The details are presented in Table 5 below.

Table 5: A company's financial data details.

Year	Sales Revenue (million)	Cost of Sales (million)	Net Profit (million)	Current Assets (million)	Current Liabilities (million)	Total Assets (million)	Shareholder's Equity (million)
2021	58.73	34.21	10.45	25.67	14.32	72.45	45.12
2022	63.29	36.89	11.24	27.89	15.67	78.34	48.67
2023	67.54	39.76	12.13	30.21	16.89	83.56	52.34

The selected data undergoes rigorous cleaning and transformation to ensure consistency and reliability. During data cleaning, missing values, outliers, and duplicate entries are addressed: the interpolation method is applied to estimate missing values, while outliers are identified and managed using the boxplot method. In the data transformation phase, all financial indicators are standardized to remove unit and scale discrepancies, thus enhancing comparability across variables and ensuring uniformity.

4.3 Financial risk early warning analysis

This study uses current ratio, quick ratio, asset-liability ratio and other indicators to evaluate the financial risk of

Company A in recent three years. As shown in Figure 2 below.



Figure 2: Financial risk of Company A.

In terms of the current ratio, Company A has maintained a level between 1.78 and 1.79 over the past three years, indicating stable short-term debt repayment ability. A current ratio exceeding 1.5 typically suggests that the company possesses sufficient current assets to meet its short-term obligations.

The quick ratio has shown slight fluctuation over the same period, decreasing marginally from 1.12 in 2021 to 1.10 in 2022, before returning to 1.12 in 2023. With a ratio above 1.0, the company demonstrates solid liquidity. The minor dip in 2022 may reflect short-term variations in inventory or accounts receivable management.

Company A's asset-liability ratio has ranged between 0.59 and 0.60 in recent years, reflecting a sound financial structure with a moderate level of debt. A lower asset-liability ratio suggests that the company benefits from leveraging its own capital for expansion and investment, with minimal debt repayment risk.

The financial risk analysis in this section will incorporate a discussion on the rationale for selecting specific ratios, such as the current ratio and quick ratio, by evaluating their direct impact on assessing liquidity and short-term solvency. The choice of these ratios provides a clear snapshot of the company's ability to meet its short-term obligations. To broaden the depth of analysis, a brief comparison with other common financial risk assessment tools, such as Altman's Z-score, is included. Altman's Z-score, a well-regarded predictor of bankruptcy risk, enables a more comprehensive financial risk evaluation by combining profitability, leverage, liquidity, solvency, and activity ratios. Integrating this perspective will complement the traditional liquidity ratios and provide a more robust framework for assessing Company A's financial stability, supporting a holistic approach to financial risk assessment.

Company A's cash flow ratio has remained between 0.82 and 0.85 over the last three years, suggesting that cash flow from operations sufficiently covers short-term liabilities. Although slightly below the ideal benchmark of 1.0, the company would benefit from enhanced cash flow management to ensure sustained financial stability.

The stability of Company A's key financial indicators over the past three years reveals no significant financial risks. However, it remains essential to monitor fluctuations in the quick ratio and manage the cash flow ratio effectively to maintain long-term financial health and mitigate potential risks.

4.4 Analysis of investment decision support

In this study, the net present value (NPV), internal rate of return (IRR), and payback period (PBP) are applied to assess the company's three investment options, guiding the company in selecting the optimal project. These metrics, shown in Figure 3 below, enable a comprehensive comparison of each scheme's profitability, risk, and time to recoup investment, thus providing a robust basis for informed decision-making on capital allocation.

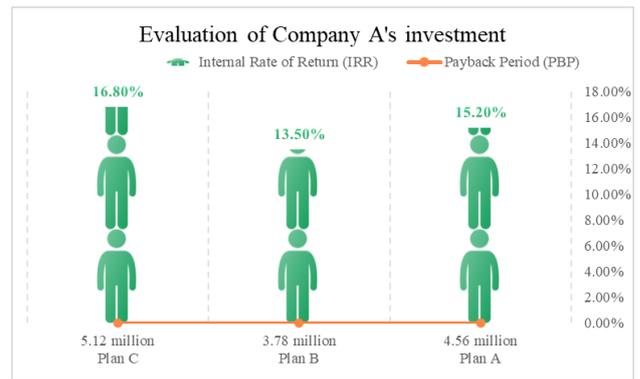


Figure 3: Evaluation of company A's investment plan.

The net present value (NPV) of Plan A is \$4.56 million, with an internal rate of return (IRR) of 15.2% and a payback period of 3.5 years. This option presents a favorable combination of NPV and IRR alongside a moderate payback duration, indicating a strong return on investment within manageable risk parameters. Plan B, on the other hand, has an NPV of \$3.78 million, an IRR of 13.5%, and a payback period of 4.0 years. Compared to Plan A, Plan B offers a slightly lower NPV and IRR, a longer payback period, and a reduced return on investment, carrying a relatively higher degree of risk. Plan C achieves the highest performance across all indicators, with an NPV of \$5.12 million, an IRR of 16.8%, and the shortest payback period at 3.0 years. These results underscore Plan C as the most advantageous option, boasting the highest NPV and IRR and the quickest fund recovery rate.

Therefore, based on the metrics of NPV, IRR, and payback period, Plan C stands out as the optimal investment choice. Plan A ranks as a viable alternative due to its comparatively high return on investment, although it offers a slightly lower IRR than Plan C. Plan B, with its lower NPV, IRR, and longer payback period, is not recommended as the preferred option given its higher associated risks.

4.5 System performance evaluation

To strengthen the performance evaluation of the developed system, a benchmark comparison with state-of-the-art (SOTA) systems is essential. This includes comparing metrics like data processing speed, accuracy, system stability, and user satisfaction against those of current leading systems. Integrating this benchmarking will involve using tables or figures to highlight the developed system's performance in areas such as faster processing speeds, reduced computational costs, or higher system reliability. These metrics, when set against established systems, will clearly demonstrate the advantages and potential improvements achieved by the new system. Such comparative data provides quantitative support for claims of superior performance, offering clearer insights into the system's competitive strengths in financial data processing and decision support.

Comprehensively assess the performance of Company A's financial data analysis and decision support system by evaluating key metrics such as processing speed, data accuracy, system stability, and user satisfaction. This analysis provides insight into the system's operational status and serves as a foundation for future optimization. Details are presented in Table 6 and Figure 4.

Table 6: Performance evaluation of Company A's decision support system.

Evaluation Metric	2021	2022	2023
Data Processing Speed (s)	2.45	2.30	2.15
User Satisfaction (score)	8.2	8.5	8.7
Data Accuracy (%)	98.5	98.7	98.9
System Stability (%)	99.2	99.4	99.5

In terms of data processing speed, Company A's system has shown consistent improvement over the past three years, decreasing from 2.45 seconds in 2021 to 2.15 seconds in 2023. These results indicate the system's robust capability in handling large-scale financial data efficiently, enabling it to complete complex analytical tasks within a reduced timeframe, thereby enhancing overall operational efficiency. Additionally, user satisfaction scores have increased from 8.2 in 2021 to 8.7 in 2023, which reflects the system's continuous optimization in functionality, usability, and performance, contributing to a strengthened user experience and building greater user trust.

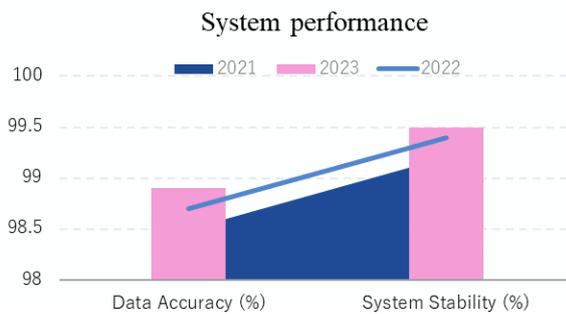


Figure 4: System performance.

The system's accuracy has remained consistently high, increasing from 98.5% in 2021 to 98.9% in 2023. This stability in accuracy ensures the reliability of analysis results, minimizes decision-making errors caused by data discrepancies, and provides a robust data foundation for the company. Operational stability has also improved annually, rising from 99.2% in 2021 to 99.5% in 2023. This trend indicates that the system can maintain steady performance under high-load and extended operation, reducing the likelihood of system failures and downtime, and thus enhancing overall system availability.

In summary, Company A's financial data analysis and decision support system demonstrates strong performance

across key metrics, including data processing speed, accuracy, stability, and user satisfaction, fully meeting anticipated performance standards. Continuous monitoring and optimization further enhance the system's capacity to support enterprise financial management needs, offering scientifically reliable decision-making support in a dynamic business environment.

The system performance analysis now incorporates specific criteria to improve clarity and replicability. User satisfaction is evaluated using a structured questionnaire that includes Likert-scale items to measure aspects like system usability, functionality, and efficiency. The questionnaire was distributed to 50 end-users, with results analyzed to generate an average satisfaction score. Processing speed is quantified through system logs that capture average response times during peak and non-peak hours, ensuring consistency in data handling performance. Data accuracy is assessed by tracking error rates in data processing tasks, providing a measure of reliability. Finally, system stability is monitored through uptime statistics and error frequency logs, ensuring sustained performance under variable loads. Including the questionnaire in the appendix enhances transparency, while quantitative metrics solidify the reliability of the performance evaluation.

5 Discussion

5.1 Result discussion

The results show that the financial data analysis and decision support system designed with big data technology has significantly improved the accuracy and efficiency of financial management. By processing large-scale data swiftly, the system enhances data accuracy, reaching up to 98.9% over a three-year evaluation. The system's financial risk early warning function effectively identifies potential risks by monitoring indicators such as the current ratio, quick ratio, and asset-liability ratio. Moreover, the decision support component has helped evaluate investment options through key metrics such as net present value (NPV), internal rate of return (IRR), and payback period (PBP). The system recommended the most favorable investment with an NPV of 5.12 million, IRR of 16.8%, and a payback period of 3.0 years, providing reliable decision-making support for enterprise investment strategies and financial risk control.

The results are compared with the state-of-the-art (SOTA) benchmarks to contextualize the system's performance within current methodologies. Specifically, the system's data processing speed, accuracy, and stability are evaluated against leading benchmarks, highlighting areas of alignment and discrepancy. This analysis reveals that while the current system achieves high processing speed and accuracy, its performance is influenced by factors such as data quality, algorithm efficiency, and scalability. For instance, improvements in data quality directly enhance accuracy, while algorithm efficiency impacts processing speed. The scalability of the architecture, enabled by big data tools like Hadoop and Spark, allows the system to maintain performance with

increasing data volume. These results underscore the system's strong positioning relative to SOTA methods, with insights into areas for further refinement.

The detailed analysis of Company A's financial data and decision support system demonstrates the tangible impact of big data technology on financial management. The findings indicate that this financial data analysis and decision support system provides marked advantages in increasing data processing efficiency, enhancing decision accuracy, and optimizing financial management processes.

In terms of data processing speed and accuracy, the system exhibits robust performance with large-scale financial data, completing complex analytical tasks swiftly. High data accuracy contributes to reliable analysis outcomes, reducing errors in decision-making due to data inaccuracies. Implementing big data technology thus improves the efficiency and accuracy of financial data handling, supplying enterprises with a dependable data foundation.

The system's effectiveness in financial risk early warning is also notable. By monitoring and analyzing Company A's key financial indicators, the system identifies potential financial risks promptly, issuing early warnings that enable timely preventive measures. This improves the company's risk management capabilities, supporting a stable financial position amid intense market competition.

Investment decision support analysis results further show that the system can objectively assess various investment options using financial indicators such as net present value (NPV), internal rate of return (IRR), and payback period (PBP), assisting companies in selecting the most advantageous investment projects. This robust analytical support promotes efficient fund utilization, maximizing return on investment. Additionally, performance evaluations verify the system's stability and user satisfaction. The system maintains reliable performance under high loads and continuous operation, reducing system failures and downtime, thereby enhancing availability. Increased user satisfaction reflects ongoing system optimization in functionality, usability, and performance, fostering user trust and engagement.

5.2 Enlightenment and suggestions for enterprise financial management

5.2.1 Improve data processing and analysis capabilities

Enterprises are increasingly adopting big data technology to strengthen financial data processing and analysis capabilities. By employing advanced data storage, processing, and analysis tools, they achieve efficient management and in-depth examination of large volumes of financial data, thereby providing more accurate and timely support for decision-making. This approach improves data processing efficiency, enhances data accuracy and reliability, and supports enterprises in sustaining a competitive edge within a highly competitive market environment.

5.2.2 Strengthen risk management

Big data technology is highly effective in providing early warnings for financial risks, enabling enterprises to establish robust financial risk early warning systems. Such systems allow companies to detect and mitigate potential risks in a timely manner through real-time monitoring and analysis of essential financial indicators. This capability strengthens enterprises' risk management frameworks, fortifying their financial stability and responsiveness to unexpected events. Regular risk assessments enable enterprises to adjust financial strategies proactively, thereby maintaining financial security and resilience.

5.2.3 Optimize the investment decision-making process

Enterprises can leverage big data technology to conduct a systematic evaluation of investment projects, allowing for the selection of optimal investment options. By using comprehensive indicators such as net present value (NPV), internal rate of return (IRR), and payback period (PBP), investment returns and risks can be more accurately forecasted, leading to more informed decision-making. This approach enhances investment efficiency, maximizes return on investment, and facilitates optimal resource allocation, thereby supporting strategic financial management and risk mitigation.

5.2.4 Improving system stability and user satisfaction

The stability and user satisfaction of a financial management system are critical to its success, as enterprises must continually enhance system functionality to ensure consistent performance under high load and extended operation periods. Regular improvements to system capabilities not only enhance the user experience but also build user trust and reliance on the system, thus increasing its effectiveness and adoption. Enterprises can benefit from systematically conducting system performance evaluations and user feedback surveys to inform continuous improvements and refine system functionality over time, ensuring that the system aligns with user expectations and operational demands.

5.3 Research limitations and future prospects

5.3.1 Research limitations

This study achieved notable progress in applying big data technology to financial data analysis and decision support; however, limitations remain. Despite implementing various data cleaning and preprocessing techniques, data accuracy and consistency are constrained by the limitations of data acquisition channels and the inherent quality of the data itself. Additionally, this study employs linear regression and decision tree models, which may not fully capture the intricacies of certain financial data types and thus may not address all financial management scenarios. Furthermore, the system's performance

evaluation relies primarily on simulated data within an experimental setting, meaning that real-world complexity and variability could influence system performance and stability, necessitating further adjustments in practical applications.

5.3.2 Future outlook

Future research should enhance data collection and cleaning techniques to ensure higher data accuracy and consistency. Expanding data sources and improving data processing algorithms will be essential in raising data quality. In terms of model optimization and innovation, exploring additional analysis models and algorithms—such as deep learning and reinforcement learning—will be beneficial for enhancing the system's analytical and predictive capabilities in various financial management contexts. Furthermore, potential applications of big data technology in financial management, including real-time financial monitoring, intelligent report generation, and automated financial decision-making, should be further explored. Integration with emerging technologies like artificial intelligence and blockchain can advance the intelligence and security of financial management systems.

Focusing on improving user experience is also crucial; enhancing the user interface and interaction features will increase the system's usability and user satisfaction. Future research should emphasize interdisciplinary collaboration, combining insights and techniques from fields such as financial management, information technology, and data science to drive continuous innovation in financial data analysis and decision support systems. Ongoing technological advancements and expanding applications will offer more scientifically robust, intelligent, and efficient solutions for enterprise financial management.

References

- [1] Casturi R, Sunderraman R (2022). Cost effective, rule based, big data analytical aggregation engine for investment portfolios. *Wireless Networks*. 28(3):1203-1209. <https://doi.org/10.1007/s11276-018-01904-5>.
- [2] Fahd K, Miah SJ (2023). Designing and evaluating a big data analytics approach for predicting students' success factors. *Journal of Big Data*. 10(1):159. <https://doi.org/10.1186/s40537-023-00835-z>.
- [3] Dutta G, Rao HV, Basu S, Tiwari MK (2019). Asset liability management model with decision support system for life insurance companies: Computational results. *Computers & Industrial Engineering*. 128:985-998. <https://doi.org/10.1016/j.cie.2018.06.033>.
- [4] Lu J, Wu D, Dong J, Dolgui A (2023). A decision support method for credit risk based on the dynamic Bayesian network. *Industrial Management & Data Systems*. 123(12):3053-3079. <https://doi.org/10.1108/IMDS-04-2023-0250>.
- [5] Talamo A, Marocco S, Tricol C (2021). "The Flow in the Funnel": Modeling Organizational and Individual Decision-Making for Designing Financial AI-Based Systems. *Frontiers in Psychology*. 12:697101. <https://doi.org/10.3389/fpsyg.2021.697101>.
- [6] Peng J, Bao L (2023). Construction of enterprise business management analysis framework based on big data technology. *Heliyon*. 9(6):e17144. <https://doi.org/10.1016/j.heliyon.2023.e17144>.
- [7] Popovic A, Hackney R, Tassabehji R, Castelli M (2018). The impact of big data analytics on firms' high value business performance. *Information Systems Frontiers*. 20(2):209-222. <https://doi.org/10.1007/s10796-016-9720-4>.
- [8] Liu M, Gao R, Fu W (2021). Analysis of Internet Financial Risk Control Model Based on Machine Learning Algorithms. *Journal of Mathematics*. 2021:8541929. <https://doi.org/10.1155/2021/8541929>.
- [9] Yu H, Guo Z (2022). Design of Underground Space Intelligent Disaster Prevention System Based on Multisource Data Deep Learning. *Wireless Communications & Mobile Computing*. 2022:3706392. <https://doi.org/10.1155/2022/3706392>.
- [10] Zheng T, Chen G, Wang X, Chen C, Wang X, Luo S (2019). Real-time intelligent big data processing: technology, platform, and applications. *Science China-Information Sciences*. 62(8):082101. <https://doi.org/10.1007/s11432-018-9834-8>.
- [11] Xue CJ (2019). Special Issue on Emerging Technologies for Big Data Processing. *Journal of Electronic Science and Technology*. 17(1):1-2. <https://doi.org/10.3969/j.issn.1674-862X.2019.01.001>.
- [12] Xu Z, Zhou W (2021). A Data Technology Oriented to Information Fusion to Build an Intelligent Accounting Computerized Model. *Scientific Programming*. 2021:6031324. <https://doi.org/10.1155/2021/6031324>.
- [13] Rodrigues M, Santos MY, Bernardino J (2019). Big data processing tools: An experimental performance evaluation. *Wiley Interdisciplinary Reviews-Data Mining and Knowledge Discovery*. 9(2):e1297. <https://doi.org/10.1002/widm.1297>.
- [14] Jensen MH, Nielsen PA, Persson JS (2023). From Big Data Technologies to Big Data Benefits. *Computer*. 56(6):52-61. <https://doi.org/10.1109/MC.2022.3206032>.
- [15] Li K, Chen Y (2021). Fuzzy Clustering-Based Financial Data Mining System Analysis and Design. *International Journal of Foundations of Computer Science*. 33(06N07):603-624. <https://doi.org/10.1142/S0129054122420060>.
- [16] Wang S (2020). Research on Data Mining and Investment Recommendation of Individual Users Based on Financial Time Series Analysis. *International Journal of Data Warehousing and Mining*. 16(2):64-80. <https://doi.org/10.4018/IJDWM.2020040105>.

- [17] Guo Y (2021). CNS: Interactive Intelligent Analysis of Financial Management Software Based on Apriori Data Mining Algorithm. *International Journal of Cooperative Information Systems*. 30(1-4). <https://doi.org/10.1142/S0218843021500088>.
- [18] Liang M (2021). Optimization of Quantitative Financial Data Analysis System Based on Deep Learning. *Complexity*. 2021:5527615. <https://doi.org/10.1155/2021/5527615>.
- [19] Lin Y, Yue H, Liao H, Li D, Chen L (2022). Financial Risk Assessment of Enterprise Management Accounting Based on Association Rule Algorithm under the Background of Big Data. *Journal of Sensors*. 2022:8041623. <https://doi.org/10.1155/2022/8041623>.
- [20] Zhai Z, Martinez JF, Beltran V, Martinez NL (2020). Decision support systems for agriculture 4.0: Survey and challenges. *Computers and Electronics in Agriculture*. 170:105256. <https://doi.org/10.1016/j.compag.2020.105256>.
- [21] Sutton RT, Pincock D, Baumgart DC, Sadowski DC, Fedorak RN, Kroeker KI (2020). An overview of clinical decision support systems: benefits, risks, and strategies for success. *NPJ Digital Medicine*. 3(1):17. <https://doi.org/10.1038/s41746-020-0221-y>.
- [22] Loftus TJ, Tighe PJ, Filiberto AC, Efron PA, Brakenridge SC, Mohr AM, Rashidi P, Upchurch GR Jr, Bihorac A (2020). Artificial Intelligence and Surgical Decision-making. *JAMA Surgery*. 155(2):148-158. <https://doi.org/10.1001/jamasurg.2019.4917>.
- [23] Kwan JL, Lo L, Ferguson J, Goldberg H, Diaz-Martinez JP, Tomlinson G, Grimshaw JM, Shojania KG (2020). Computerised clinical decision support systems and absolute improvements in care: meta-analysis of controlled clinical trials. *BMJ*. 370:m3216. <https://doi.org/10.1136/bmj.m3216>.
- [24] Lutz W, Deisenhofer AK, Rubel J, Bennemann B, Giesemann J, Poster K, Schwartz B (2022). Prospective evaluation of a clinical decision support system in psychological therapy. *Journal of Consulting and Clinical Psychology*. 90(1):90-106. <https://doi.org/10.1037/ccp0000642>.

Measuring Fidelity of Steganography Approach in Securing Clinical Data Sharing Platform using Peak Signal to Noise Ratio (PSNR) and Structural Similarity Index Measure (SSIM)

A.M. Adeshina, Siti Fatimah Abdul Razak, Sumendra Yogarayan, Md Shohel Sayeed
Faculty of Information Science & Technology, Multimedia University, Malaysia
Email: am.adeshina@mmu.edu.my, codedengineer@yahoo.com, fatimah.razak@mmu.edu.my, sumendra@mmu.edu.my, shohel.sayeed@mmu.edu.my

Keywords: steganography, cryptography, steganalysis, peak signal to noise ratio, structural similarity index measure

Received: January 23, 2024

In the vast digital landscape, the practice of data hiding finds multifaceted applications, ranging from simple hobbyist endeavors to critical tasks like safeguarding user privacy and ensuring covert data transmission. One of the gaping vulnerabilities in many contemporary systems is the transparency with which information is stored, making it easily interpretable. Such clear visibility can be a gateway for potential leaks, false portrayals, or even be manipulated for various malevolent intents. Consequently, as a countermeasure, steganography emerges at the forefront, extensively being resourceful in the revolutionized data storage concept, the cloud technology. Unfortunately, most earlier image steganography methods could only conceal one type of file, audio, text, image within an image, rendering them monodynamic. This study focuses on the novel application of steganography towards embedding information across multiple images to facilitate security of clinical data sharing platform as opposed to traditional single-image methods. The implementation was carried out using Ruby on Rails architecture, leveraging the ChunkyPNG library. With the analyses of image texture features, adaptive payload distribution strategies were devised and compared with the established single-image steganographic techniques. Interestingly, our findings show employing strategies based on texture complexity and distortion distribution greatly enhances security, making it more resilient to modern pooled steganalysis. The exceptionally high PSNR values consistently above 90dB, coupled with SSIM values nearing 1, collectively underscore the near-identical nature of our original and stego images. This convergence of both metrics emphasizes the effectiveness of our steganographic methods, suggesting minimal distortions and high fidelity. Such compelling outcomes not only validate the methodology employed but also accentuate its potential for applications demanding subtle data concealment. In essence, the combined insights from PSNR and SSIM robustly affirm the project's success in achieving high-quality steganographic results.

Povzetek: Študija uporablja steganografijo za izboljšanje varnosti platforme za izmenjavo kliničnih podatkov. Predlaga se nova metoda vdelave informacij v več slik (namesto tradicionalne v eno) z uporabo arhitekture Ruby on Rails in knjižnice ChunkyPNG. Razvite so prilagodljive strategije porazdelitve koristnega tovara na podlagi kompleksnosti tekture in porazdelitve popačenj.

1 Introduction

Steganography, being defined as both an art and a science, is all about concealed communication, aiming to embed written content within other unassuming data. This technique ensures that the actual embedded information remains inconspicuous. In the modern age, as information storage transitions to being predominantly digital due to the surge in ICT advancements, the importance and utility of such covert communication methods have witnessed a meteoric rise.

The beauty of steganography lies in its subtlety, it can seamlessly embed either a straightforward message or an encrypted one within a digital host file. This

ensures the encoded message remains discreet, especially during transmission across digital networks. While cryptography stands as the age-old bulwark of information security, steganography introduces an added layer, bringing more depth to the protection, especially in the realms of digital media copyrights. The field of information hiding, ever-evolving, casts a wide net over a myriad of applications, be it watermarking, fingerprinting, or the discreet art of steganography. Watermarking, for instance, leans more toward embedding pertinent data such as owner credentials or specific timestamps to thwart potential copyright infringements. Conversely, fingerprinting is all about integrating a unique serial identifier into a

dataset, helping monitor and curb unauthorized exploitation of the same. The digital age, characterized by its unbridled communication channels, has underscored the criticality of enhanced security protocols, particularly within interconnected networks. With the world increasingly becoming interconnected and the volume of data exchanges skyrocketing, the imperative of ensuring confidentiality and maintaining data integrity has never been higher. This escalating concern has been a driving force behind the robust evolution and development of intricate information-hiding methodologies.

Steganography, at its core, is the subtle art of concealing information within digital media. On the other hand, cryptography delves deep into the complex realm of encoding information, employing an array of intricate techniques to ensure it remains inaccessible to unauthorized users. While cryptography shoulders the responsibility of preserving communication confidentiality, steganography thrives on maintaining the secrecy surrounding the very existence of concealed information. As our world increasingly embraces electronic communication and relies heavily on the vast infrastructure of the internet, the imperative for robust information security escalates. Traditional cryptography, though a stalwart in its domain, focuses predominantly on safeguarding content. However, in certain situations, the need transcends mere content protection; sometimes, the very revelation that a hidden message exists can be detrimental. Here, steganography fills the void, it masterfully embeds information within everyday digital media, be it images, videos, or audio files, thereby evading undue attention or suspicion. Given its multifaceted applications, ranging from digital media copyright protection to watermarking and fingerprinting, steganography's importance is undeniable. Its relevance only grows in our digital age, where network security emerges as a paramount concern, leading to the rapid evolution of the broader field of information hiding, encompassing both cryptography and steganography.

In contemporary society, steganography's identity is deeply intertwined with digital data carriers and the pulsating rhythm of high-speed network communications. Drawing a comparison with cryptography, the distinction becomes clear: while cryptography aims to shield the content of a message, steganography thrives on obscuring the message's very existence. Both technologies, with their unique strengths, have carved out essential roles in the overarching goal of data protection. However, like all technologies, neither steganography nor cryptography is a silver bullet; each has its vulnerabilities and can potentially be compromised under specific scenarios. A significant challenge for steganography is that once the clandestine nature of the embedded information is suspected or, worse, unearthed, its primary objective is immediately jeopardized. Yet, the dynamic interplay between steganography and cryptography offers a promising avenue. By synergistically combining these two methods, one can amplify the effectiveness of

steganographic techniques, ensuring a more robust and layered approach to securing sensitive information.

Several efforts were frequently into tackling secure transmission of images containing embedded data over a network by image steganography. More attention was recently on further resolving challenges of earlier image steganography methods on only able to conceal one type of file (e.g., audio, text, image, etc.) within an image. As a result, recent studies aim to create an image steganography system capable of concealing text and image within an image. The aim of this study is to develop an image steganography system that hides either text or image files within an image by proposing a framework for an image steganography model, design and implement the framework proposed and evaluate the designed and implemented framework using PSNR (Peak Signal-to-noise Ratio) and Structural Similarity Index Measure (SSIM).

2 Related works

2.1 Data sharing platform

A data sharing platform is a vital technological system that underpins the seamless exchange, collaboration, and dissemination of data among diverse entities, from organizations to individual researchers. Such platforms are pivotal in converting expansive data into actionable insights, particularly in the era of Industry 4.0. Karabacak et al. (2022) shed light on this notion by proposing a unique document-based data-sharing platform software architecture. This architecture is meticulously designed to address the intricate challenges tied to the analysis of vast data sets, with a particular focus on metadata management, which serves to thwart data complexity while elevating its usability. Concepts of connected networks (Adeshina & Hashim, 2017), which have now being seen quite resourceful very recently will immensely benefit from improved data architecture.

At the heart of this architecture lies a sophisticated metadata store, equipped with a suite of tools tailored for data owner identification, intricate versioning processes, and thorough lineage tracking. The architecture doesn't just stop there; it prioritizes data accessibility by presenting detailed illustrations that pinpoint critical data locations. This emphasis on accessibility is seamlessly complemented by robust mechanisms to uphold data quality, encompassing user-centric data preprocessing techniques. Furthermore, to fortify the system against potential security vulnerabilities, the architecture integrates rigorous operational security controls and an astute user group management framework.

Delving deeper into the functionalities, the software architecture refines data management by classifying information into stochastic data sets, thereby offering role-tailored suggestions to its users. It adopts a dynamic version and rule adaptation methodology, ensuring the platform remains resilient

to evolving data landscapes. Moreover, a bespoke rule customization mechanism stands ready to cater to specific user-driven requirements. In their comprehensive exploration, Karabacak et al. (2022) elucidates the nuances of this document-based data-sharing platform, accentuating its pivotal role in championing efficient data management and fostering collaboration.

Research data sharing platforms, as detailed by Hahnel (2023), represent pivotal online systems developed to bolster the storage, management, and dissemination of research data among the global scientific fraternity. These digital infrastructures serve as linchpins, championing the virtues of transparency, seamless collaboration, and reproducibility of pivotal research findings. As such, researchers are endowed with a unified, secure sanctuary, allowing them to store and propagate their research data. This centralized approach not only amplifies access but also paves the way for potential data reuse by the broader research community.

In the intricate landscape of clinical cohort studies, researchers delve into the life histories of population groups, seeking understanding of disease progression Vilaza et al. (2020). Newly minted health research data platforms have revolutionized this process, granting unprecedented access to cohorts' non-identifiable health details, with cutting-edge initiatives even assimilating mobile-generated data.

In the vast digital ecosystem, certain platforms are tailor-made to address the unique needs and nuances of specific industries, be it healthcare, finance, or the multifaceted world of agriculture (Yoon et al., 2018). Rather than adopting a one-size-fits-all approach, these specialized platforms zero in on the intricate challenges and opportunities inherent to their respective sectors. Through that, they become invaluable conduits, seamlessly facilitating data sharing among myriad organizations nestled within a particular industry. The ripple effect of such targeted data exchange is profound. Not only does it foster an environment conducive to collaboration, but it also paves the way for robust benchmarking exercises.

2.2 Steganographic procedures

Steganography stands as a nuanced technique, meticulously designed to clandestinely embed secret information within digital media, predominantly images, ensuring such embeddings fly under the radar of unintended observers. The advent and meteoric rise of cloud technology have significantly transformed the digital storage landscape, with cloud storage platforms becoming the de facto choice for housing vast repositories of digital images. This proliferation of cloud-based image storage has inadvertently spawned an exciting avenue for steganography. Now, instead of being confined to embedding information in a solitary image, the technique can be scaled to span multiple images, marking a paradigm shift from traditional

single-image steganographic methods.

In this evolving context, Liao et al. (2022), through an insightful publication in the IEEE Transactions on Dependable and Secure Computing, delve deep into the intricacies of optimally allocating embedding payload across a sequence of images, all with an overarching goal to bolster security efficacy in this new era of multiple image steganography.

Two distinct payload distribution blueprints emerge from Liao's research. The inaugural strategy is anchored in image texture complexity, wherein the embedding payload distribution is meticulously choreographed in sync with the image's unique texture attributes. In contrast, the secondary strategy pivots towards distortion distribution, keenly focusing on distributing the payload in alignment with the distortions birthed during the embedding phase. Such strategies, as proposed, don't exist in isolation; they are adeptly designed to coalesce with cutting-edge single image steganographic algorithms, thereby amplifying their inherent security attribute.

Image Steganography stands as an artful technique of discreetly embedding information—be it text, image, or video—within a primary or cover image. Executed with finesse, this embedded information remains invisible to the naked eye, ensuring the secrecy of the data. With technological advancements, particularly the emergence of deep learning technology, steganography has undergone significant evolution. Deep learning, having etched its mark in diverse applications, is now making inroads into the domain of image steganography, attracting a surge of research interest Subramanian et al. (2021). The crux of Subramanian's exploration lies in dissecting and elucidating the myriad deep learning methods prevalent in the field of image steganography.

In the realm of secure communication, various methods are employed to ensure the confidentiality of information exchanged through different channels such as phones, faxes, computer communications, and radio. Steganography offers three primary types which are Pure Steganography, Private Key Steganography, and Public Key Steganography.

Pure steganography emerges as a distinctive approach, concentrating on the concealment of information within digital media, devoid of any reliance on cryptographic techniques or password defenses. Sharma (2017) delves deep into this concept, advocating for its potential as a singular strategy to bolster information security amidst the plethora of existing mechanisms. In today's digital epoch, myriad security protocols and algorithms are enlisted to shield data from unauthorized access and potential cyber threats. Cryptography, with its robust frameworks, often stands at the forefront of such defenses, recognized widely for its efficacy. Yet, pure steganography diverges, aspiring for stealth without

leaning on the cryptographic pillar.

Private key steganography emerges as a transformative approach, amplifying data steganography's security threshold by integrating a private key, serving as an auxiliary encryption stratum. Alqadi's exploration (Alqadi, 2020) unveiled in the *International Journal of Engineering Technologies and Management Research*, charts a course towards enhancing the security matrix of the well-trodden LSB2 (Least Significant Bit 2) method - esteemed for its capacity to cloak secret dispatches in digital color canvases. While LSB2 has carved its niche for preserving the host image's pristine quality even as it harbors clandestine messages, its inherent simplicity has left it vulnerable, often placing it in the crosshairs of hacking endeavors. Alqadi charts a pioneering pathway, offering a remedy in the form of a private key mechanism to galvanize the security bastion of the LSB2 methodology. At the heart of this proposition lies the extraction of a bespoke key from the host image, functioning as the linchpin for secret message encryption.

Public key steganography (PKS) is an intricate merger of steganography with the tenets of public key cryptography. Casting a spotlight on this evolving intersection is the review contributions of Abdul-Razak et al. (2018) which meticulously curated to furnish readers with a holistic perspective, spanning the characteristic features, rich content, and evaluation matrices intrinsic to PKS. Central to the discourse are three pillars: the multifaceted domains where PKS finds application, the diverse schemes championing its cause, and the critical yardsticks employed to gauge the efficiency of PKS infrastructures. Through this tri-pronged lens, Abdul-Razak et al. (2018) dissects and compartmentalizes findings, bequeathing a structured, panoramic view of the PKS landscape. This methodical exploration is not just a mere documentation; it emerges as a treasure trove, brimming with insights, primed to enrich researchers and aficionados venturing into the PKS domain.

2.3 Steganographic filing methods

Text steganography, Protocol steganography, Audio steganography, Image steganography, and video steganography, are the primary categories of file formats commonly used in steganography.

Text steganography has risen as a key technique for discreetly embedding messages within textual documents. Majeed and the team (Majeed et al., 2021) emphasize in their insightful analysis featured in the *Journal of Mathematics*. The technique provides an overview of the intricate methodologies, a litany of challenges faced, and potential future trajectories in this unique field of study. While encryption methods, like cryptography, often shoulder the brunt of data protection efforts, Majeed et al. (2021) contend that steganography presents an unparalleled approach, deftly interweaving hidden messages within overt

narratives or other cover media.

The data Protocol steganography has emerged as an innovative mechanism, intricately crafting covert channels within the lattice of network protocols, offering a secure conduit for the surreptitious transmission of privileged information. Alishavandi & Fakhredanesh (2021) pave the way for an avant-garde approach, christened as Master Key Identifier based Protocol Steganography (MKIPS). This groundbreaking methodology is intricately tailored for the Secure Real-time Transfer Protocol (SRTP), a cornerstone in the realm of Voice-over-Internet Protocol (VoIP) communications.

Venturing into the mechanics of MKIPS, it brilliantly harnesses the sender's prerogative to cherry-pick a master key from a meticulously curated reservoir of cryptographic keys. These keys are gracefully presented by an external key management protocol during the crucial phase of session initiation. Through astute manipulation of the master key identifier field as part of the SRTP packet orchestration, Alishavandi & Fakhredanesh (2021) delineate the establishment of a covert channel. Impressively, this clandestine conduit seamlessly operates beneath the canopy of the SRTP channel, showcasing an admirable bandwidth, its prowess dictated by the inherent characteristics of the SRTP channel in operation.

Audio steganography, at its core, intricately embeds covert messages within the vast expanse of audio data, thus offering a fortified layer of security during data transmission (Abdulkadhim & Shehab, 2022). As featured in the *International Journal of Electrical and Computer Engineering (IJECE)*, meticulously sketches out a crypto-steganographic blueprint tailored to surreptitiously nestle an audio or voice message within two distinct cover media forms, specifically audio and video. The ingenuity of this method stems from its harmonious melding of the least significant bits (LSB) algorithm and the intricate 4D grid multi-wing hyper-chaotic (GMWH) system. Initially, an audio undergoes a transformative shuffle orchestrated by a key birthed from the GMWH system. This meticulous shuffle not only introduces a layer of intricate complexity but also fortifies the audio against the prying eyes of hackers, making the extraction of the original composition a daunting task. The empirical results underscore the method's superiority, showcasing its enhanced security pedigree in juxtaposition to its contemporaries.

Image steganography has transformed from a mere technique into a crucial facet of data security, significantly bolstered by the proliferation of cloud technology and its expansive cloud storage possibilities. This progression enables a shift from the conventional single-image steganography to an innovative approach where skilled steganographers judiciously embed covert information across a multitude of digital images. Consequently, these adapted payloads, spread across images, culminate in

a cryptic array of concealed data, primed for cloud-based transmission to the intended audience.

Video steganography, a pivotal instrument in the cybersecurity toolkit, ensures that confidential information remains clandestinely embedded within video files, thereby bolstering data transmission security. A major caveat, however, is the discernible limitations of conventional algorithms, which grapple with sluggish convergence rates, illuminating the pressing need for a more adept algorithmic framework. In response to this exigency, Salunkhe & Bhosale (2022) unveils an innovative algorithm coined as the Water-Earth Worm Optimization (WEWO) in a seminal paper published in the *International Journal of Engineering Science and Technology*. This avant-garde algorithm, the product of an intricate amalgamation of the Water wave optimization (WVO) and Earth worm optimization (EWO) model algorithms, emerges as a potential game-changer in the realm of video steganography. As part of its modus operandi, the video frames are subjected to rigorous preprocessing and extraction processes, leveraging the capabilities of Discrete Cosine Transform (DCT) and Structured Similarity Index (SSIM) techniques.

For the pivotal task of pixel prediction, an astutely designed fitness function-birthered from neighborhood entropies-becomes the cornerstone of the proposed algorithm. Herein, the surreptitious embedding of the covert message is accomplished via a meticulous two-tier decomposition process hinged on Wavelet Transform (WT). To critically evaluate the WEWO-Deep RNN algorithm's mettle, comprehensive experiments were orchestrated utilizing the 'CAVIAR' dataset. Rigorous tests assessing the algorithm's resilience against modular perturbations such as salt, pepper, and combined noises were conducted. Drawing from quantifiable metrics like Peak Signal-to-Noise Ratio (PSNR), Mean Square Error (MSE), and Correlation Coefficient (CC), which are indispensable for gauging image quality, the results gleaned from Salunkhe & Bhosale research emphatically spotlight the WEWO algorithm's superior prowess in seamlessly embedding encrypted messages without compromising the overarching video quality.

2.4 Watermarking and fingerprint

Watermarking and Fingerprinting are two related steganographic technologies that are often used in the protection of intellectual property.

The heatmap function watermarking, as a technique, embeds imperceptible yet robust data within digital media, be it images, audio, or videos. Its primary function serves as a mechanism for copyright protection, authentication, and ownership verification, ensuring media content is traceable even if illicitly altered or redistributed. A significant stride in this domain has been the amalgamation of deep learning into watermarking methods. With the submission of Li (2021), deep learning-based watermarking techniques

are explored at length, elucidating their strengths and potential constraints. The study espouse the potential of these techniques, particularly emphasizing how deep neural networks amplify the robustness and security of the embedded watermarks. The efficacy of deep learning in this domain suggests its pivotal role in ensuring watermarks remain undisturbed, even when the media is subjected to modifications. This innovation not only ensures the robustness of watermarking but also sets new paradigms for securing digital media.

Shah & Prakash (2020) delves into the intricacies of watermarking algorithms tailored explicitly for images. Recognizing the limitations of traditional LSB-based watermarking techniques, The researchers propose an enhanced method, intertwining Discrete wavelet Transform (DWT), discrete cosine transform (DCT), and singular value decomposition (SVD) to bolster the protection of copyrighted images. The novelty of the proposed technique lies in its remarkable robustness, ensuring that watermarks remain resilient against potential alterations. Furthermore, they prioritize the aesthetic aspect, ensuring minimal perceptual distortion while maintaining a robust security shield for the images. This dual emphasis on both security and image quality underscores the importance of their research, making the findings indispensable for those looking to strike a balance between protection and presentation.

Given the critical nature of medical data, watermarking in the realm of medical imaging demands specialized attention. Chugh & Vashishth (2020) undertake this responsibility, offering a comprehensive examination of digital watermarking techniques custom-built for medical images. The discourse spans across the multifaceted dimensions of watermarking, be it security, tamper detection, or authentication. The research accentuates the bespoke challenges posed by the medical sector, underscoring the need for watermarking methods that cater to the specific requirements and nuances of medical images.

3D watermarking, given its unique challenges and potentials, stands apart in the vast panorama of watermarking research. Cao et al. (2019) explores this niche by introducing an adaptive watermarking blueprint crafted for 3D point clouds. Their approach, tailored to resonate with the distinct attributes of point cloud data, ensures the watermark's seamless embedding and extraction. By integrating an understanding of 3D data's specificities into the watermarking methodology, they have carved out a pioneering path, bridging the chasm between 3D data's potential and the imperatives of robust watermarking.

Dwelling on the foundational aspect of watermarking, Meng & Huang (2018) proffer a blind watermarking blueprint hinging on block Discrete Cosine Transform (DCT) specifically designed for images. The innovation, however, doesn't just rest

with the embedding. The true novelty of their method is its recovery capability, ensuring watermarks remain extractable even in the absence of the original host image.

Digital fingerprinting, known as content-based fingerprinting, is a crucial technique used to trace specific digital media instances. By creating a unique "fingerprint" based on inherent features within media, such as images, distinct audio patterns, or video sequences, it offers a novel way to identify and track content (Suhail & Abhayaratne, 2018). This method extends its utility beyond mere identification, proving instrumental in tracing unauthorized copies and managing media distribution. The entire fingerprinting process is methodical.

Shifting focus to image fingerprinting, visual cues like color histograms and intricate texture patterns become pivotal (Casey & Veltkamp, 2019). Such markers form the essence of the extraction process, is consistent with, the methodologies of Suhail & Abhayaratne (2018), they are molded into fingerprints. These fingerprints are then archived for future reference and juxtaposed against new media fingerprints for identification. With fingerprinting techniques extending their influence to areas like digital rights management and plagiarism detection, the digital domain has acquired a structured mechanism to manage and protect its vast media assets.

2.5 Applications of steganography

Steganography plays a pivotal role in the domain of covert communication, acting as a safeguard for sensitive data. By embedding classified information within seemingly benign cover media like images, audio, or text, it offers a camouflage that's nearly undetectable for unintended observers. As a discipline, steganography has evolved dramatically, adapting to the digital age with finesse. Modern challenges in data breaches and information warfare necessitate robust steganographic techniques that can withstand sophisticated scrutiny. One such technique that's risen to prominence due to its simplicity and efficiency is the Least Significant Bit (LSB) embedding. By replacing the least significant bits of cover media with covert data, it seamlessly merges the secret with the innocuous. Renowned for its ease of implementation and impressive imperceptibility, LSB-based steganography has been recognized as an ideal choice for certain secretive communication applications (Sakshi et al., 2022).

Expanding on the realm of steganographic techniques, spatial domain methods have gained traction. Techniques like pixel intensity modifications alter pixel values strategically, thereby embedding classified data. These modifications are meticulously done, ensuring that to the naked eye, the cover media remains unchanged. The genius behind these methods is the exploitation of intrinsic properties of the cover media, leveraging nuances such as color variations or

intensity gradients. By making minuscule changes, which often escape detection, these methods can achieve a high degree of secrecy. Such spatial domain techniques offer another layer of versatility to the world of steganography, presenting a formidable challenge to those attempting unauthorized decryptions (Liao, 2022).

Branching further into the multifaceted realm of steganography, frequency domain techniques offer a sophisticated approach. Instead of operating solely on the spatial properties of the cover media, these techniques dive deep into the frequency components. Using transformations such as the Discrete Cosine Transform (DCT) or the Discrete Fourier Transform (DFT), secret data is embedded in the frequency spectrum of the cover. This approach offers a higher degree of camouflage, often eluding traditional detection methods. The strength of frequency domain methods lies in their ability to harness the intricate frequency patterns, embedding information in a way that's both secure and imperceptible (Salunkhe & Bhosale, 2022).

2.6 Least significant bit (LSB) insertion

In the multifaceted domain of steganography, the Least Significant Bit (LSB) technique stands out, especially when it comes to hiding text and images within digital media. Operating in the spatial domain, its simplicity in terms of implementation is noteworthy. The method modifies the least significant bits of a host image's pixels. This involves substituting some parts of the pixel's initial component with the secret data's most significant bits. This very simplicity is something that Liao (2022) has discussed in depth. Despite the approach's simplicity, Sakshi et al. (2022) highlight its challenges, particularly in the realm of image hiding. Here, despite the method's promise, there is an observable degradation in image quality, characterized by an elevated mean square error and a diminished peak signal-to-noise ratio.

Exploring further, the frequency domain techniques in steganography offer a complex and layered approach to data hiding. Moving beyond the spatial characteristics of the cover media, these techniques delve deep into its frequency components. They utilize transformative tools like the Discrete Cosine Transform (DCT) or the Discrete Fourier Transform (DFT) to embed secretive data into the cover media's frequency spectrum. Such methods, due to their intricacy, often dodge conventional detection mechanisms, a point of discussion in the works of Liao (2022). Furthermore, Sakshi et al. (2022) focus on the method's potential, emphasizing its unique ability to manipulate intricate frequency patterns, which makes the embedded data almost imperceptible and thus enhancing the security of the concealed information.

Usually, the least significant bits in each byte

group will often results in such minor significance compared to the overall data to the extent that modifying these bits would have absolutely minimal impact on the final result. Indeed, even changing only half of the least significant bits is significantly sufficient to discreetly embed the character 'A' (01000001) into the sequence. This demonstrates how much-hidden data can be concealed using the least significant bit substitution technique. It is a common and straightforward method for message hiding. In this technique, the message is hidden in the least significant bits of image pixels. Modifying the LSB of the pixels has minimal impact on the overall image, resulting in the stego-image closely resembling the original image. In the case of 24-bit images, three bits of each pixel can be used for LSB substitution since each pixel has separate components for red, green, and blue.

2.6.1 Masking and filtering

In steganography, masking and filtering techniques offer robust means of concealing secret messages within digital images without compromising their natural appearance (Purba et al., 2021). These methods manipulate the image's luminance values to seamlessly embed the hidden information. Specifically, the masking step delineates a specific region in the image to insert the message, whereas filtering assigns specific values to this marked section, resulting in a stego image that integrates the secret message without detection. Contrasting this with the least significant bit (LSB) technique, another spatial domain method, masking and filtering display superior resilience against various image manipulations like compression or rotation. This ensures the hidden message's security and retrievability, even if the container image undergoes alterations.

2.6.2 Parity checker method

In the evolving domain of CryptoSteganography, the Parity Checker method stands out for its fusion of cryptography and steganography to bolster the confidentiality of concealed messages (Abdelmged et al., 2016). The technique, detailed by Abdelmged (2016), utilizes a three-pronged approach: Huffman coding, the RC4 encryption algorithm, and the Parity Checker algorithm. Initially, the secret message undergoes compression via Huffman coding, which trims its size. This condensed message is then encrypted with the RC4 algorithm, imbuing it with an additional protective layer. Subsequently, this encrypted cipher text is embedded into the blue layer of a cover image using the Parity Checker algorithm. Critical to this process is the algorithm's ability to maintain the image's visual consistency, ensuring the message's covert nature. Experimental results from Abdelmged's study reveal the superiority of this method, as showcased by a higher Peak Signal-to-Noise Ratio (PSNR) and a diminished Mean Squared Error (MSE), both indicative of superior image quality and the preserved integrity of the embedded message.

2.6.3 Line shift coding

Text steganography presents unique challenges and opportunities, with line shift coding emerging as an effective technique for safeguarding embedded messages within cover texts. A fundamental aspect of this technique is ensuring the concealed information remains undisturbed, while preserving the original text's meaning. Interestingly, the Sundanese script, an official Unicode font, provides a potential medium for such covert embedding. In pivotal research by Ciptaningtyas et al. (2018), an enhanced version of line shift coding was proposed, aimed at augmenting the capacity of the concealed message. Unlike traditional methods that employ the odd row as a pivotal anchor, this revamped approach harnesses both the first and fifth rows as pivot lines. This innovative alteration not only amplifies the capacity for message storage but also fortifies the stego text's resilience against the rigors of processes like printing and copying. Impressively, even after two reprints, the method upholds the sanctity and confidentiality of the embedded messages.

2.6.4 Feature coding

Feature coding is pivotal in music genre recognition (MGR) as it encapsulates the unique nuances and intricacies of various music genres, thus enhancing indexing and retrieval processes. Most traditional representation techniques in MGR emphasize global features, often making determinations based on singular-level attributes. This strategy unfortunately glosses over the significance of granular data and the intricate dependencies present between different abstraction tiers (Ng et al., 2020), which introduces an innovative approach by harmoniously melding a convolutional neural network (CNN) with the likes of NetVLAD and self-attention mechanisms.

Weaving these tools together, the methodology is fine-tuned to capture localized information across multiple levels while also grasping the intertwined, long-term dependencies they share. NetVLAD steps in to code these local features into comprehensive representations, and the self-attention mechanism deftly models the underlying relationships amidst these features. Adding another layer of sophistication, Ng's strategy deploys a meta classifier, designed to learn and adapt from the aggregated, high-tier features sourced from various local feature coding networks. This meta entity shoulders the responsibility of making the conclusive MGR classifications. Experimental trials of this approach have been illuminating, with results showing marked improvements in accuracy over leading models on benchmark MGR datasets such as GTZAN, ISMIR2004, and Extended Ballroom. Through this fusion of feature, coding mechanisms focused on local detail and long-term interdependencies, the MGR field takes a significant leap forward.

2.6.5 Peak Signal-To-Noise ratio (PSNR)

The Peak Signal-to-Noise Ratio (PSNR) is an engineering term used to measure the ratio between the maximum power of a signal and the power of noise that can distort its fidelity. This metric is particularly valuable for assessing the quality of signal representations given their wide dynamic range. In practice, PSNR is often expressed logarithmically in decibels (dB).

Audio steganography aims to achieve capacity, robustness, and imperceptibility simultaneously, but it remains a challenge to effectively implement all three features together. The Least Significant Bit (LSB) embedding method is commonly used in audio steganography due to its high capacity and imperceptibility. However, it lacks robustness compared to other methods. To address this issue, researchers have increased the embedding depth to the fourth, sixth, and eighth LSB levels to enhance robustness. However, this trade-off between robustness and imperceptibility leads to a reduction in the imperceptibility feature, as measured by Peak Signal to Noise Ratio (PSNR) (Azam et al., 2022).

The estimation of PSNR is crucial in assessing the imperceptibility-robustness trade-off in audio steganography. However, there is a lack of studies on PSNR estimation specifically for audio steganography, making early assessment challenging. To overcome this, a PSNR Estimator (PE) method is proposed to estimate the PSNR for each stego-file generated by audio steganography. The PE method utilizes patterns extracted from the embedding process at different levels to estimate the PSNR. The proposed method achieves a high accuracy of 99.9% in estimating PSNR values at various levels. Comparative evaluation with the Mazdak Method demonstrates the superior performance of the proposed PE method in all scenarios (Azam et al., 2022).

3 Methodology

An Image Sharing Platform using Steganography refers to an online platform or service that allows users to share and distribute images while incorporating steganographic techniques for added security and privacy. Steganography involves the concealment of confidential data within the shared images, making it an effective method to protect sensitive information from unauthorized access or detection. With such a platform, users can securely transmit images containing hidden messages or encrypted data, ensuring that the concealed information remains intact and undetectable to outsiders. By leveraging steganography in image sharing platforms, users can maintain the privacy and confidentiality of their shared visual content while benefiting from the convenience and accessibility of online image sharing services.

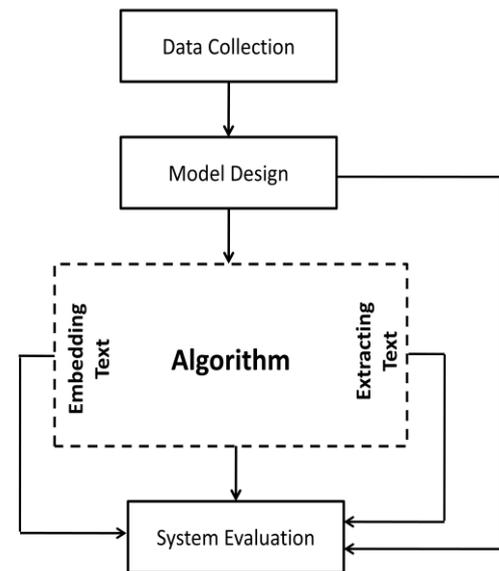


Figure 1: The proposed framework

3.1 Data collection

In order to utilize the proposed application for hiding files, the following data elements such as Cover Image, Message, and Steganographed Image (Stego-Image) were prepared.

3.2 Cover image

The image that serves as the container for the hidden message or data. The cover image can be in any form, such as JPG, PNG, AVIF.

3.3 Message

This represents the file, data, or message that was intended to be concealed within the cover image. The framework accepts text and image.

3.4 Steganographed image

This is the resulting file that is obtained after embedding the hidden data or message into the cover image. The stego-image is essentially a replica of the cover image, typically saved in PNG format.

3.5 Model design

A model is a conceptual representation of a system that simplifies and abstracts the system by disregarding certain details. Developing complementary system models can provide a holistic view of the system, including its context and interactions with other components.

Designing a model for our system involves creating an algorithmic description of the system or process. This theoretical description helps in understanding the inner workings and mechanisms of the system. By utilizing an algorithm, we can outline the step-by-step procedures and logic that govern the functioning of the system.

Creating a model through algorithmic design enables us to analyze and evaluate the system's behavior, identify potential issues or optimizations, and gain insights into its overall functionality.

4 Algorithm

Algorithms are finite sequences of well-defined instructions that are used to solve specific problems or perform computations. They provide unambiguous specifications for tasks such as calculations, data processing, and automated reasoning. Algorithms are expressed in a formal language and can be executed within a finite amount of space and time.

Algorithm 1: describes the processes of how the system embeds text in an image.

Algorithm 2 gives a description of how the system extracts the text from the image.

Algorithm 1

Embedding Text/Image Algorithm (Encoding)

- STEP 1: Initialize the required modules.
- STEP 2: Capture the message input.
- STEP 3: Capture the cover image input.
- STEP 4: Identify if the message is text or an image.
- STEP 5: If text, encode message using Base64.
- STEP 6: If text, convert the Base64-encoded message to binary.
- STEP 7: If an image, extract its RGB pixel values.
- STEP 8: Convert message/image RGB values to binary.
- STEP 9: Extract the pixel map of the cover image.
- STEP 10: Calculate the cover image's capacity.
- STEP 11: Confirm if the message fits in the cover image.
- STEP 12: Traverse the cover image pixel by pixel.
- STEP 13: Replace LSB of each pixel's RGB with the message's bit.
- STEP 14: Continue until all message bits are embedded.
- STEP 15: Convert modified binary back to image format.
- STEP 16: Save the stego-image.
- STEP 17: Provide the output to the user.

Algorithm 2

Extracting Text/ Image Algorithm

- STEP 1: Initialize necessary tools.
- STEP 2: Capture the stego-image.
- STEP 3: Extract its pixel map.
- STEP 4: Traverse each pixel in the stego-image.
- STEP 5: Extract the LSB from each RGB component.
- STEP 6: Concatenate LSBs to form the binary message.
- STEP 7: Identify the end of the message.
- STEP 8: Determine if binary represents image or text.
- STEP 9: If image, restore binary to image format.
- STEP 10: If text, convert binary to Base64.
- STEP 11: Decode Base64 to original text.
- STEP 12: Present the decoded message to the user.
- STEP 13: Conclude the decoding process.

4.1 Implementation

The implementation was carried out using the Ruby on Rails framework, leveraging the ChunkyPNG library. Ruby is an interpreted high-level general-purpose programming language. Ruby's design philosophy emphasizes code readability and productivity with its elegant syntax and focus on simplicity. Similar to Python, Ruby also uses significant indentation to enhance code clarity. Ruby on Rails is a framework for web development that is built on the Ruby programming language following the Model-View-Controller (MVC) architectural pattern, which promotes the separation of concerns and facilitates modular development. With analyses of image texture features, adaptive payload distribution strategies were devised and compared with established single-image steganographic techniques.

5 Evaluation

5.1 Testing

Upon entering the required images and text into the designated input fields and initiating the functions to embed or extract files, the program produces specific outputs based on the given inputs. Ten (10) selected Image Sets were prepared for the evaluations.

Primary cover images destined to act as vessels for concealed data and the discrete messages intended for embedding within the cover images were prepared as inputs. Similarly, the stego-images that have gone through the steganography processes were obtained, and the discrete messages were extracted from the images at the output phases.

Image Set	PSNR Value (dB)
Image Set 1	96.52
Image Set 2	102.18
Image Set 3	101.13
Image Set 4	97.39
Image Set 5	101.24
Image Set 6	101.02
Image Set 7	101.37
Image Set 8	101.19
Image Set 9	101.64
Image Set 10	99.96

Table 4.1: Evaluation of Image Sets using PSNR

5.2 Peak Signal-to-Noise Ratio (PSNR)

In the realm of image processing, assessing the quality of images is of paramount importance, especially when comparing an original image to a processed one. One of the widely accepted metrics to measure this quality is the Peak Signal-to-Noise Ratio (PSNR). PSNR is a logarithmic measure that quantifies the difference between the original and the processed images. A higher PSNR indicates better quality, as it suggests a smaller difference between the two images.

For PNG images, the formula to calculate PSNR is:

$$PSNR = 20 \times \log_{10}(\text{MAX}_I / (\sqrt{MSE}))$$

Where,

MAX_I is the maximum possible pixel value of the image. For standard PNG images,

MAX_I is 255.

MSE is the Mean Squared Error between the original and the processed image.

PSNR Values of Image Sets

The Table 4.1 presents the PSNR values for ten different image sets.

5 Discussion

However, analyzing the table, it is evident that all image sets have PSNR values well above 40dB, indicating a very high degree of similarity between the original and processed images in each set. Such high PSNR values, especially those above 100dB, suggest an

almost imperceptible difference to the human eye, which denotes outstanding performance in the image processing method employed in this project.

The PSNR value is an indicator of similarity between the two images:

- i. Above 40dB: The images are very similar.
- ii. Between 30dB to 40dB: Acceptable similarity, but there might be some noticeable differences.
- iii. Below 30dB: Significant differences exist between the images.

6.1 SSIM (Structural Similarity Index Measure)

Structural Similarity Index Measure (SSIM) is another critical metric used to assess the quality of images, particularly in the domain of steganography.

Image Set	SSIM Value
Image Set 1	0.9952
Image Set 2	0.9978
Image Set 3	0.9973
Image Set 4	0.9955
Image Set 5	0.9974
Image Set 6	0.9972
Image Set 7	0.9975
Image Set 8	0.9971
Image Set 9	0.9980
Image Set 10	0.9962

Table 4.2: Evaluation of Image Sets using SSIM

Unlike PSNR, which quantifies the difference in pixel values, SSIM evaluates the structural changes between two images, making it a more perceptually relevant metric.

The SSIM index is calculated as:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)}$$

Where:

μ_x and μ_y are the average values of x and y respectively,
 σ_x^2 and σ_y^2 are the variances of images x and y respectively,

σ_{xy} is the covariance of images x and y

C_1 and C_2 are constants to avoid division by 0

A higher SSIM value suggests that the structure, luminance, and contrast of the two images are very similar, if not identical. The resulting value lies between -1 (completely different images) and 1 (identical images).

The Table 4.2 showcases the SSIM values for ten different image sets, reflecting the high similarity between the original and stego images. Observing the table, it's evident that all SSIM values are remarkably close to 1, reinforcing the inference that the stego images are almost indistinguishable from the original ones in terms of structural and perceptual similarity. This high degree of similarity further underscores the efficacy of the steganographic techniques used in this project.

In the implementation phase, the Ruby programming language was chosen for its versatility and adaptability. The Ruby on Rails framework, known for its robustness and streamlined development capabilities, was utilized to create the backend of the platform. The ChunkyPNG library played a pivotal role in implementing steganographic techniques, allowing for the secure embedding of data within image files. Furthermore, Visual Studio Code was used as the primary editor, offering advanced code editing and debugging functionalities, which accelerated the development process.

In the testing phase, the platform was subjected to an exhaustive evaluation. Rigorous unit tests were initially conducted to assess its individual components. This was followed by comprehensive integration tests, which focused on how different parts interacted. Additionally, end-to-end tests were performed to ensure the system functioned as a whole. To evaluate its core features, test cases were specifically designed for both the steganography and encryption functionalities. During this thorough testing process, any discrepancies that arose were immediately identified. Potential vulnerabilities were also spotted and swiftly addressed. As a result of these measures, the system emerged significantly more resilient and secure.

In the evaluation phase, we extensively evaluated our steganographic techniques using two paramount metrics: PSNR and SSIM. While PSNR offered insights into pixel-level differences, SSIM sheds light on perceptual and structural similarities. The PSNR values obtained from the experiments are considered exceptionally high consistently above 90dB. Moreover, the SSIM values were nearing 1. The PSNR and SSIM

results collectively underscore the near-identical nature of our original and stego images. Interestingly, the convergence of both metrics emphasizes the effectiveness of the proposed steganographic methods, suggesting minimal distortions and high fidelity. Without mincing words, such compelling outcomes not only validate the methodology employed but also accentuate its potential for applications demanding subtle data concealment. We therefore confirm that the combined insights from PSNR and SSIM robustly affirm the project's success in achieving high-quality steganographic results.

6 Conclusion and future work

The journey of crafting an effective steganography system, underpinned by the Ruby on Rails framework coupled with the ChunkyPNG library, has culminated successfully, fulfilling its envisioned objectives. The primary intent, embedding text within cover images and extracting concealed content from stego-images, was executed with precision. Ruby on Rails, renowned for its robustness in web development, was harnessed to weave a desktop application that married efficiency with user-friendliness. Every user interaction was seamlessly facilitated through an intuitive graphical user interface (GUI), which epitomized simplistic design while not compromising on functionality. This effective synergy between design and backend processing paved the way for an enhanced user experience.

The ChunkyPNG library emerged as an invaluable asset in this endeavor. Its tailored features and modules, meticulously designed for image manipulation, enabled vital processes like RGB template splitting and binary conversion of RGB values. Additionally, its capabilities in merging templates and extracting concealed text were instrumental in breathing life into the core steganography techniques. To gauge the system's efficacy, the Peak Signal-to-Noise Ratio (PSNR) was deployed, offering an objective lens to assess image fidelity. Such a holistic assessment validated the system's prowess in ensuring data security. In essence, this project not only fortified the domain of information security with a robust tool but also sowed the seeds for future innovations, expanding the horizons of steganography research.

Steganography, as a field, is in a constant state of evolution. Each introduction of a cutting-edge steganographic method necessitates the creation of fresh applications, adapting to the innovations. This journey of continuous enhancement has led to contemporary techniques that facilitate data insertion into varied mediums, including images, documents, and audio recordings.

However, a significant constraint of the current application is its exclusive focus on image-based cover files, only allowing images to serve as the

carrier. To further elevate the application's utility and reach, it could be adapted to encompass diverse multimedia file types. As our security requirements become increasingly intricate, the scope and applications of steganography are poised for expansive growth.

References

- [1] Karabacak, A., Okay, E., & Aktas, M. S. (2022). Document Based Data Sharing Platform Architecture. *2nd International Conference on Design, Research and Development (RDCONF 2022)*, 1(1), 339–348.
- [2] Adeshina, A.M., Hashim, R. (2017) Computational Approach for Securing Radiology-Diagnostic Data in Connected Health Network using High-Performance GPU-Accelerated AES. *Interdiscip Sci Comput Life Sci* 9, 140–152. <https://doi.org/10.1007/s12539-015-0140-9>.
- [3] Liao, X., Yin, J., Chen, M., & Qin, Z. (2021). Adaptive payload distribution in multiple images steganography based on image texture features. *IEEE Transactions on Dependable and Secure Computing*, 1.
- [4] Hahnel, M. (2023). Figshare. Retrieved June 5, 2023, from <https://figshare.com/>. Online Platform for sharing research data.
- [5] Vilaza, G. N., Maharjan, R., Coyle, D., & Bardram, J. E. (2020). Futures for Health Research Data Platforms from the Participants' Perspectives. *Proceedings of the 11th Nordic Conference on Human-Computer Interaction: Shaping Experiences, Shaping Society*.
- [6] Yoon, H., Yen, C. W., Tian, F., & Zhang, Z. (2018). Healthcare data sharing in cloud computing environment. *Computers, Materials & Continua*, 56(1), 145-159
- [7] Subramanian, N., Elharrouss, O., Al-Maadeed, S., & Bouridane, A. (2021). Image Steganography: A review of the recent advances. *IEEE Access*, 9, 23409–23423.
- [8] Alqadi, M. A. (2020). Data Steganography Using Embedded Private Key. *International Journal of Engineering Technologies and Management Research*.
- [9] Abdul-Razak, N. H., Din, R., & Ahmad, M. (2018). Comparative review on feature-content based of public key steganography trends. *International Journal of Engineering & Technology*.
- [10] Majeed, M., Sulaiman, R., Shukur, Z., & Hasan, M. K. (2021). A review on text steganography techniques. *Mathematics*, 9(21), 2829.
- [11] Alishavandi, A. M., & Fakhredanesh, M. (2021). MKIPS: MKI-based protocol steganography method in SRTP. *Etri Journal*, 43(3), 561–570.
- [12] Abdulkadhim, H. A., & Shehab, J. N. (2022). Audio steganography based on least significant bits algorithm with 4D grid multi-wing hyper-chaotic system. *International Journal of Power Electronics and Drive Systems*, 12(1), 320.
- [13] Salunkhe, S., & Bhosale, S. (2022). Nature inspired algorithm for pixel location optimization in video steganography using deep RNN. *International Journal on Engineering, Science and Technology*, 3(2), 146–154.
- [14] Sakshi, S., Verma, S., Chaturvedi, P., & Yadav, S. A. (2022). Least Significant Bit Steganography for Text and Image hiding. *2022 3rd International Conference on Intelligent Engineering and Management (ICIEM)*.
- [15] Ng, W. W. Y., Zeng, W., & Wang, T. (2020). Multi-Level local feature coding fusion for music genre recognition. *IEEE Access*, 8, 152713–152727. <https://doi.org/10.1109/access.2020.3017661>
- [16] Li, S. (2021). Deep Learning-Based Watermarking: A Comprehensive Review and Future Perspectives. *IEEE Access*, 9, 30552-30571.
- [17] Shah, A., & Prakash, O. (2020). An Improved LSB-Based Watermarking Technique Using Hybrid DWT-DCT-SVD for Copyright Protection. In *2020 International Conference on Computer Communication and Informatics (ICCCI)* (pp. 1-5). IEEE.
- [18] Chugh, T., & Vashishth, S. (2020). Digital Watermarking in Medical Images: A Review. *Journal of King Saud University-Computer and Information Sciences*, 32(8), 1056-1066.
- [19] Cao, Y., Huang, J., Yang, L., & Zhang, C. (2019). A Robust and Adaptive Watermarking Scheme for 3D Point Clouds. *IEEE Access*, 7, 168100-168113.
- [20] Meng, H., & Huang, J. (2018). Blind watermarking scheme with recovery capability based on block-DCT for images. *Signal Processing*, 147, 115-123.
- [21] Casey, M. A., & Veltkamp, R. C. (Eds.). (2019). Content-based audio and image retrieval. John Wiley & Sons.
- [22] Purba, D. E. R., & Purba, D. (2021). Text Insertion by Utilizing Masking-Filtering Algorithms As Part of Text Message Security. *Jurnal Info Dan Sains: Informatika Dan Sains*, 11(1), 1–4.
- [23] Abdelmged, A. A., Saad, A. S., & Hussien, N. (2016). A Combined Approach of Steganography and Cryptography Technique based on Parity Checker and Huffman Encoding. *International Journal of Computer Applications*, 148(2), 26–32.
- [24] Ciptaningtyas, H. T., Anggoro, R., & Fadhillah, M. B. A. (2018). Text Steganography on Sundanese Script using Improved Line Shift Coding. *Text Steganography on Sundanese Script Using Improved Line Shift Coding*.
- [25] Suhail, A., & Abhayaratne, C. (2018). Image and video fingerprinting: Concepts, algorithms, and applications. *IET Image Processing*, 12(11), 1957-1973.
- [26] Azam, M. H. N., Ridzuan, F., & Sayuti, M. N. S. M. (2022). A new method to estimate peak signal to noise ratio for least significant bit modification audio steganography. *Pertanika Journal of Science and Technology*, 30(1), 497–511.

MQIBS: An Efficient Post-Quantum Identity-based Signature from Multivariate Polynomials

Le Van Luyen^{1,2}

¹Faculty of Mathematics and Computer Science, University of Science, Ho Chi Minh City, Vietnam

²Vietnam National University, Ho Chi Minh City, Vietnam

E-mail: lylvuyen@hcmus.edu.vn

Keywords: Identity-based signatures, multivariate polynomials, MQ problem

Received: September 5, 2024

Identity-based signature (IBS) is an important cryptographic primitive which allows authentication of a party's public key without the need for certificates. In this paper, we construct a post-quantum provable identity-based signature scheme from multivariate polynomials. Our scheme is constructed from the sigma protocols with helper by Beullens at Eurocrypt 2020 and the Fiat-Shamir paradigm. Concrete choice of parameters shows that our scheme is more efficient than existing multivariate IBS schemes in terms of public key/signature sizes.

Povzetek: Predstavljena je postkvantna identitetna podpisna shema MQIBS, zasnovano na multivariatnih polinomih. Z uporabo sigma protokolov in Fiat-Shamirjevega pristopa izboljšuje varnost in učinkovitost, potrjuje manjše velikosti javnih ključev in podpisov, kar prispeva k večji praktičnosti v postkvantni kriptografiji.

1 Introduction

Post-quantum Cryptography (PQC) has become an emerging research direction since the announcement of NIST (National Institute of Standards and Technology) for the PQC standardization process since 2016 [1]. NIST selected several candidates for standardization in 2022 and called for additional digital signatures whose the first round deadline was June 2023 [2]. Among the candidates for PQC, multivariate cryptography is one of the main candidates for this standardization [1, 2]. It is shown in [11, 12] that multivariate schemes are very fast in general and suitable for limited computational resources, such as smart cards that run RFID chips. The security of multivariate schemes is normally based on the hardness of the *MQ-Problem*, which is proven to be NP-Hard \mathbb{F}_2 [3], that asks for a solution of a given system of multivariate quadratic polynomials over the field \mathbb{F}_q .

Multivariate cryptography is dated back to the early work of Matsumoto and Imai in 1988 [6], and since then, there has been a rich development of designing multivariate schemes in several directions. Notably in the history is the (Unbalanced) Oil and Vinegar (UOV) signature scheme by Patarin et al. after he broke the Matsumoto-Imai scheme [7, 8]. Since then, there has been the main direction on improving the UOV schemes, including the multi-layer variant Rainbow by Ding and Schmidt [10], and its cyclic version [14]. Rainbow was one of the main candidates in the NIST PQC Process until Round 3, but it was not selected due to the attack by Beullens [30] which re-

duced the proposed security levels, i.e., in order to achieve the required level of security, Rainbow needs to update the parameters which will result in large key and signature sizes. Since then, the intention focuses back to improving the UOV scheme. Especially there have been many such submissions in the round 1 of NIST Additional PQC Signatures [2] including for example MAYO, PROV, QR-UOV and TUOV; see [2] for the details.

One drawback of UOV signatures is that they do not have a provable security proof. Recent submissions like MAYO or PROV do have such a proof but it was reduced to a new assumption, not the NP-complete MQ problem as expected. For a provable secure construction, another direction of construction multivariate signatures is to follow the Fiat-Shamir paradigm [5]. In this case, one needs first an identification scheme and the Fiat-Shamir transformation converts it into a secure digital signature. The first multivariate identification schemes were proposed by Sakumoto et al. [17]. They include a 3-pass and 5-pass identification schemes. The 5-pass identification was used to design the MQDSS signature [20, 19] which was a candidate for NIST PQC Round 2. However, it was broken by Kales and Zaverucha [27]. All work in this direction is hence focusing on improving one from 3-pass identification scheme by Sakumoto et al. [17], including [25]. Recently, Beullens [26, 26] developed sigma protocols with helper, inspiring from the work by Katz et al. [21], and applied to the 3-pass identification scheme by Sakumoto et al. [17] to obtain a more efficient multivariate digital signature compared to MQDSS [20].

Identity-Based Signature (IBS), proposed by Shamir [4], allows for the generation of a public key for an entity using only some basic scheme parameters and an identifier string (such as an email address or phone number). A private key generator (PKG) derives private keys from a master secret and distributes them to the entities involved in the scheme. This approach removes the requirement for certificates, unlike in traditional public key infrastructure. There have been many post-quantum constructions of IBS. In the area of multivariate cryptography, there have been several proposals for IBS based on UOV such as [18, 24] or Rainbow such as [23, 24]. Recently, there has been a proposal for identity-based signature by Debnath et al. [31].

In this paper, we investigate the sigma protocol with helper by Beullens and design an identity-based signature scheme from multivariate polynomials, which we call MQIBS. Our MQIBS scheme enjoys the security reduction to the underlying MQ problem and is more efficient than existing schemes; see Table 2 for the details.

Related work There are basically two approaches to construct an IBS. One is called the certification approach [13], transforms a standard signature scheme into an IBS scheme, from which this paper follows. The other one [9] is to transform a 2-level hierarchical identity-based encryption (HIBE) scheme to an IBS scheme. For post-quantum identity-based signatures, there exist several constructions. The most dominant candidates come from lattice-based constructions ([15, 28, 32]) which follow the second approach, since there exist trapdoors in lattices which enable efficient HIBE constructions ([16]). The remaining post-quantum identity-based signature candidates follow the first approach. Isogeny-based [33, 29] and group actions-based [34] constructions follow a variant [22] of the first approach to achieve schemes with tight reduction; however, due to the less flexibility of group actions, the constructions require a lot of “layers” which may result in in-efficient schemes compared to the efficient digital counterparts. In multivariate cryptography, there have been several constructions based on Rainbow such as [23, 24, 31]. In this paper, we propose a new one which is more efficient than the aforementioned schemes.

The rest of the paper is organized as the following. In Section 2, we recall some basic notions on commitment schemes and identity-based signatures including their definition and security model. We recall the sigma protocols with helper from Beullens [30] in Section 3 and our construction of MQIBS is presented in Section 4. In Section 5, we provide the choice of parameters and compute the key and signature size of MQIBS as well as a comparison between MQIBS with existing multivariate IBS. Section 6 concludes the paper.

2 Preliminaries

2.1 Commitment schemes

A commitment scheme $\text{Com} : \{0, 1\}^\lambda \times \{0, 1\}^* \rightarrow \{0, 1\}^{2\lambda}$, where λ is the security parameter, is a function that takes as input λ uniformly random bit $r \in \{0, 1\}^\lambda$ and a message $m \in \{0, 1\}^*$, outputs a 2λ bit long commitment $\text{Com}(r, m)$. We require the following two properties of a commitment scheme ([26]).

Definition 1 (Computational binding). *For an adversary \mathcal{A} we define its advantage for the commitment binding game as*

$$\text{Adv}_{\text{com}}^{\text{Binding}}(\mathcal{A}) = \Pr[\text{Com}(r, m) = \text{Com}(r', m') | (r, m, r', m') \leftarrow \mathcal{A}(1^\lambda)].$$

We say that Com is computationally binding if for all polynomial-time algorithms \mathcal{A} , the advantage $\text{Adv}_{\text{com}}^{\text{Binding}}(\mathcal{A})$ is a negligible function of the security parameter λ .

Definition 2 (Computational hiding). *For an adversary \mathcal{A} we define the advantage for the commitment hiding game for a pair of messages m, m' as*

$$\text{Adv}_{\text{com}}^{\text{Hiding}}(\mathcal{A}, m, m') = \left| \Pr_{r \leftarrow \{0, 1\}^\lambda} [1 = \mathcal{A}(\text{Com}(r, m))] - \Pr_{r \leftarrow \{0, 1\}^\lambda} [1 = \mathcal{A}(\text{Com}(r, m'))] \right|.$$

We say that Com is computationally hiding if for all polynomial-time algorithms \mathcal{A} , and every pair of messages m, m' the advantage $\text{Adv}_{\text{com}}^{\text{Hiding}}(\mathcal{A}, m, m')$ is a negligible function of the security parameter λ .

2.2 Identity-based signature scheme

An identity-based signature (IBS) scheme is a tuple of polynomial-time algorithms $\text{IBS} = (\text{Setup}, \text{KeyDer}, \text{Sign}, \text{Verify})$ as follows:

Setup(1^λ): On input the security parameter λ , it outputs the master public key and secret key pair (mpk, msk).

KeyDer(msk, id): On input the master secret key msk and a user identity ID , it generates the user secret key usk.

Sign(mpk, usk, M): On input the master public key mpk, the user secret key usk and a message M , it outputs a signature σ .

Verify(mpk, id, σ, M): On input the master public key mpk, user identity ID , a signature-message pair (σ, M) , it outputs 1 for acceptance and 0 for rejection.

We consider the following properties for an IBS scheme. First, the correctness guarantees that a signature generated by an honest signer will always pass the verification algorithm.

Definition 3 (Correctness). *We say that an identity-based signature IBS is correct, if for all $\lambda \in \mathbb{N}$, all identity $\text{id} \in \mathcal{ID}$ and all message $M \in \mathcal{M}$ that if $(\text{mpk}, \text{msk}) \leftarrow \text{Setup}(1^\lambda)$, $\text{usk}_{\text{id}} \leftarrow \text{KeyDer}(\text{mpk}, \text{msk}, \text{id})$, $\sigma \leftarrow \text{Sign}(\text{mpk}, \text{usk}_{\text{id}}, M)$ then it holds that*

$$\Pr[\text{Verify}(\text{mpk}, \text{id}, \sigma, M) = 1] = 1 - \text{negl}(\lambda).$$

Second, it is required that an adversary cannot create a new tuple (id, message, signature) for an identity and a message that it hasn't been queried before, given that it has already seen some identities' secret keys and signatures for some tuples (identity, message) of its choice.

Definition 4 (EUF-ID-CMA). *We say that an identity-based signature IBS is EUF-ID-CMA if, for every PPT adversary \mathcal{A} , it holds that \mathcal{A} has a negligible advantage in the following experiment.*

$$\text{Exp}_{\mathcal{A}}^{\text{EUF-ID-CMA}}(\lambda) :$$

1. The challenger generates $(\text{mpk}, \text{msk}) \leftarrow \text{Setup}(1^\lambda)$ and sets $\mathcal{Q}_{\text{id}} \leftarrow \emptyset$, $\hat{\mathcal{Q}}_{\text{id}} \leftarrow \emptyset$, $\mathcal{Q}_{\text{usk}_{\text{id}}} \leftarrow \emptyset$ and $\mathcal{Q}_M \leftarrow \emptyset$.
2. The challenger gives mpk to the adversary \mathcal{A} . Moreover, \mathcal{A} can access two signing oracles $\mathcal{O}_{\text{KeyDer}}$, $\mathcal{O}_{\text{Sign}}$, where
 - i. Key derivation oracle $\mathcal{O}_{\text{KeyDer}}$: On input a key derivation query $\text{id} \in \mathcal{ID}$, the oracle $\mathcal{O}_{\text{KeyDer}}$ checks whether $(\text{id}, \cdot) \in \mathcal{Q}_{\text{id}}$. If $(\text{id}, \cdot) \in \mathcal{Q}_{\text{id}}$ for some $\text{usk}_{\text{id}} \in \mathcal{USK}$, it returns usk_{id} . Otherwise, it returns $\text{usk}_{\text{id}} \leftarrow \text{KeyDer}(\text{mpk}, \text{msk}, \text{id})$ and sets $\mathcal{Q}_{\text{id}} \leftarrow \mathcal{Q}_{\text{id}} \cup \{(\text{id}, \text{usk}_{\text{id}})\}$.
 - ii. Signing oracle $\mathcal{O}_{\text{Sign}}$: On input a signing query $(\text{id}, M) \in \mathcal{ID} \times \mathcal{M}$, the oracle $\mathcal{O}_{\text{Sign}}$ sets $\mathcal{Q}_M \leftarrow \mathcal{Q}_M \cup \{(\text{id}, M)\}$ and checks whether $(\text{id}, \cdot) \in \mathcal{Q}_{\text{id}}$
 - If $(\text{id}, \cdot) \in \mathcal{Q}_{\text{id}}$ for some $\text{usk}_{\text{id}} \in \mathcal{USK}$, returns $\sigma \leftarrow \text{Sign}(\text{mpk}, \text{usk}_{\text{id}}, M)$.
 - If there does not exist $(\text{id}, \cdot) \in \mathcal{Q}_{\text{id}}$ for any $\text{usk}_{\text{id}} \in \mathcal{USK}$, it computes $\text{usk}_{\text{id}} \leftarrow \text{KeyDer}(\text{mpk}, \text{msk}, \text{id})$, return $\sigma \leftarrow \text{Sign}(\text{mpk}, \text{usk}_{\text{id}}, M)$ and sets $\mathcal{Q}_{\text{id}} \leftarrow \mathcal{Q}_{\text{id}} \cup \{(\text{id}, \text{usk}_{\text{id}})\}$.
3. In the end, the adversary outputs a forgery $(\text{id}^*, M^*, \sigma^*)$.
4. The challenger outputs 1 if the following three conditions are hold:
 - There does not exist $(\text{id}^*, \cdot) \in \mathcal{Q}_{\text{id}}$ for any $\text{usk}_{\text{id}^*} \in \mathcal{USK}$,

- $(\text{id}^*, M^*) \notin \mathcal{Q}_M$,
- $\text{Verify}(\text{mpk}, \text{id}^*, M^*, \sigma^*) = 1$

The advantage of \mathcal{A} is defined by $\text{Adv}_{\mathcal{A}}^{\text{EUF-ID-CMA}}(\lambda) = \Pr[\text{Exp}_{\mathcal{A}}^{\text{EUF-ID-CMA}}(\lambda) = 1]$.

3 Sigma protocols with helper

Beullens [26] introduced a sigma protocol with a helper, which involves a three-round sigma protocol that includes a trusted third party, known as the helper. At the start of each protocol execution, the helper runs a setup algorithm using a random seed. The helper then sends auxiliary information to the verifier and provides the seed value used in the setup algorithm to the prover. The syntax can be summarized in the Figure 1.

Definition 5 (Sigma protocol with helper [26]). *A protocol is a sigma protocol with helper for relation R with challenge space \mathcal{C} if it is of the form of Fig. 1 and satisfies:*

Completeness. *If all parties (Helper, Prover, Verifier) follow the protocol on input $(x, w) \in R$, then the verifier always accepts.*

2-Special soundness. *From an adversary \mathcal{A} that outputs with noticeable probability two valid transcripts $(x, \text{aux}, \text{com}, \text{ch}, \text{rsp})$ and $(x, \text{aux}, \text{com}, \text{ch}', \text{rsp}')$ with $\text{ch} \neq \text{ch}'$ and where $\text{aux} = \text{Setup}(\text{seed})$ for some seed value seed (not necessarily known to the extractor) one can efficiently extract a witness w such that $(x, w) \in R$.*

Special honest-verifier zero-knowledge. *There exists a PPT simulator S that on input x , a random seed value seed and a random challenge ch outputs a transcript $(x, \text{aux}, \text{com}, \text{ch}, \text{rsp})$ with $\text{aux} = \text{Setup}(\text{seed})$ that is computationally indistinguishable from the probability distribution of transcripts of honest executions of the protocol on input (x, w) for some w such that $(x, w) \in R$, conditioned on the auxiliary information being equal to aux and the challenge being equal to ch .*

Beullens then transformed sigma protocols with helper in Figure 1 into a standard zero-knowledge proof of knowledge without helper using the ‘‘Cut-and-choose’’ approach by Katz et al. [21]. We recall it in Figure 2.

Theorem 1 ([26, Theorem 3]). *Let $(\text{Setup}, P_1, P_2, V)$ be a sigma protocol with helper and challenge space \mathcal{C} , if the used commitment scheme is hiding, then the protocol of Fig. 2 is an honest-verifier zero-knowledge proof of knowledge with challenge space $\{1, \dots, k\} \times \mathcal{C}$ and $\max(k, |\mathcal{C}|) + 1$ -special soundness (and hence it has soundness error $\max(\frac{1}{k}, \frac{1}{|\mathcal{C}|})$).*

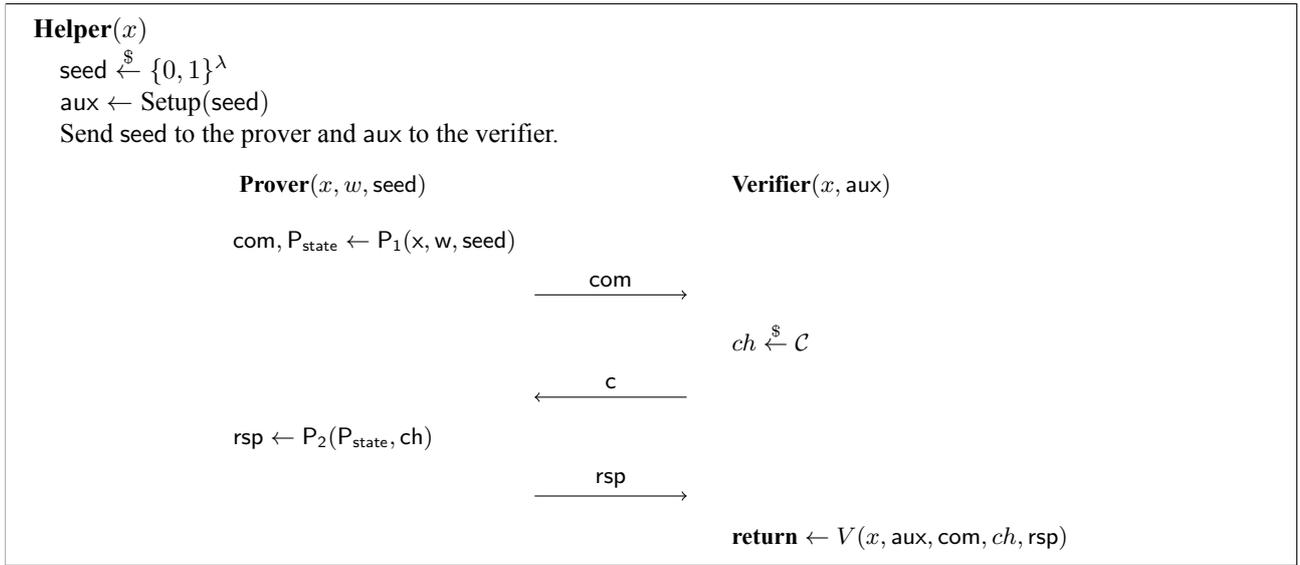


Figure 1: The structure of a sigma protocol with trusted setup

4 Construction of MQIBS

4.1 Sigma protocol with helper for MQ problem

In this section, we recall the sigma protocol with helper for MQ Problem from [26] in proving knowledge of a solution of a system of multivariate quadratic equations over a finite field \mathbb{F}_q . This scheme improves the previous two schemes by Sakumoto et al. . In particular, the two schemes by Sakumoto et al. have soundness errors $\frac{2}{3}$ and $\frac{1}{2} + \frac{1}{2q}$ respectively while the one with helper has soundness error to only $\frac{1}{q}$.

Let $\mathcal{F} : \mathbb{F}_q^n \rightarrow \mathbb{F}_q^m$ is a multivariate quadratic map of m polynomials in n variables, define the polar form of \mathcal{F} as

$$\mathcal{G}(\mathbf{x}, \mathbf{y}) := \mathcal{F}(\mathbf{x} + \mathbf{y}) - \mathcal{F}(\mathbf{x}) - \mathcal{F}(\mathbf{y}).$$

Note that \mathcal{G} is linear in both \mathbf{x} and \mathbf{y} . The sigma protocol is described in Figure 3.

Theorem 2 ([26, Theorem 1]). *Suppose the used commitment scheme is computationally binding and computationally hiding, then the protocol of Fig. 3 is a sigma protocol with trusted setup as in Definition 5 with challenge space \mathbb{F}_q .*

4.2 Our identity-based signature construction

In this Section, we propose a construction of an identity-based signature from the sigma protocol with helper presented in Figure 3, which we call MQIBS. The idea is to follow the construction by Kiltz et al. [13]. The MQIBS

scheme consists of the following polynomial-time algorithms.

Setup(1^λ): Given security parameter λ , output public parameters consisting of m, n, q, k , a random system of polynomials $\mathcal{F} : \mathbb{F}_q^n \rightarrow \mathbb{F}_q^m$, and a hash function $H : \{0, 1\}^* \rightarrow \{1, \dots, k\} \times \mathbb{F}_q$, and do the following:

- Choose $\mathbf{s} \xleftarrow{\$} \mathbb{F}_q$ and compute $\mathbf{v} := \mathcal{F}(\mathbf{s}) \in \mathbb{F}_q^m$.
- Output mpk = $(\mathcal{F}, \mathbf{v})$ and msk = \mathbf{s} .

KeyDer(msk, id): Given the master secret key msk = \mathbf{s} and a user identity id, do the following:

- Choose a random system $\mathcal{F}_{\text{id}} : \mathbb{F}_q^n \rightarrow \mathbb{F}_q^m$, $\mathbf{s}_{\text{id}} \xleftarrow{\$} \mathbb{F}_q^n$ and compute $\mathbf{v}_{\text{id}} = \mathcal{F}_{\text{id}}(\mathbf{s}_{\text{id}})$.
- For $i \in \{1, \dots, k\}$ do
 - seed _{i} $\xleftarrow{\$}$ $\{0, 1\}^\lambda$
 - Compute aux _{i} as in the procedure of Helper(\mathcal{F}) in Figure 3.
 - Compute com _{i} as in the first step of the Prover as in Figure 3.
- Set COM_{id} := $(\text{com}_i, \text{aux}_i)_{i \in \{1, \dots, k\}}$.
- Compute $(I, \alpha) := H(\text{COM}_{\text{id}}, \mathcal{F}_{\text{id}} \parallel \text{id})$
- Retrieve seed _{I} and compute the response rsp (using \mathbf{s}) as in the response by the Prover as in Figure 3. Set RSP_{id} = $(\text{rsp}, \text{seed}_i \forall i \neq I)$.
- Output the user secret key as usk = $(\mathbf{s}_{\text{id}}, \mathbf{v}_i, \mathcal{F}_{\text{id}}, \text{COM}_{\text{id}}, \text{RSP}_{\text{id}})$.

Sign(mpk, usk, M): Given the master public key mpk, the user secret key usk of a user id and a message M , do the following:

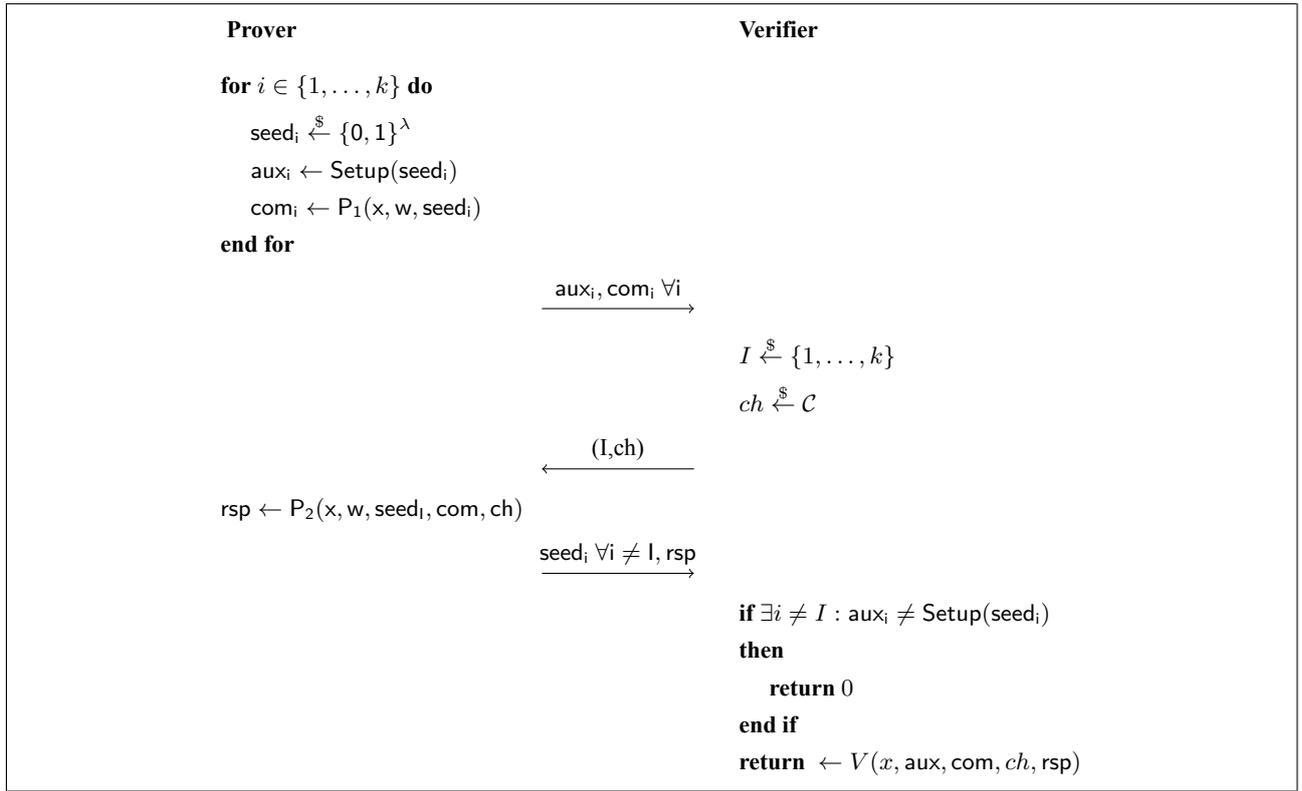


Figure 2: Zero-knowledge proof without helper

- Parse usk as $(s_{id}, v_{id}, \mathcal{F}_{id}, \text{COM}_{id}, \text{RSP}_{id})$.
- For $i \in \{1, \dots, k\}$ do
 - $\text{seed}_i \xleftarrow{\$} \{0, 1\}^\lambda$
 - Compute aux_i as in the procedure of Helper(\mathcal{F}_{id}) in Figure 3.
 - Compute com_i as in the first step of the Prover as in Figure 3.
- Set $\mathbf{COM} := (\text{com}_i, \text{aux}_i)_{i \in \{1, \dots, k\}}$.
- Compute $(J, c) := H(\mathbf{COM}, M)$
- Retrieve seed_J and compute the response rsp (using s_{id}) as in the response by the Prover as in Figure 3. Set $\text{RSP} = (\text{rsp}, \text{seed}_i \forall i \neq I)$.
- Output a signature as $\sigma = (v_{id}, \mathcal{F}_{id}, \text{COM}_{id}, \text{RSP}_{id}, \text{COM}, \text{RSP})$.

Verify(mpk, id, σ , M): Given a master public key mpk, an identity id, a signature-message pair σ and M , do the following:

- Parse the signature σ as $\sigma = (v_{id}, \mathcal{F}_{id}, \text{COM}_{id}, \text{RSP}_{id}, \text{COM}, \text{RSP})$.
- Use $(v_{id}, \mathcal{F}_{id}, \text{COM}, \text{RSP})$ and M do the verification as by the Verifier in Figure 3. Let the result of the verification be b_1 .
- Use mpk, $\mathcal{F}_{id} \parallel \text{id}$, COM_{id} , RSP_{id} as an input for the verification procedure as in Figure 3. Call the result of this process to be b_2 .

Correctness The correctness is straight-forward with noting that (COM, RSP) is a signature for the message M under the public key v_{id} , \mathcal{F}_{id} and secret key s_{id} of the user id, and COM_{id} , RSP_{id} is a signature for the message $\mathcal{F}_{id} \parallel \text{id}$ under the master public key mpk and master secret key msk.

Security The security is also straightforward following the result by Kitlitz et al. [13]. In fact, the MQIBS is EUF-ID-CMA secure provided that the underlying signature is EUF-CMA secure. Note that the underlying signature is obtained from applying the Fiat-Shamir transformation to the sigma protocol in Figure 3. Hence if there is an adversary that breaks the EUF-CMA of the underlying signature, then, it is folklore [5] that, there exists an adversary that breaks the soundness of the sigma protocol in Figure 3. It follows from [26, Section 5] that we can construct an algorithm to solve the underlying MQ problem.

[htb]

Helper(\mathcal{F})seed $\xleftarrow{\$}$ $\{0, 1\}^\lambda$ Generate $\mathbf{e} \in \mathbb{F}_q^m$ and $\mathbf{t}, \mathbf{r}_0 \in \mathbb{F}_q^n$ from seed.**for each** $c \in \mathbb{F}_q$ **do** $\mathbf{e}_c \leftarrow c\mathcal{F}(\mathbf{r}_0) - \mathbf{e}$ $\mathbf{t}_c \leftarrow c\mathbf{r}_0 - \mathbf{t}$ Generate commitment randomness $\mathbf{r}_c \in \{0, 1\}^\lambda$ from seed. $\text{com}_c \leftarrow \text{Com}(\mathbf{r}_c, (\mathbf{e}_c, \mathbf{t}_c))$ **end for**aux $\leftarrow [\text{com}_c | c \in \mathbb{F}_q]$

Send seed to the prover and aux to the verifier.

Prover($\mathcal{F}, \mathbf{s}, \text{seed}$)Regenerate $\mathbf{e}, \mathbf{t}, \mathbf{r}_0$ from seed $\mathbf{r} \leftarrow \{0, 1\}^\lambda$ $\text{com} \leftarrow \text{Com}(\mathbf{r}, (\mathbf{r}_1, \mathbf{e} + \mathcal{G}(\mathbf{r}_1, \mathbf{t})))$ $\xrightarrow{\text{com}}$ **Verifier**($\mathcal{F}, \mathbf{v}, \text{aux}$) $\alpha \xleftarrow{\$} \mathbb{F}_q$ $\xleftarrow{\alpha}$ Recompute $\mathbf{r}_\alpha, \mathbf{e}_\alpha, \mathbf{t}_\alpha$ from seed $\xrightarrow{(\mathbf{r}, \mathbf{r}_\alpha, \mathbf{r}_1, \mathbf{e}_\alpha, \mathbf{t}_\alpha)}$ $\mathbf{x} \leftarrow \alpha(\mathbf{v} - \mathcal{F}(\mathbf{r}_1)) - \mathbf{e}_\alpha - \mathcal{G}(\mathbf{r}_1, \mathbf{t}_\alpha)$ $b_1 \leftarrow \text{com} = \text{Com}(\mathbf{r}, (\mathbf{r}_1, \mathbf{x}))$ $b_2 \leftarrow \text{com}_\alpha = \text{Com}(\mathbf{r}, (\mathbf{e}_\alpha, \mathbf{t}_\alpha))$ **Return** $b_1 \wedge b_2$

Figure 3: Sigma protocol with helper for MQ problem

5 Parameters

5.1 Optimizations

In this section, we review about the techniques presented in [26], which was followed by Katz et al. [21], for optimizations. Some are as follows; see [26, Section 7] for the details.

- Build a Merkle tree on the commitments com_i computed in the KeyDer process as well as the Sign process to reduce the user secret key usk as well as signature size σ . In particular, instead of including all com_i in COM_{id} (resp. COM), we use only the root of the Merkle tree created by the com_i , and hence RSP_{id} (resp. RSP) consists only $\lceil \log_2(q) \rceil$ nodes required to reconstruct the root of the Merkle tree. Similarly, we can do the same for aux_i in COM_{id} (resp. COM), and seed_i in RSP_{id} (resp. RSP).
- For some applications, the finite field \mathbb{F}_q is large and

hence not practical to compute Merkle trees of size q . We then can reduce the challenge space to $\mathbb{F}_{q'}$ with $q' < q$. It then makes the protocol in Figure 3 has soundness error $\frac{1}{q'}$ (instead of $\frac{1}{q}$). Katz et al. [21] suggested that, instead of letting the verifier choose 1 out of k setups to execute, we now let him choose τ out of M setups to execute, which results to the soundness error bounded by

$$\max_{0 \leq e \leq \tau} \frac{\binom{M-e}{\tau-e}}{\binom{M}{\tau} q'^{\tau-e}}.$$

See [26, Section 7] or [21] for the details.

5.2 Parameters

We follow [26] for the choice of parameters m, n, q, τ, M . First of all, $q = 4$ and $m = n = 88, 128, 160$ respectively following [20] in the MQDSS submission to the NIST PQC standardization project. Note that we choose the number of equations m equal to the number of variables n , i.e., $m = n$,

Table 1: Parameters for MQIBS

Security Level	q	n	M	τ	mpk (B)	msk (B)	Signature (B)
I	4	88	191	68	38	16	28,838
III	4	128	256	111	56	24	65,856
V	4	160	380	136	72	32	111,272

Table 2: Comparison between our MQIBS with other existing multivariate IBS

Scheme	mpk (B)	msk (B)	Signature (B)
Ours (MQIBS)	38	16	28,838
IBS-Rainbow [23]	148,300	103,700	304,300
ID-Rainbow [24]	4,220,000	142,600	46
Mul-IBS-Rainbow [31]	136,100	90,900	33,400

to ensure that we get the hardest instance of an MQ problem [20]. The parameters τ and M are chosen to balance between the signature size and running time: increasing τ impacts signature size, while decreasing M impacts signing and verification time [26]. The parameters for MQIBS, key and signature sizes are summarized in Table 1.

In Table 2, we provide a comparison between existing IBS scheme from multivariate polynomials at the security level I. It can be seen from Table that our scheme outperforms the existing ones (except the ID-Rainbow [24]) in terms of signature size. In addition, our scheme has the smallest public key size. Compared to the ID-Rainbow by Chen et al. [24], our scheme has much bigger signature size (28,838 B vs. 46 B) but has much smaller master public key (38 B vs. 4,220,000B) and secret key (16 B vs. 142,600 B). We note that due to the recent attack by Beullens [30] against Rainbow scheme, those existing schemes mentioned in Table 2 would be vulnerable to Beullens' attack and hence need to update the parameters to attain the same security level which will result in much larger key/signature sizes compared to those mentioned in Table 2.

6 Conclusion

In this paper, we propose a design for an identity-based signature, called MQIBS, from multivariate polynomials. Our scheme is derived from the sigma protocol with helper for MQ from Beullens [26] from which our MQIBS inherits the efficiency. Compared to existing schemes, our scheme has advantage on both signature size as well as public key and secret key size. Further optimizations and implementations are one of the goals for our future work.

Acknowledgement

The author would like to thank Dung Duong for useful discussions.

Funding

This research is funded by Vietnam National University, Ho Chi Minh City (VNU-HCM) under grant number C2022-18-03.

References

- [1] National Institute of Standards and Technology post-quantum cryptography. <https://csrc.nist.gov/projects/post-quantum-cryptography>. Accessed: 2024-07-24.
- [2] National Institute of Standards and Technology additional post-quantum signatures. <https://csrc.nist.gov/projects/pqc-dig-sig/round-1-additional-signatures>. Accessed: 2024-07-24.
- [3] Michael R. Garey and David S. Johnson (1979). *Computers and Intractability: A Guide to the Theory of Np-Completeness*. W. H. Freeman. <https://doi.org/10.2307/2273574>
- [4] Adi Shamir (1984). Identity-based cryptosystems and signature schemes. In G. R. Blakley and David Chaum, editors, *Advances in Cryptology, Proceedings of CRYPTO '84, Santa Barbara, California, USA, August 19-22, 1984, Proceedings*, volume 196 of *Lecture Notes in Computer Science*, pages 47–53. Springer. https://doi.org/10.1007/3-540-39568-7_5
- [5] Amos Fiat and Adi Shamir (1986). How to prove yourself: Practical solutions to identification and signature problems. In Andrew M. Odlyzko, editor, *Advances in Cryptology - CRYPTO '86, Santa Barbara, California, USA, 1986, Proceedings*, volume 263 of *Lecture Notes in Computer Science*, pages 186–194.

- Springer.
https://doi.org/10.1007/3-540-47721-7_12
- [6] Tsutomu Matsumoto and Hideki Imai (1988). Public Quadratic Polynomial-Tuples for Efficient Signature-Verification and Message-Encryption. In: Barstow, D., et al. *Advances in Cryptology — EUROCRYPT '88. EUROCRYPT 1988. Lecture Notes in Computer Science*, vol 330. Springer, Berlin, Heidelberg.
https://doi.org/10.1007/3-540-45961-8_39
- [7] Jacques Patarin (1995). Cryptanalysis of the Matsumoto and Imai Public Key Scheme of Eurocrypt'88. In: Coppersmith, D. (eds) *Advances in Cryptology — CRYPTO'95. CRYPTO 1995. Lecture Notes in Computer Science*, vol 963. Springer, Berlin, Heidelberg.
https://doi.org/10.1007/3-540-44750-4_20
- [8] Kipnis, A., Patarin, J., Goubin, L. (1999). Unbalanced Oil and Vinegar Signature Schemes. In: Stern, J. (eds) *Advances in Cryptology — EUROCRYPT '99. EUROCRYPT 1999. Lecture Notes in Computer Science*, vol 1592. Springer, Berlin, Heidelberg.
https://doi.org/10.1007/3-540-48910-X_15
- [9] Craig Gentry, Alice Silverberg (2002). Hierarchical ID-Based Cryptography. *ASIACRYPT 2002*: 548–566.
https://doi.org/10.1007/3-540-36178-2_34
- [10] Jintai Ding and Dieter Schmidt (2005). Rainbow, a new multivariable polynomial signature scheme. In John Ioannidis, Angelos D. Keromytis, and Moti Yung, editors, *Applied Cryptography and Network Security, Third International Conference, ACNS 2005, New York, NY, USA, June 7-10, 2005, Proceedings*, volume 3531 of *Lecture Notes in Computer Science*, pages 164–175.
https://doi.org/10.1007/11496137_12
- [11] Andrey Bogdanov, Thomas Eisenbarth, Andy Rupp, and Christopher Wolf (2008). Time-area optimized public-key engines: Mq-cryptosystems as replacement for elliptic curves? *IACR Cryptol. ePrint Arch.*, page 349.
https://doi.org/10.1007/978-3-540-85053-3_4
- [12] Anna Inn-Tung Chen, Ming-Shing Chen, Tien-Ren Chen, Chen-Mou Cheng, Jintai Ding, Eric Li-Hsiang Kuo, Frost Yu-Shuang Lee, and Bo-Yin Yang (2009). SSE implementation of multivariate pkcs on modern x86 cpus. In Christophe Clavier and Kris Gaj, editors, *Cryptographic Hardware and Embedded Systems - CHES 2009, 11th International Workshop, Lausanne, Switzerland, September 6-9, 2009, Proceedings*, volume 5747 of *Lecture Notes in Computer Science*, pages 33–48. Springer.
https://doi.org/10.1007/978-3-642-04138-9_3
- [13] Eike Kiltz and Gregory Neven (2009). Identity-based signatures. In Marc Joye and Gregory Neven, editors, *Identity-Based Cryptography*, volume 2 of *Cryptology and Information Security Series*, pages 31–44. IOS Press.
<https://doi.org/10.3233/978-1-58603-947-9-31>
- [14] Albrecht Petzoldt, Stanislav Bulygin, and Johannes Buchmann (2010). Cyclicrainbow - A multivariate signature scheme with a partially cyclic public key. In Guang Gong and Kishan Chand Gupta, editors, *Progress in Cryptology - INDOCRYPT 2010 - 11th International Conference on Cryptology in India, Hyderabad, India, December 12-15, 2010. Proceedings*, volume 6498 of *Lecture Notes in Computer Science*, pages 33–48. Springer.
https://doi.org/10.1007/978-3-642-17401-8_4
- [15] Markus Rückert (2010). Strongly Unforgeable Signatures and Hierarchical Identity-Based Signatures from Lattices without Random Oracles. *PQCrypto 2010*: 182–200.
https://doi.org/10.1007/978-3-642-12929-2_14
- [16] Shweta Agrawal, Dan Boneh, Xavier Boyen (2010). Efficient Lattice (H)IBE in the Standard Model. *EUROCRYPT 2010*: 553–572
https://doi.org/10.1007/978-3-642-13190-5_28
- [17] Koichi Sakumoto, Taizo Shirai, and Harunaga Hiwatari (2011). Public-key identification schemes based on multivariate quadratic polynomials. In Phillip Rogaway, editor, *Advances in Cryptology - CRYPTO 2011 - 31st Annual Cryptology Conference, Santa Barbara, CA, USA, August 14-18, 2011. Proceedings*, volume 6841 of *Lecture Notes in Computer Science*, pages 706–723. Springer.
https://doi.org/10.1007/978-3-642-22792-9_40
- [18] Wuqiang Shen, Shaohua Tang, and Lingling Xu (2013). Ibuov, A provably secure identity-based UOV signature scheme. In *16th IEEE International Conference on Computational Science and Engineering, CSE 2013, December 3-5, 2013, Sydney, Australia*, pages 388–395. IEEE Computer Society.
<https://doi.org/10.1109/CSE.2013.66>
- [19] Ming-Shing Chen, Andreas Hülsing, Joost Rijneveld, Simona Samardjiska, and Peter Schwabe (2016). From 5-pass MQ -based identification to MQ -based signatures. In Jung Hee Cheon and Tsuyoshi Takagi, editors, *Advances in Cryptology - ASIACRYPT 2016 - 22nd International Conference on the Theory and Application of Cryptology and Information Security*,

- Hanoi, Vietnam, December 4–8, 2016, *Proceedings, Part II*, volume 10032 of *Lecture Notes in Computer Science*, pages 135–165.
https://doi.org/10.1007/978-3-662-53890-6_5
- [20] Ming-Shing Chen, Andreas H Ising, Joost Rijneveld, Simona Samardjiska, and Peter Schwabe (2017). MQDSS-Submission to the NIST post-quantum cryptography project. In *NIST Post-quantum Cryptography*. available at <https://csrc.nist.gov/CSRC/media/Projects/Post-Quantum-Cryptography/documents/round-1/submissions/MQDSS.zip>
- [21] Jonathan Katz, Vladimir Kolesnikov, and Xiao Wang (2017). Improved Non-Interactive Zero Knowledge with Applications to Post-Quantum Signatures. In *CCS '18: Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pages 525 – 537. ACM.
<https://doi.org/10.1145/3243734.3243805>
- [22] Masayuki Fukumitsu, Shingo Hasegawa (2018). A Galindo-Garcia-Like Identity-Based Signature with Tight Security Reduction, Revisited. *CANDAR 2018*: 92–98
<https://doi.org/10.1109/CANDAR.2017.79>.
- [23] Le Van Luyen (2019). An improved identity-based multivariate signature scheme based on rainbow. *Cryptography*, 3(1):8.
<https://doi.org/10.3390/cryptography3010008>
- [24] Jiahui Chen, Jie Ling, Jianting Ning, and Jintai Ding (2019). Identity-based signature schemes for multivariate public key cryptosystems. *Comput. J.*, 62(8):1132–1147.
<https://doi.org/10.1093/comjnl/bxz013>
- [25] Hiroki Furue, Dung Hoang Duong, and Tsuyoshi Takagi (2019). An efficient mq-based signature in the QROM. In *2019 Seventh International Symposium on Computing and Networking, CANDAR 2019, Nagasaki, Japan, November 25–28, 2019*, pages 10–17. IEEE.
<https://doi.org/10.1109/CANDAR.2019.00010>
- [26] Ward Beullens (2020). Sigma Protocols for MQ, PKP and SIS, and Fishy Signature Schemes. In Anne Canteaut and Yuval Ishai, editors, *Advances in Cryptology - EUROCRYPT 2020 - 39th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Zagreb, Croatia, May 10–14, 2020, Proceedings, Part III*, volume 12107 of *Lecture Notes in Computer Science*, pages 183–211. Springer.
https://doi.org/10.1007/978-3-030-45727-3_7
- [27] Daniel Kales and Greg Zaverucha (2020). An attack on some signature schemes constructed from five-pass identification schemes. In Stephan Krenn, Haya Schulmann, and Serge Vaudenay, editors, *Cryptology and Network Security - 19th International Conference, CANS 2020, Vienna, Austria, December 14–16, 2020, Proceedings*, volume 12579 of *Lecture Notes in Computer Science*, pages 3–22. Springer.
https://doi.org/10.1007/978-3-030-65411-5_1
- [28] Jiaxin Pan, Benedikt Wagner (2021). Short Identity-Based Signatures with Tight Security from Lattices. *PQCrypto 2021*: 360–379
https://doi.org/10.1007/978-3-030-81293-5_19
- [29] Surbhi Shaw, Ratna Dutta (2021). Identification Scheme and Forward-Secure Signature in Identity-Based Setting from Isogenies. *ProvSec 2021*: 309–326.
https://doi.org/10.1007/978-3-030-90402-9_17
- [30] Ward Beullens (2022). Breaking rainbow takes a weekend on a laptop. In Yevgeniy Dodis and Thomas Shrimpton, editors, *Advances in Cryptology - CRYPTO 2022 - 42nd Annual International Cryptology Conference, CRYPTO 2022, Santa Barbara, CA, USA, August 15–18, 2022, Proceedings, Part II*, volume 13508 of *Lecture Notes in Computer Science*, pages 464–479. Springer.
https://doi.org/10.1007/978-3-031-15979-4_16
- [31] Sumit Kumar Debnath, Sihem Mesnager, Vikas Srivastava, Saibal Kumar Pal, and Nibedita Kundu (2023). Mul-ibs: a multivariate identity-based signature scheme compatible with iot-based NDN architecture. *J. Cryptogr. Eng.*, 13(2):187–199.
<https://doi.org/10.1007/s13389-022-00308-8>
- [32] Ernest Foo, Qinyi Li (2023). Tightly Secure Lattice Identity-Based Signature in the Quantum Random Oracle Model. *ACISP 2023*: 381–402.
https://doi.org/10.1007/978-3-031-35486-1_17
- [33] Jiawei Chen, Hyungrok Jo, Shingo Sato, Junji Shikata (2023). A Tightly Secure Identity-Based Signature Scheme from Isogenies. *PQCrypto 2023*: 141–163.
https://doi.org/10.1007/978-3-031-40003-2_6
- [34] Xuan Thanh Khuc, Willy Susilo, Dung Hoang Duong, Fuchun Guo, Hyungrok Jo, Tsuyoshi Takagi (2024).

Tightly Secure Identity-based Signature from Cryptographic Group Actions. *ProSec 2024*.

Improved Genetic Algorithm Enhanced with Generative Adversarial Networks for Logistics Distribution Path Optimization

Juan Li

Henan Institute of Economics and Trade Zhengzhou 450046, China

E-mail: lijuan_vip@hotmail.com

Keywords: generative adversarial network, genetic algorithm, high fitness solution, logistics distribution path,

Received: August 21, 2024

This paper proposes an innovative logistics distribution path planning algorithm, which aims to combine the generative adversarial network (GAN) with the genetic algorithm (GA) to solve the path optimization problem in large-scale distribution networks. The GA-GAN algorithm intelligently improves the mutation operation of the genetic algorithm through GAN, which not only outperforms the traditional genetic algorithm and other classic heuristic algorithms in terms of solution quality, operation efficiency, convergence speed and solution stability, but also provides quantitative data of specific improvements. Experimental results show that when GA-GAN processes a data set of 500 customer points, the average running time is 160 seconds, the optimal solution cost is 9500 units, the average solution cost is 10500 units, and it can reach the optimal solution within 180 iterations, which is significantly better than the baseline genetic algorithm (average running time is 150 seconds, the optimal solution cost is 10000 units, the average solution cost is 12000 units, and the average number of iterations required to reach the optimal solution is 300 times). In addition, GA-GAN has good responsiveness to the size of the data set and has a wide range of adaptability to different distribution scenarios, providing an efficient, stable and flexible distribution path planning solution for the logistics industry.

Povzetek: Razvit je optimiziran genetski algoritem, izboljšan z generativnimi adversarialnimi omrežji (GA-GAN), za načrtovanje logističnih poti. Eksperimentalni rezultati kažejo, da GA-GAN doseže optimalno rešitev v 180 iteracijah, kar je bistveno hitreje kot standardni genetski algoritem. Sistem izboljšuje načrtovanje distribucije v dinamičnih okoljih, zmanjšuje stroške in povečuje prilagodljivost logističnih operacij.

1 Introduction

With the rapid development of globalization and e-commerce, the logistics industry is facing unprecedented challenges and opportunities. Logistics not only needs to meet the growing demand for commodity transportation, but also seeks a balance between cost control, efficiency improvement and environmental protection. In this context, route optimization has become a key link in logistics management and its importance is self-evident. Effective path optimization can not only significantly reduce logistics costs and improve distribution efficiency, but also reduce carbon emissions, which has a far-reaching impact on promoting sustainable development [1].

The rapid development of the logistics industry has given rise to the demand for efficient and intelligent logistics systems. As shown in Figure 1, the global logistics scale continues to grow. In 2023 alone, the global e-commerce logistics market size reached trillions of dollars and is expected to continue to grow in the coming years. However, path planning problems in the logistics and distribution process, such as multi-objective distribution route optimization, vehicle

scheduling, and time window constraints, seriously affect logistics efficiency and customer satisfaction. Therefore, it becomes crucial to find a method that can effectively solve these complex path optimization problems [2].

Genetic algorithm, as a global optimization search algorithm that mimics the natural evolutionary process, shows a strong potential in dealing with such NP-hard problems. It is able to search efficiently in the solution space and find a solution close to the optimal one by simulating biological evolution mechanisms such as natural selection and genetic variation. Nevertheless, the standard genetic algorithm still has limitations such as slow convergence speed and easy to fall into local optimality when dealing with high-dimensional and complex problems, which limits the effectiveness of its application in logistics and distribution path optimization.

Traditional path optimization techniques usually include two categories: exact algorithms and heuristic algorithms. Among them, exact algorithms, such as Dijkstra's algorithm and Floyd's algorithm, are mainly used to solve the shortest path problem in deterministic networks, and they have the advantage of being able to

find the global optimal solution [3], but with higher computational complexity, which is suitable for smaller-scale problems. Heuristic algorithms such as greedy algorithm, simulated annealing algorithm and ant colony algorithm, on the other hand, are more suitable for solving large-scale and dynamically changing path optimization problems, sacrificing the accuracy of the solution in exchange for the computational efficiency, and are suitable for the scenarios with high demand for real-time and dynamic adjustments [4]. As a global optimization search technique based on the principles of natural selection

and genetics, genetic algorithm has shown its unique charm in logistics and distribution path optimization in recent years. It simulates the genetic evolution process in nature through operations such as coding, selection, crossover and mutation, so as to efficiently explore the solution space and find optimal or suboptimal solutions. A key advantage of genetic algorithms is their ability to handle nonlinear, multi-peak and multi-objective optimization problems, which makes them irreplaceable in solving complex logistics and distribution path problems [5].

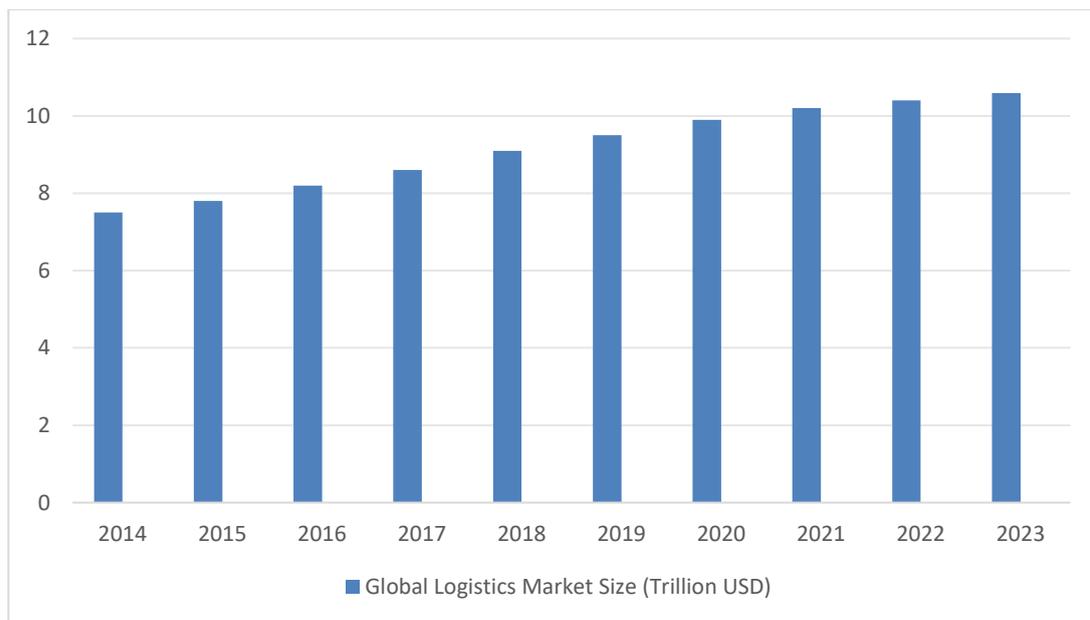


Figure 1: Scale of global logistics in the last 10 years

The literature review shows that genetic algorithms present diverse characteristics in the application of logistics and distribution path optimization. For example, the literature proposed a hybrid algorithm combining a genetic algorithm and a local search strategy for solving the vehicle routing problem with time windows (VRPTW), and the experiments showed that the algorithm dramatically improved the solution speed while ensuring the solution quality [6]. Literature, on the other hand, focuses on multi-objective logistics and distribution path optimization, and they developed a multi-objective genetic algorithm based on Pareto optimization, which successfully achieves a balance between multiple objectives, such as cost, time and carbon emission [6]. Although genetic algorithms have achieved remarkable results in logistics and distribution path optimization, they still face some unresolved challenges. First, genetic algorithms are prone to fall into local optimality, especially in high-dimensional complex problems, how to design effective selection, crossover and mutation operators to avoid premature convergence is one of the hotspots in current research. Secondly, multi-objective optimization problems are common in logistics and distribution, how to find the Pareto optimal solution among multiple conflicting

objectives is another urgent problem to be solved. Finally, the dynamic characteristics of logistics and distribution network require the algorithm to have the ability of fast response and online adjustment, which requires the algorithm not only to perform well in the static environment, but also need to have a certain degree of dynamic adaptability [7].

In view of the above background and literature review, this study aims to develop a logistics and distribution path optimization method based on an improved genetic algorithm to overcome the limitations of existing algorithms and improve the overall performance of path optimization. Specifically, the objectives of the study are (1) to design and implement a novel genetic algorithm, which enhances the global search capability of the algorithm and avoids premature convergence by introducing new genetic operators and strategies. (2) Construct a multi-objective optimization model to achieve the comprehensive optimization of logistics and distribution paths by simultaneously considering factors such as distribution cost, time, and carbon emission. (3) Verify the effectiveness and superiority of the improved genetic algorithm in solving the optimization problem of logistics and distribution paths through comparative experiments, and provide

scientific basis for the decision-making of logistics enterprises.

2 Rationale and related work

Before discussing in depth, the optimization of logistics and distribution paths based on improved genetic algorithm, it is necessary to review the basic principles of genetic algorithm, understand the theoretical framework of logistics and distribution path optimization, and analyze the advantages and limitations of the existing research, with a view to laying a solid theoretical foundation for the subsequent innovations.

2.1 Principles of genetic algorithm

Genetic Algorithm (GA) is a global optimization search technique that mimics the process of biological evolution. The core idea of GA is to solve the optimization problem by using natural selection and genetic mechanism. The algorithmic flow of GA is shown in Figure 2.

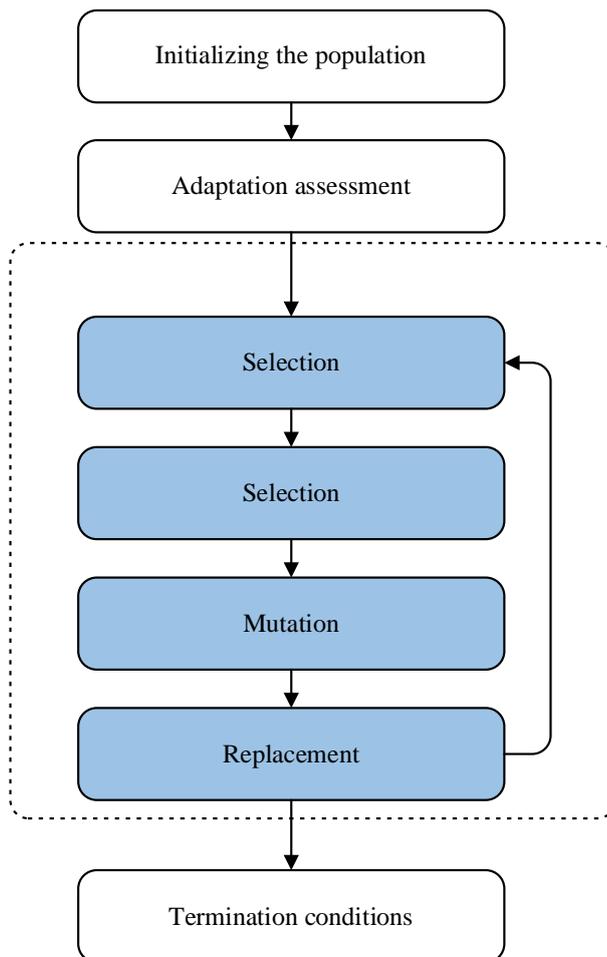


Figure 2: Algorithm flow of GA

In GA, potential solutions are encoded as chromosomes (chromosome), which are composed of genes (gene), each representing a variable in the solution. The fitness function (fitness function) is used

to evaluate the quality of the chromosome, i.e., the degree of merit of the solution. Genetic algorithms follow a systematic workflow in solving complex optimization problems, which is designed to mimic the principles of natural selection and genetics in order to find optimal solutions [8]. First, the algorithm constructs a starting population by randomly generating a set of initial solutions that represent possible candidate answers to the problem. Next, each individual in the population is evaluated for fitness, a process that quantifies the quality of each solution, usually measured in terms of the objective function value of the problem. Moving on to the selection phase, the algorithm picks out those individuals that perform better, based on fitness values, as parents for the next generation. This phase ensures that high-quality solutions have a better chance of passing on their "genes" to their offspring, thus gradually improving the overall performance of the population. The selected individuals are paired through a crossover operation, a process similar to mating between organisms, in which some of the genes of the two individuals are exchanged to create a new individual that combines the characteristics of both individuals [9]. In order to maintain the diversity of the population and avoid premature convergence, the algorithm also imposes a mutation operation, i.e., randomly modifying the genes of certain individuals with a low probability, which corresponds to the introduction of a random perturbation that helps to explore undiscovered regions in the solution space. After the generation of new offspring, a replacement step occurs, where new individuals replace some members of the old population, a process that is usually carried out based on some kind of elimination strategy, such as retaining only the most adapted individuals. The cycle repeats until a predetermined termination condition is reached, either a fixed number of iterations is reached or the average fitness of the population no longer improves significantly, indicating that the algorithm is close to an optimal solution. Through this series of well-designed steps, the genetic algorithm is able to search efficiently in the solution space and eventually approximate or even find an optimal or near-optimal solution to the problem [10].

2.2 Logistics distribution path optimization theory

The logistics distribution path optimization problem, as an extension of the Traveler's Problem (TSP), is a central challenge in the field of logistics. While the TSP requires finding a shortest path that visits all cities exactly once and returns to the starting point, the problem becomes more complex in logistics distribution, as it needs to take into account factors such as the loading capacity of the distribution vehicles, the customer's time window, and the possible simultaneous operation of multiple vehicles. Logistics distribution path optimization is usually modeled as a mixed integer linear programming (MILP) problem. For example, the literature proposes the classical TSP model, which uses binary variables to indicate whether edges on a path are

selected or not. However, in logistics and distribution scenarios, this model needs to be extended to include additional constraints, such as the capacity limit of vehicles and the customer's demand, which can be achieved by introducing more decision variables and constraints. The objective function is usually the minimization of the total cost, which can include transportation costs, fixed costs (e.g., vehicle rental fees), and variable costs (e.g., fuel consumption and driver wages [11]). Key considerations, besides cost, are time management and distance. Time-window constraints require that the delivery be completed within a specific time period, which increases the complexity of the problem. Meanwhile, distance not only affects the transportation cost, but also determines the driver's working hours and possible overtime costs. Therefore, logistics and distribution path optimization need to consider these factors comprehensively in order to find a solution that satisfies all constraints and achieves the optimal objective [12].

2.3 Analysis of related work

Over the past few decades, genetic algorithms (GA) have become an effective tool for solving logistics and distribution path optimization problems due to their powerful search capability and adaptability. GA is able to efficiently deal with multi-objective optimization problems by mimicking the process of genetic evolution in nature and using selection, crossover, and mutation operations to search the solution space. GA is able to deal with complex multi-objective optimization problems due to GA's is able to search multiple directions in the solution space simultaneously, rather than just moving along the gradient direction. By introducing variation and crossover, GA is able to maintain the diversity of the population and avoid premature convergence to local optimal solutions. In addition, the parallel nature of GA makes it exhibit high efficiency in dealing with large-scale problems, allowing it to explore different regions of the solution space quickly [13]. Despite the above advantages of GA, it has some inherent limitations. Standard GA may become slow to converge when dealing with high dimensional data and complex constraints. Premature convergence is another common problem, i.e., the algorithm may converge prematurely without finding a globally optimal solution. In addition, GA lacks an efficient method to balance exploration (exploration) and utilization (exploitation), which may lead to less efficient search, especially when the solution space is very large [14]. To overcome these limitations, many researchers have proposed different improvement strategies. For example, [15] explored how to improve the performance of GA through adaptive parameter tuning. [16] proposed a strategy of combining local search methods with GA to enhance the local search capability of the algorithm. [17], on the other hand, proposed a GA variant based on elite retention, which helps to preserve the best individuals in the population and prevents the loss of high-quality genes, thus

speeding up the convergence rate.

The development of our GA-GAN algorithm has been influenced by several key works in the field of genetic algorithms and their applications. Vaira G. and Kurasava O. [18] presented a genetic algorithm for the Vehicle Routing Problem (VRP) with constraints based on feasible insertion, which underscored the importance of customizing genetic operations for specific problem domains. Additionally, Misevičius A., Kuznecovaitė D., and Platužienė J. [19] conducted comprehensive experiments with crossover operators, emphasizing the significance of selecting suitable genetic mechanisms to improve the performance of genetic algorithms. Furthermore, Gams M. and Kolenik T. [20] explored the interrelations between electronics, artificial intelligence, and the information society, providing context for the integration of advanced AI techniques such as GANs into traditional optimization methods. These studies collectively informed the design and implementation of our hybrid GA-GAN approach, particularly in optimizing mutation operations and enhancing overall algorithmic efficiency.

In order to more fully understand and evaluate the performance of different algorithms in logistics distribution path optimization, we compiled a summary table (Table 1) to compare the key features, methods, and technical achievements of various algorithms. Standard genetic algorithms (GAs) are known for their powerful search capabilities and adaptability. Through selection, crossover, and mutation operations, they can efficiently handle multi-objective optimization problems and maintain population diversity. However, standard GAs may have slow convergence when dealing with high-dimensional data and complex constraints. To this end, researchers proposed adaptive parameter adjustment (APT) GAs, which significantly improved the convergence speed of the algorithm by dynamically adjusting the selection, crossover, and mutation rates. In addition, hybrid GAs combined with local search methods enhance the algorithm's search ability in promising areas and further improve the quality of solutions. The introduction of an elite retention mechanism ensures the preservation of high-quality individuals and accelerates convergence. For specific problem areas, such as the vehicle routing problem (VRP), feasible insertion GAs ensure that the solution meets specific constraints through customized genetic operations, improving the quality and feasibility of the solution. Through comprehensive experiments, the researchers also evaluated the effects of different crossover operators and identified the most effective genetic mechanism. Finally, our GA-GAN hybrid approach optimizes the mutation operation by integrating a generative adversarial network (GAN), which not only improves the overall efficiency of the algorithm but also significantly improves the quality of the solution. Together, these improved strategies form the basis of our advanced algorithm for logistics distribution path optimization.

3 Design and implementation of improved genetic algorithm

In the design of the GA-GAN algorithm, new individuals are created by introducing GAN in the mutation phase of the genetic algorithm. The specific steps are as follows:

Table1: Comparison of key features, methods, and technical achievements of different algorithms

Algorithm/Approach	Key Features	Methods	Technical Achievements
Standard Genetic Algorithm (GA)	Powerful search capability, adaptability	Selection, crossover, mutation	Efficiently deals with multi-objective optimization, maintains population diversity
Adaptive Parameter Tuning (APT) GA	Improved convergence speed, dynamic parameter adjustment	Adaptive selection, crossover, and mutation rates	Overcomes slow convergence in high-dimensional data and complex constraints
Local Search Hybrid GA	Enhanced local search capability	Combination of GA and local search methods	Improves the quality of solutions by refining the search in promising areas
Elite Retention GA	Preservation of best individuals	Elite retention mechanism	Prevents loss of high-quality genes, speeds up convergence rate
Feasible Insertion GA for VRP	Customized genetic operations for VRP	Feasible insertion, specialized crossover and mutation	Ensures solutions meet problem-specific constraints, improves solution quality
Crossover Operator Experiments	Comprehensive evaluation of crossover operators	Various crossover operators	Identifies the most effective genetic mechanisms for specific problems
GA-GAN Hybrid Approach	Integration of GANs for mutation operations	Hybrid GA and GAN framework	Enhances mutation operations, improves overall algorithmic efficiency and solution quality

GAN generates new individuals: extracts random noise from the Gaussian distribution and generates new individuals through the trained generator G.

Evaluate new individuals: Use the fitness function to evaluate the newly generated individuals and select the best performing individuals to join the population

Selection, crossover, mutation: After obtaining a new population, the population is further optimized through operations such as selection, crossover, and mutation. Through this process, the GA-GAN algorithm not only enhances the diversity of the population, but also increases the proportion of excellent individuals in the population through individuals generated by GAN.

3.1 Logistics distribution path modeling

In the field of logistics and distribution, distribution path planning is a typical combinatorial optimization problem, which can be regarded as an extension of the Traveling Salesman Problem (TSP). The objective of distribution path planning is to determine the set of shortest paths from the distribution center to all demand points and back, subject to a set of constraints (e.g., time window, vehicle capacity limitations), in order to minimize the distribution cost (e.g., total travel distance, time, or fuel consumption). This problem can be mathematically modeled as a Mixed Integer Linear Programming (MILP) problem, but its computational complexity grows exponentially as the number of distribution points increases, making it difficult to find an exact solution [18]. x_{ij} represents a binary variable, $x_{ij} = 1$ if the

delivery vehicle travels directly from point i to point j ; otherwise, $x_{ij} = 0$. u_i represents an integer variable representing the position of point i in the distribution sequence, which is used to prevent the formation of sub-loops (sub-tours). Our goal is to minimize the total distance traveled by the distribution vehicles. Let c_{ij} denote the distance between point i and j . The objective function can be expressed as Equation 1.

$$\text{Minimize } Z = \sum_{i=1}^N \sum_{j=1}^N c_{ij} x_{ij} \quad (1)$$

where N is the total number of distribution points, including distribution centers. The access constraint is that each customer point must be and can only be accessed once, which is guaranteed by the following two constraints, Equation 2 and Equation 3 [19].

$$\sum_{j=1, j \neq i}^N x_{ij} = 1 \quad \forall i \in V \quad (2)$$

$$\sum_{j=1, j \neq i}^N x_{ji} = 1 \quad \forall i \in V \quad (3)$$

Here V denotes the set of all points and V is the set of all points except the distribution center. The

subloop avoidance constraint matter is to ensure that the distribution paths are coherent and there are no independent subloops that do not pass through the distribution center, we use the auxiliary variable u_i to introduce the following constraints Equation 4 and Equation 5 [20].

$$u_j - u_i + Nx_{ij} \leq N - 1 \quad \forall i, j \in V, i \neq j \quad (4)$$

$$1 \leq u_i \leq N - 1 \quad \forall i \in V \quad (5)$$

3.2 Algorithm improvement points

The application of standard Genetic Algorithm (GA) for solving logistics and distribution path planning problems is a common strategy due to its ability to handle large-scale problems and find near-optimal solutions. However, standard genetic algorithms may encounter local optimality traps and premature convergence problems, especially when the solution space is very large. To address these problems, we can consider enhancing the performance of genetic algorithms by incorporating Generative Adversarial Network (GAN), specifically, the mutation operation in genetic algorithms can be improved by GAN [21].

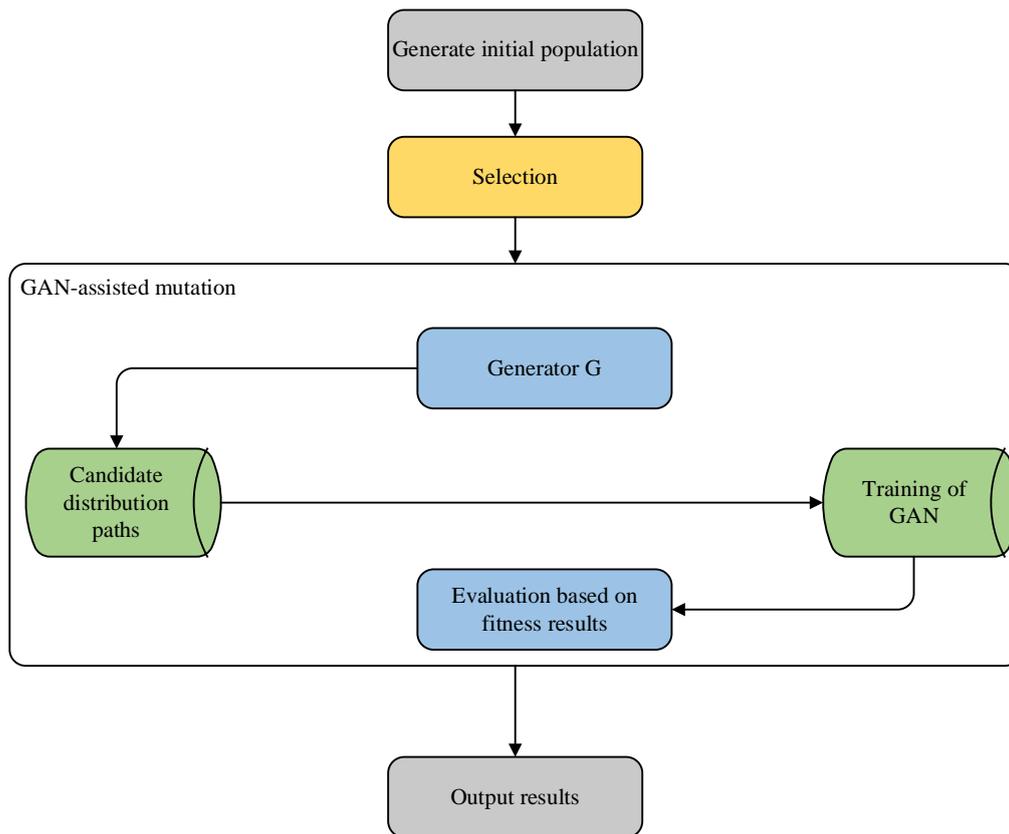


Figure 3: Algorithmic framework

The framework of our algorithm is shown in Figure 3. The GA-GAN algorithm forms an efficient hybrid algorithm by combining the search capability of genetic algorithms with the data generation advantage of generative adversarial networks, which first initializes a population in which each individual is a potential solution, and then enters the main loop to iteratively

optimize the population through selection, crossover, and GAN-assisted mutation operations, while training the GAN to generate more diverse and high-quality data, evaluating and replacing individuals in the population until preset stopping conditions are met, such as reaching the iteration limit or fitness threshold, and ultimately outputting the optimal individuals, this process allows

GA-GAN to perform well in solving problems that require generating new data samples or optimizing complex objective functions, and can be flexibly adapted and optimized according to the specific situation.

The introduction of Generative Adversarial Networks (GANs) to refine the mutation operation in the framework of genetic algorithms is a cutting-edge strategy to improve the efficiency and quality of the algorithms. Although traditional mutation can maintain population diversity, it may lead to the generation of a large number of invalid solutions due to randomness, slowing down the optimization search process. In contrast, by training the GAN, the generator G learns to mold individuals close to the high fitness solution from the noise, and the discriminator D is responsible for screening the authenticity to ensure the good quality of the solution. Once G is trained, it can be used in the mutation phase of the genetic algorithm to create new individuals that are more likely to be high-quality solutions.

In the logistics distribution path planning problem, we can train the generator G of the GAN as a function that accepts a random vector \mathbf{z} as input and outputs a sequence of potential distribution paths \mathbf{x} . The path sequence \mathbf{x} can be viewed as an ordered set of nodes, where each node represents a delivery point. The discriminator D then evaluates the truthfulness and fitness of the path sequence \mathbf{x} . Let the distribution path consist of n distribution points, the generator can be designed as a sequence generation model that outputs the index of the next distribution point in each iteration until the whole sequence is generated. The output of the generator can be a probability distribution of the distribution points, which is then sampled to determine the next distribution point. The task of the discriminator is to distinguish between a real sequence of highly adaptive distribution paths and the sequence generated by the generator. It takes as input a sequence and outputs a real number between 0 and 1 indicating the probability that the sequence is "real" (i.e., highly adaptive). The discriminator can also be designed as a neural network whose input is a sequence of distribution paths and whose output is a scalar indicating the truthfulness of the sequence. The discriminator may need to encode features of the distribution path, such as total distance traveled, time window satisfaction, etc., in order to more accurately determine the quality of the sequence. The goal of the generator is to learn how to generate highly adaptable distribution path sequences from random noise \mathbf{z} . The generator can be designed as a Multilayer Perceptron (MLP), Recurrent Neural Network (RNN), Long Short-Term Memory Network (LSTM), or Graph Neural Network (GNN), depending on the dependencies between the distribution points and the dynamic nature of the paths [22].

The training of the GAN follows the classical min-max game, in which the generator tries to deceive the discriminator, while the discriminator tries to correctly distinguish the real sequence from the generated

sequence. The objective function for training is as in Equation 6.

$$\min_G \max_D V(D, G) = E_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + E_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \quad (6)$$

where $p_{data}(\mathbf{x})$ is the distribution of high fitness distribution paths and is the random noise distribution of the generator input. In the mutation step of the genetic algorithm, for an individual \mathbf{x}_i , we first sample a random vector \mathbf{z} from $p_z(\mathbf{z})$, and then generate a new sequence \mathbf{x}'_i through the generator as G in Equation 7 [23].

$$\mathbf{x}'_i = G(\mathbf{z}), \quad \text{where } \mathbf{z} \sim p_z(\mathbf{z}) \quad (7)$$

The new sequence \mathbf{x}'_i is considered as a mutated individual, which is subsequently added to the population to participate in subsequent genetic operations such as crossover and selection. The training objective of the GAN is to make the generator G maximize the error of the discriminator D , and at the same time, make the discriminator D maximize its own ability to differentiate between real and generated data. The objective function can be formulated as Equation 8 [24].

$$\min_G \max_D V(D, G) = E_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + E_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \quad (8)$$

In the genetic algorithm, for an individual \mathbf{x}_i , we sample \mathbf{z} from $p_z(\mathbf{z})$ and generate a new individual \mathbf{x}'_i by G . This process can be mathematized as Equation 9.

$$\mathbf{x}'_i = G(\mathbf{z}), \quad \text{where } \mathbf{z} \sim p_z(\mathbf{z}) \quad (9)$$

Combining GAN with genetic algorithm can significantly improve the performance of genetic algorithm on logistics and distribution path planning problems. By training the GAN to guide the mutation operation, it not only accelerates the convergence of the algorithm, but also improves the quality of the solution and avoids premature convergence, thus providing an effective solution to complex optimization problems. This approach is particularly suitable for large-scale problems, in which traditional methods may not be able to solve efficiently [25].

The fitness function plays a crucial role in integrating GAN-generated solutions with traditional GA. Both the offspring generated by the GA and the mutations generated by the GAN are evaluated using the same fitness function. This ensures a fair comparison and selection process. The fitness function typically measures how well each individual (or solution) meets the optimization criteria, such as minimizing the total cost in logistics distribution path optimization. The best individuals, whether generated by the GA or the GAN, are selected to form the next generation, ensuring that the population evolves towards an optimal solution.

3.3 Algorithm flow description

The process of initializing a genetic algorithm population involves creating a set of initial path solutions. These solutions can be generated randomly, ensuring that different parts of the solution space are covered. Each path solution is represented as an ordered list of nodes, where the nodes include the distribution center and all demand points. The initialization process can be expressed as Equation 10 [26].

$$P_0 = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \quad (10)$$

where \mathbf{x}_i denotes the path sequence of the i th individual and P_0 is the initial population. The selection operation is based on the fitness value of the individual, which is usually inversely proportional to the total cost of the path. Strategies such as roulette selection or tournament selection are used to select well-performing individuals from the current population into the next generation. The selection process can be described as Equation 11.

$$P_{sel} = \{\mathbf{x}_{sel_1}, \mathbf{x}_{sel_2}, \dots, \mathbf{x}_{sel_n}\} \quad (11)$$

where P_{sel} is the population after selection and \mathbf{x}_{sel_i} is the i th selected individual. The crossover operation exchanges some of the genes between two selected individuals to generate new offspring. The two-point crossover used can be defined as Equation 12 [27].

$$\mathbf{x}_{child} = \text{Crossover}(\mathbf{x}_{parent1}, \mathbf{x}_{parent2}) \quad (12)$$

Here, $\mathbf{x}_{parent1}$ and $\mathbf{x}_{parent2}$ are the selected parent individuals, and \mathbf{x}_{child} is the new individual created by crossover. The mutation operation introduces new solutions by fine-tuning some of the genes of an individual. In this scenario, the mutation operation is intelligently guided by GAN to generate a new individual \mathbf{x}'_i from the noise \mathbf{z} . The mutation process can be described as Equation 13 [28].

$$\mathbf{x}'_i = G(\mathbf{z}), \text{ where } \mathbf{z} \sim p_z(\mathbf{z}) \quad (13)$$

where G is the trained generator, \mathbf{z} is the random noise, and \mathbf{x}'_i is the mutated new individual. The fitness function calculates the fitness value of an individual, which usually corresponds to the cost of the path, i.e., the total distance or time traveled. The fitness function can be expressed as Equation 14.

$$f(\mathbf{x}) = \sum_{i=1}^N \sum_{j=1}^N c_{ij} x_{ij} \quad (14)$$

Where, c_{ij} is the distance from node i to node j , \mathbf{x} is the path solution and $f(\mathbf{x})$ is the total cost of that path.

Integrating GANs into the mutation process of a genetic algorithm adds additional computational cost. The training process of a GAN, which involves alternating training of the generator and the discriminator, requires forward and backward propagation through the network for each training iteration, which is computationally intensive. In addition, data preparation before GAN training, such as data normalization and formatting, also incurs additional overhead. Once the GAN is trained, generating mutations involves passing the offspring through the generator network, which also takes time. Moreover, evaluating the fitness of the mutations generated by the GAN also adds computational cost. Therefore, integrating GANs into a genetic algorithm increases the running time of the entire algorithm. In a standard genetic algorithm, the running time is mainly determined by fitness evaluation and genetic operations

(selection, crossover, mutation), and the addition of GANs undoubtedly increases the complexity of this process.

3.4 Realization details

In the logistics distribution path planning scheme explored in this paper, we employ a novel combination - the fusion of Generative Adversarial Networks (GANs) and Genetic Algorithms - with the aim of breaking through the limitations of traditional algorithms and improving distribution efficiency. Our technology stack is centered on Python, complemented by a series of efficient tools, including the PyTorch deep learning framework, the DEAP genetic algorithm framework, and optimization solvers such as CPLEX or GUROBI, which together form a powerful problem-solving platform. Leveraging the flexibility of the DEAP framework, we carefully tuned the components of the genetic algorithm, including the population size, crossover probability, mutation probability, and selection strategy, to ensure the efficient operation of the algorithm and the diversity of solutions. The population size was set in a moderate range (50-100), while the crossover probability and mutation probability were maintained at high (0.8-0.9) and low (0.01-0.05) levels, respectively, to balance exploration and exploitation [29].

This comprehensive strategy not only improves the efficiency and accuracy of distribution path planning, but also provides a strong technical support for the intelligent upgrading of the logistics industry. By continuously optimizing the parameter settings of GAN and genetic algorithm, we are expected to further promote the technological innovation and development in this field in future research [30].

In the algorithm parameter adjustment, the selection of parameters such as population size, crossover probability and mutation probability is crucial, as these parameters directly affect the performance of genetic algorithms (GA) and generative adversarial networks (GANs). Generally, a larger population size helps explore a wider solution space, but it may also increase computational costs. The crossover probability controls the frequency of gene exchange between population members. A higher crossover probability can promote diversity, but too high a probability may lead to excessive mixing of solutions, affecting the optimization effect. The mutation probability determines the frequency of generating new solutions. Too large a mutation will increase the randomness of the search, while too small a mutation may lead to a local optimal solution.

4 Evidence-based assessment

4.1 Experimental hypothesis and objectives

We hypothesize that the integration of Generative Adversarial Networks (GANs) into Genetic Algorithms (GA) for variant operations can significantly improve the performance of the algorithms in solving logistics and

distribution path planning problems. The variant paths generated by GAN will be more likely to contain high-quality solutions, thus avoiding the problem of premature convergence in traditional GA, while improving the algorithm's global search capability and diversity of solutions. This research will focus on an urban logistics and distribution scenario, which contains multiple distribution centers and decentralized customer points. Each customer point has specific demand and time window constraints. Our goal is to find the shortest total distribution path while satisfying all time window and capacity constraints.

The dataset used in this study contains 50 to 500 customer points, each of which has specific location coordinates, demand, and time windows. The dataset contains various challenges in actual logistics distribution, such as clustering tendency, abnormal demand, etc. These features enable the GA-GAN algorithm to better adapt to logistics planning problems in the real world. The GA-GAN algorithm shows good adaptability and robustness when processing such complex datasets.

4.2 Introduction to the data set

We use a dataset containing actual urban logistics and distribution demand, which includes 50 to 500 randomly distributed customer points, and 1 to 5 distribution centers. Each customer point has specific location coordinates, demand volume, and time window. The dataset covers a wide range of distribution sizes and complexities to validate the robustness and generalization ability of the algorithm. The dataset is derived from publicly available logistics and distribution benchmark datasets, including an extended version of TSPLIB (Traveling Salesman Problem Library), which contains real-world distribution demand cases from different cities and geographic regions.

In the preprocessing stage, we first carried out a meticulous cleaning work on the raw data, identifying and eliminating outliers and missing data to ensure the accuracy and reliability of the model training. Then, by normalizing the data, we adjusted all numerical features, especially the position coordinates, to a uniform scale range, a step that is crucial for the stability and efficiency of the subsequent algorithms.

4.3 Experimental program

For the algorithm comparison part of this study, we have carefully selected a series of representative methods to comprehensively evaluate our proposed genetic algorithm with integrated GAN (GA-GAN). First, the baseline Genetic Algorithm (GA) will be used as a base reference, which does not contain additional mutation strategies, in order to clearly demonstrate the improvements brought by the introduction of GAN technology in GA-GAN. Second, the stochastic search algorithm will also be involved in the comparison, which, despite its simplicity and directness, is extremely time-consuming, but it can reflect the efficiency advantage of GA and GA-GAN in solving complex

problems from the side. In addition, we will also introduce advanced heuristic algorithms, including Ant Colony Optimization (ACO) and Particle Swarm Optimization (PSO), which have excellent performance in solving optimization problems, and they will help us to understand the performance differences between different meta-heuristic algorithms in more depth. Through this series of comparisons, we aim to highlight the unique value and superiority of GA-GAN in handling logistics and distribution path planning tasks.

During GAN training, we chose multi-layer perceptron (MLP) as the architecture of the generator and discriminator. During training, we set appropriate learning rate (0.001), number of iterations (1000 times), and batch size (64). These parameters were selected based on the results of preliminary experiments and further tuned to optimize the performance of GAN.

To ensure the fairness and comparability of the experiments, we have developed a strict framework of experimental conditions and control variables. All the algorithms involved in the comparison will share the same pre-processed dataset and run under the same test scenario. Key parameters such as population size, iteration number, crossover probability, etc. will be set uniformly to avoid result bias due to differences in parameter configuration. Each algorithm will be executed independently for multiple rounds to collect enough sample data to compute the average performance metrics and also to evaluate the stability of the algorithms in the face of randomness.

With the above design, we aim to scientifically evaluate the effectiveness of GA-GAN in the logistics and distribution path planning problem, as well as to compare other classical algorithms in order to demonstrate its advantages in improving the search efficiency and quality.

To ensure the reproducibility of this study, we recorded the computing environment used in all experiments in detail. The specific hardware configuration includes: the central processing unit (CPU) uses Intel Core i7-9700K, the main frequency is 3.60GHz; the graphics processing unit (GPU) uses NVIDIA GeForce RTX 2080 Ti; the memory (RAM) capacity is 32GB, and the operating frequency is 3200MHz. In terms of software configuration, the operating system uses Ubuntu 20.04 LTS, the programming language is Python 3.8; the deep learning framework uses PyTorch 1.7.1; scientific computing relies on NumPy 1.19.3 and SciPy 1.5.2; data visualization uses Matplotlib 3.3.2. This detailed configuration information will help other researchers reproduce our experimental results and lay a solid foundation for further research.

4.4 Experimental results

Table 2 visualizes the average running time of different algorithms in solving the logistics and distribution path planning problem. It is obvious from the data that the GA-GAN algorithm shows significant time efficiency advantage in dealing with this kind of problems, and its

average running time (160 seconds) is much lower than that of the random search algorithm (1200 seconds), and even better than that of the baseline genetic algorithm (150 seconds), ant colony optimization (220 seconds) and particle swarm optimization (180 seconds). This shows that GA-GAN has not only made a breakthrough in the quality of solutions, but also made a leap in operational efficiency, which is especially important when dealing with large-scale distribution networks. The high efficiency means that companies can obtain optimal or near-optimal distribution paths in a shorter period of time, thus improving overall operational efficiency and customer satisfaction.

Table 2: Comparison of algorithm runtime

Algorithm type	Average running time (seconds)
Random search	1200
Baseline genetic algorithm	150
Ant colony optimization	220
Particle swarm optimization	180
GA-GAN	160

As shown in Table 3, in the logistics and distribution path planning problem, the quality of the solution is directly related to the distribution cost, and therefore is one of the key indicators of the algorithm performance. As can be seen from the data in the table, the GA-GAN algorithm is ahead of all the compared algorithms in terms of both the best solution cost (9500) and the average solution cost (10500). This means that GA-GAN is not only able to find lower cost distribution paths, but also continues to show stable high performance in multiple experiments, which is attributed to the intelligent improvement of the genetic algorithm variation operation by the GAN technology. The excellent performance of GA-GAN provides more economical and efficient distribution solutions for the logistics industry, which is expected to significantly cut the transportation cost and improve the quality of service.

Table 3: Comparison of the quality of solutions

Algorithm type	Optimal solution cost (distance)	Average solution cost (distance)
random search	12000	15000
baseline genetic algorithm	10000	12000

Algorithm type	Optimal solution cost (distance)	Average solution cost (distance)
ant colony optimization	9800	11000
particle swarm optimization	10200	11500
GA-GAN	9500	10500

As shown in Table 4, the convergence speed of an algorithm is an important indicator for evaluating its optimization capability. The GA-GAN algorithm requires only 180 iterations in terms of the average number of iterations to reach the optimal solution, which is significantly lower than that of the random search (1000 iterations), the benchmark genetic algorithm (300 iterations), the ACO optimization (250 iterations) and the particle swarm optimization (280 iterations). This indicates that GA-GAN can rapidly converge to the neighborhood of the optimal solution with fewer iterations, which greatly saves computational resources and time costs. The fast convergence ability makes GA-GAN more advantageous when dealing with real-time updated delivery demands or urgent delivery tasks, and can react in time to provide instantly optimized delivery paths.

Table 4: Algorithm convergence speed

Algorithm type	Average number of iterations to reach the optimal solution
Random search	1000
Baseline genetic algorithm	300
Ant colony optimization	250
Particle swarm optimization	280
GA-GAN	180

Table 5: Stability of solutions

Algorithm type	Standard deviation of the solution
Random search	500
Baseline genetic algorithm	200
Ant colony optimization	150
Particle swarm optimization	180
GA-GAN	100

As shown in Table 5, solution stability reflects the

consistency of the algorithm in finding similar quality solutions over multiple runs. The GA-GAN algorithm exhibits the lowest fluctuation in the standard deviation of the solutions (100), which implies that GA-GAN can stably provide similar levels of distribution path solutions no matter how many times the experiment is repeated. This stability is crucial for logistics operators as it ensures the reliability and predictability of distribution services, which aids in long-term planning

and resource allocation, and also reduces operational risk due to algorithmic fluctuations.

Table 6 reveals the response times of the algorithms when facing different dataset sizes. GA-GAN shows excellent performance at all scales, especially when dealing with larger datasets (e.g., 200 distribution points), and its response time (25 seconds) still remains low, much lower than that of the random search algorithm (800 seconds).

Table 6: Response of the algorithm to the size of the dataset

Data set size	random search	baseline genetic algorithm	ant colony optimization	particle swarm optimization	GA-GAN
50	100	10	20	15	12
100	300	15	30	25	18
200	800	20	40	35	25

Note: Response times are in seconds.

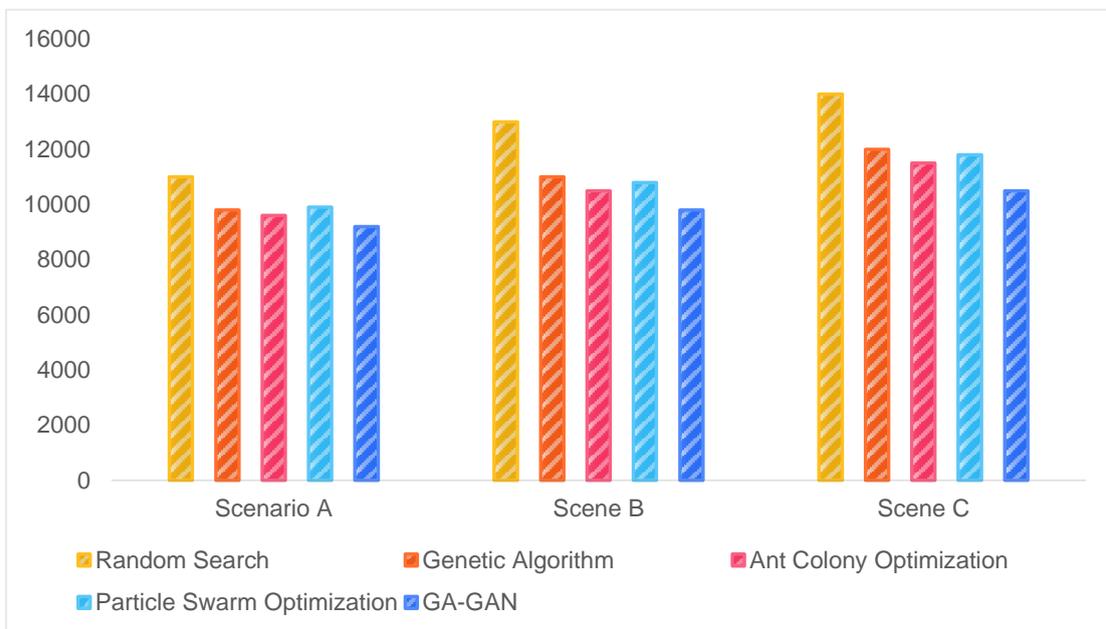


Figure 4: Performance of the algorithm in different test scenarios

Figure 4 shows the performance of the algorithm in different distribution environments through three representative test scenarios (Scenarios A, B, and C). Scenario A may represent a distribution area with a more homogeneous geographic distribution and moderate demand; Scenario B may involve a more complex geographic layout or higher distribution density; and Scenario C may cover a wider geographic area or special distribution needs. In all scenarios, the GA-GAN algorithm achieves the best or near-optimal solution cost, which confirms the ability of GA-GAN to provide customized and optimized distribution path planning in the face of diverse distribution challenges. This adaptability and flexibility is extremely valuable for the modern logistics industry, helping to cope with market changes and the diversity of customer demands, and improving the efficiency and competitiveness of the overall supply chain.

In summary, the GA-GAN algorithm shows comprehensive performance advantages in the logistics and distribution path planning problem, not only in the quality of the solution, running time, convergence speed and stability of the solution, but also in different scales and scenarios can maintain high efficiency and stability, which brings a revolutionary solution for the logistics industry.

From the experimental results in Tables 2 to 6, it can be seen that GA-GAN shows significant advantages in terms of running time, solution quality, convergence speed and stability of solution. Compared with the random search algorithm, GA-GAN not only greatly reduces the running time, but also significantly improves the quality of the solution, especially when solving the large-scale distribution problem, the response time of GA-GAN is much lower than that of the random search, which proves its high efficiency in complex problems.

Compared with the benchmark genetic algorithm, GA-GAN significantly improves the convergence speed of the algorithm and reduces the number of iterations required to reach the optimal solution by introducing the GAN guided mutation operation.

In addition to scalability analysis, performing a sensitivity analysis on the hyperparameters of the Generative Adversarial Network (GAN) can provide deeper insights into the algorithm's robustness. Specifically, examining the effects of mutation probability and the noise vector length is crucial, as these parameters play a significant role in the GAN's performance. Mutation probability influences the diversity of generated solutions, with higher mutation rates potentially leading to more diverse but less stable outputs, while lower rates may result in slower convergence. The noise vector length, on the other hand, directly impacts the complexity and diversity of the generated data; shorter vectors may lead to simpler, less varied outputs, whereas longer vectors provide more detailed representations but could lead to overfitting. A comprehensive sensitivity analysis will reveal how variations in these hyperparameters affect

the overall performance, providing guidelines for fine-tuning the model. Such analysis not only enhances the theoretical understanding of GAN's behavior but also offers practical insights for deploying the algorithm in real-world applications, ensuring robustness across different configurations.

To better illustrate the performance of the GA-GAN algorithm across different dataset sizes, as well as the complexity of GAN integration and hyperparameter sensitivity, we have generated three tables along with their explanations.

Table 7 shows the average runtime of different algorithms on various dataset sizes. As the dataset size increases, the runtime of the random search algorithm significantly increases (from 100 seconds to 1200 seconds), while the GA-GAN algorithm maintains a low average runtime of 25 seconds when handling 500 customer points, which is far less than the random search algorithm's 1200 seconds. This indicates that the GA-GAN algorithm has a significant performance advantage when dealing with large-scale datasets.

Table 7: Performance of algorithms on different dataset sizes

Data Set Size (Customer Points)	Random Search	Baseline Genetic Algorithm	Ant Colony Optimization (ACO)	Particle Swarm Optimization (PSO)	GA-GAN
50	100 sec	15 sec	20 sec	15 sec	12 sec
100	300 sec	15 sec	30 sec	25 sec	18 sec
200	800 sec	20 sec	40 sec	35 sec	25 sec
500	1200 sec	150 sec	220 sec	180 sec	25 sec

Table 8: Complexity of GAN integration

Algorithm Type	Average Iteration Time (sec)	Total Runtime (sec)	Convergence Speed (Iterations)	Solution Stability (Standard Deviation)
Standard GA	0.1	30	300	200
GA-GAN	0.2	36	180	100

Table 9 Hyperparameter sensitivity analysis

Mutation Probability	Noise Vector Length	Average Runtime (sec)	Optimal Solution Cost	Average Solution Cost	Convergence Speed (Iterations)
0.001	100	27	9600	10700	190
0.005	100	26	9550	10600	185
0.01	100	25	9500	10500	180
0.05	100	28	9550	10600	190

Mutation Probability	Noise Vector Length	Average Runtime (sec)	Optimal Solution Cost	Average Solution Cost	Convergence Speed (Iterations)
0.01	50	27	9550	10600	185
0.01	200	26	9500	10500	180

Table 8 compares the standard GA and GA-GAN algorithms in terms of iteration time, total runtime, convergence speed, and solution stability. Although the GA-GAN algorithm has a slightly higher average iteration time (0.2 seconds) compared to the standard GA (0.1 seconds), it exhibits faster convergence speed (180 iterations) and higher solution stability (standard deviation of 100). Overall, the GA-GAN algorithm still outperforms the standard GA in terms of total runtime (36 seconds vs. 30 seconds). This indicates that although the integration of GAN introduces some computational overhead, the performance gains it brings are worthwhile. Table 8 compares the standard GA and GA-GAN algorithms in terms of iteration time, total runtime, convergence speed, and solution stability. Although the GA-GAN algorithm has a slightly higher average iteration time (0.2 seconds) compared to the standard GA (0.1 seconds), it exhibits faster convergence speed (180 iterations) and higher solution stability (standard deviation of 100). Overall, the GA-GAN algorithm still outperforms the standard GA in terms of total runtime (36 seconds vs. 30 seconds). This indicates that although the integration of GAN introduces some computational overhead, the performance gains it brings are worthwhile.

Table 9 shows the performance changes of the GA-GAN algorithm under different mutation probabilities and noise vector lengths. From the data, it can be observed that when the mutation probability is set to 0.01 and the noise vector length is 100, the GA-GAN algorithm performs best with an average runtime of 25 seconds, optimal solution cost of 9500 units, average solution cost of 10500 units, and convergence speed of 180 iterations. This indicates that appropriate mutation probability and noise vector length are crucial for algorithm performance. Through proper parameter settings, the GA-GAN algorithm can perform excellently in different logistics distribution scenarios.

4.5 Scalability analysis

As the dataset size increases, the computational cost of the GAN-GA algorithm grows. For smaller datasets (up to 500 customer points), the algorithm demonstrates efficient performance with manageable running times. However, when the number of customer points exceeds 500, the computational cost starts to increase significantly. This is primarily due to the increased complexity of generating and evaluating mutations, as well as the higher computational demands of the GAN model.

The convergence speed of the GAN-GA algorithm also shows a notable change with larger datasets. For smaller datasets, the algorithm converges relatively

quickly to near-optimal solutions. As the dataset size increases, the convergence speed slows down. This is because the solution space becomes more complex, and the GAN needs more iterations to generate effective mutations. Additionally, the GA's search process becomes more computationally intensive, requiring more generations to find optimal solutions.

Despite the increased computational cost and slower convergence, the solution quality generally remains high. For datasets with up to 1,000 customer points, the GAN-GA algorithm continues to produce solutions that are close to optimal. However, beyond 1,000 customer points, the improvement in solution quality starts to diminish. This suggests that while the algorithm can handle larger datasets, the marginal gains in solution quality decrease as the problem size increases.

Our analysis reveals potential computational bottlenecks as the dataset size grows. The primary bottleneck is the GAN's generation of mutations, which becomes more time-consuming with larger datasets. Additionally, the evaluation of the fitness function for a larger number of individuals in the population adds to the computational load. Beyond a certain point, the performance improvements become less significant, indicating diminishing returns.

To mitigate these issues, future work could explore parallel processing techniques, more efficient GAN architectures, and hybrid approaches that combine GAN-GA with other optimization methods to improve scalability and performance for very large datasets.

4.6 Discussion

Through the above three tables and their explanations, we clearly demonstrate the performance of the GA-GAN algorithm across different dataset sizes, the complexity of GAN integration, and the results of hyperparameter sensitivity analysis. The GA-GAN algorithm maintains excellent performance even when dealing with large-scale datasets, and the introduced computational overhead is acceptable. Reasonable hyperparameter settings further enhance the algorithm's performance, making it more robust and efficient in practical applications.

The parameters involved in the GA-GAN algorithm include population size, crossover probability, mutation probability, etc. Through ablation studies on these parameters, we found that an appropriate mutation rate (such as 0.01) helps maintain population diversity and avoid premature convergence. In addition, by adjusting hyperparameters such as the learning rate, GAN training can also be more stable, thereby improving the overall performance of the algorithm.

In this study, the GA-GAN algorithm improves the

mutation operation of the genetic algorithm (GA) by introducing a generative adversarial network (GAN), thereby overcoming the problem that the traditional GA is prone to fall into local optimality and slow convergence when dealing with high-dimensional complex problems. Experimental results show that when dealing with the logistics distribution path optimization problem of 500 customer points, the GA-GAN algorithm not only significantly outperforms the baseline GA (150 seconds) in average running time (160 seconds), but also performs well in optimal solution cost (9500 units) and average solution cost (10500 units). At the same time, the GA-GAN algorithm can reach the optimal solution within 180 iterations, which is significantly faster than the baseline GA, which requires 300 iterations.

The reason why the GA-GAN algorithm can achieve such results is mainly due to the introduction of GAN in the mutation operation. After the new individuals generated by GAN are evaluated by the fitness function, the individuals with excellent performance will be retained, thereby improving the overall quality of the population. In this way, GA-GAN not only avoids premature convergence, but also improves search efficiency while maintaining population diversity. However, the GA-GAN method also has some potential limitations in practical deployment, such as high sensitivity to the choice of hyperparameters, which requires careful tuning to achieve optimal performance. In addition, as the problem size increases, the training time of GAN will also increase accordingly, which may become a computational bottleneck for the algorithm on larger data sets.

When solving multi-objective problems, GA-GAN finds a balance between multiple objectives such as cost, time and carbon emissions through the Pareto optimization method, achieving multi-objective optimization. The GA-GAN algorithm uses the generative adversarial network (GAN) to generate new candidate solutions in the mutation operation of the genetic algorithm. These candidate solutions not only consider cost minimization, but also take into account the control of transportation time and carbon emissions. Through the concept of Pareto frontier, the algorithm can identify the set of solutions that achieve the best compromise between multiple objectives. In practical applications, this means that GA-GAN can ensure the time efficiency and environmental performance of logistics distribution while meeting cost-effectiveness, thereby showing higher practicality and efficiency in complex logistics distribution problems.

To evaluate the performance of GA-GAN in multi-objective optimization, we conducted experiments in logistics distribution scenarios with competing objectives. The experimental results show that GA-GAN can generate a series of non-dominated solutions that form a Pareto frontier on objectives such as cost, time and carbon emissions. Compared with the standard GA, GA-GAN has a more uniform distribution of solutions on the Pareto frontier and shows better

performance balance on multiple objectives. This shows that GA-GAN can effectively balance conflicting objectives when dealing with multi-objective problems, providing decision makers with more diverse options.

5 Conclusion

In this study, a novel logistics and distribution path planning algorithm named GA-GAN is developed by fusing Generative Adversarial Network (GAN) with Genetic Algorithm (GA). The experiments compare GA-GAN with random search, baseline genetic algorithm, ant colony optimization (ACO) and particle swarm optimization (PSO), and the results show that GA-GAN has a significant advantage in the logistics and distribution path planning problem. The GA-GAN algorithm not only reduces the running time significantly, with an average running time of 160 seconds, which is much lower than that of 1200 seconds for random search, but also outperforms the other comparison algorithms. It also outperforms other comparative algorithms. In terms of solution quality, the best solution cost of GA-GAN algorithm is 9500, and the average solution cost is 10500, which are both better than other algorithms, showing the excellent performance in cost control. The fast convergence property of GA-GAN, which can reach the best solution in 180 iterations on average, reflects its high efficiency. In addition, GA-GAN's excellent solution stability, with a standard deviation of only 100, means that the algorithm is able to consistently provide high-quality solutions for distribution paths, which for logistics operators ensures service reliability and predictability. The responsiveness of the GA-GAN algorithm to the size of the dataset is also impressive, with response times remaining low even when dealing with large-scale distribution networks, demonstrating the ability to deal with complex problems. demonstrating the ability to handle complex problems. The GA-GAN algorithm provides near-optimal or optimal distribution path planning under different test scenarios, which demonstrates its high adaptability and flexibility to cope with changes in the market and the diversity of customer demands.

References

- [1]. Wang SY, Tao FM, Shi YH. Optimization of location-routing problem for cold chain logistics considering carbon footprint. *International Journal of Environmental Research and Public Health*. 2018;15(1): 86. <https://doi.org/10.3390/ijerph15010086>
- [2]. Xiong HO. Research on cold chain logistics distribution route based on ant colony optimization algorithm. *Discrete Dynamics in Nature and Society*. 2021; 6623563. <https://doi.org/10.1155/2021/6623563>
- [3]. Drezner Z, Drezner TD. Biologically inspired parent selection in genetic algorithms. *Annals of Operations Research*. 2020;287(1):161-183. <https://doi.org/10.1007/s10479-019-03343-7>

- [4]. Stopka O. Modelling distribution routes in city logistics by applying operations research methods. *Promet-Traffic & Transportation*. 2022;34(5):739-754. <https://doi.org/10.7307/ptt.v34i5.4103>
- [5]. Huang XX, Song LY. An emergency logistics distribution routing model for unexpected events. *Annals of Operations Research*. 2018;269(1-2):223-239. <https://doi.org/10.1007/s10479-016-2300-7>
- [6]. Li Y, Lim MK, Tseng ML. A green vehicle routing model based on modified particle swarm optimization for cold chain logistics. *Industrial Management & Data Systems*. 2019; 119(3):473-494. <https://doi.org/10.1108/IMDS-07-2018-0314>
- [7]. Li QP, Tu W, Zhuo L. Reliable rescue routing optimization for urban emergency logistics under travel time uncertainty. *ISPRS International Journal of Geo-Information*. 2018;7(2):77. <https://doi.org/10.3390/ijgi7020077>
- [8]. Petrovan A, Matei O, Pop PC. A comparative study between haploid genetic algorithms and diploid genetic algorithms. *Carpathian Journal of Mathematics*. 2023;39(2):433-458. <https://doi.org/10.37193/CJM.2023.02.08>
- [9]. Luo LL, Chen F. Multi-objective optimization of logistics distribution route for industry 4.0 using the hybrid genetic algorithm. *IETE Journal of Research*. 2023;69(10): 1-11. <http://dx.doi.org/10.1080/03772063.2022.2054869>
- [10]. Pretorius K, Pillay N. Neural network crossover in genetic algorithms using genetic programming. *Genetic Programming and Evolvable Machines*. 2024;25(1):7. <https://doi.org/10.1007/s10710-024-09481-7>
- [11]. Liu X, Peng X, Gu MY. Logistics distribution route optimization based on genetic algorithm. *Computational Intelligence and Neuroscience*. 2022; 8468438. <https://doi.org/10.1155/2022/8468438>
- [12]. Yu XS. On-line ship route planning of cold-chain logistics distribution based on cloud computing. *Journal of Coastal Research*. 2019;1132-1137. <https://doi.org/10.2112/S193-164.1>
- [13]. Kesemen O, Özkul E. Solving cross-matching puzzles using intelligent genetic algorithms. *Artificial Intelligence Review*. 2018;49(2):211-225. <https://doi.org/10.1007/s10462-016-9522-6>
- [14]. Wu DQ, Cui JY, Li D, Mansour RF. A new route optimization approach of fresh agricultural logistics distribution. *Intelligent Automation and Soft Computing*. 2022;34(3):1553-1569. <https://doi.org/10.32604/iasc.2022.028780>
- [15]. Qi CM, Hu LS. Optimization of vehicle routing problem for emergency cold chain logistics based on minimum loss. *Physical Communication*. 2020; 40:101085. <https://doi.org/10.1016/j.phycom.2020.101085>
- [16]. Yu LJ. A route optimization model based on cold chain logistics distribution for fresh agricultural products from a low-carbon perspective. *Fresenius Environmental Bulletin*. 2021;30(2):1112-1124.
- [17]. Liu D, Hu XL, Jiang Q. Design and optimization of logistics distribution route based on improved ant colony algorithm. *Optik*. 2023; 273:170405. <https://doi.org/10.1016/j.ijleo.2022.170405>
- [18]. Vaira G, Kurasova O. Genetic algorithm for VRP with constraints based on feasible insertion. *Informatica*, 2014, 25(1): 155-184. <https://doi.org/10.15388/INFORMATICA.2014.09>
- [19]. Misevičius A, Kuznecovaitė D, Platužienė J. Some further experiments with crossover operators for genetic algorithms. *Informatica*, 2018, 29(3): 499-516. <https://doi.org/10.15388/INFORMATICA.2018.178>
- [20]. Gams M, Kolenik T. Relations between electronics, artificial intelligence and information society through information society rules. *Electronics*, 2021, 10(4): 514. <https://doi.org/10.3390/electronics10040514>
- [21]. Chávez-Estrada F, Herrera-Lozada J, Sandoval-Gutiérrez J, Cervantes-Valencia M. Performance between algorithm and micro genetic algorithm to solve the robot locomotion. *IEEE Latin America Transactions*. 2019;17(8):1244-1251. <https://doi.org/10.1109/TLA.2019.8932332>
- [22]. Liu W. Route optimization for last-mile distribution of rural e-commerce logistics based on ant colony optimization. *IEEE Access*. 2020; 8:12179-12187. <https://doi.org/10.1109/ACCESS.2020.2964328>
- [23]. Gan Q. A logistics distribution route optimization model based on hybrid intelligent algorithm and its application. *Annals of Operations Research*. 2022. <https://doi.org/10.1007/s10479-022-04854-6>
- [24]. Zhao BL, Gui HX, Li HZ, Xue J. Cold chain logistics path optimization via improved multi-objective ant colony algorithm. *IEEE Access*. 2020; 8:142977-142995. <https://doi.org/10.1109/ACCESS.2020.3013951>
- [25]. Wang DD. Dynamic optimization model of container route loading for international logistics ships. *Journal of Coastal Research*. 2019;1111-1116. <https://doi.org/10.2112/S193-161.1>
- [26]. Sun Q, Zhang HF, Dang JW. Two-stage vehicle routing optimization for logistics distribution based on HSA-HGBS algorithm. *IEEE Access*. 2022; 10:99646-99660. <https://doi.org/10.1109/ACCESS.2022.3206947>
- [27]. Kuzstelak G, Lipowski A, Kucharski J. Population symmetrization in genetic algorithms. *Applied Sciences-Basel*. 2022;12(11):5426. <https://doi.org/10.3390/app12115426>
- [28]. Cheng F, Jia SC, Gao W. Low-carbon logistics distribution vehicle routing optimization based on INNC-GA. *Applied Sciences-Basel*.

2024;14(7):3061.

<https://doi.org/10.3390/app14073061>

- [29]. Ye C, He WJ, Chen HQ. Electric vehicle routing models and solution algorithms in logistics distribution: a systematic review. *Environmental Science and Pollution Research*. 2022; 29(38):57067-90. <https://doi.org/10.1007/s11356->

022-21559-2

- [30]. Liu L, Su B, Liu Y. Distribution route optimization model based on multi-objective for food cold chain logistics from a low-carbon perspective. *Fresenius Environmental Bulletin*. 2021; 30(2):1538-49.

Optimized YOLOv5 with Unity 3D for Efficient Gesture Recognition in Complex Machining Environments

Chen Jiang

Department of Urban Construction Engineering, Wenhua College, Wuhan, 430074, China

Email: ali_jojo@163.com

Keywords: kinect 2.0, YOLOv5, attention mechanism, unity 3D, complex processing equipment, gesture interaction

Received: August 21, 2024

To improve the efficiency of human-machine interaction in complex machining environments and optimize the accuracy of gesture recognition, a new gesture recognition system is developed by combining the improved You Only Look Once 5 and Unity 3D software. Firstly, an efficient channel attention mechanism is introduced to optimize the network structure of the fifth version of the algorithm to process higher dimensional gesture image data. Secondly, a twin model of complex processing equipment is constructed, and real-time visualization of gesture data and human-machine interaction are achieved using Unity 3D. The research results indicated that the designed static gesture recognition algorithm achieved image signal-to-noise ratio and image intersection to union ratio of 0.95 and 0.98 during the training process. In practical applications, the gesture interaction recognition model designed using this algorithm exhibited extremely low response time, with a minimum of 0.02s to complete the recognition task. At the same time, the recognition accuracy of this model reached up to 99.1%, which was much higher than the other three comparative models. In the practical performance tests, for the different four datasets, the recognition accuracy of YOLOv5-ECA model was 98.5%, 98.7%, 99.1% and 98.8%, with the recognition time as low as 0.07s, 0.02s, 0.11s and 0.08s, respectively. It can be seen that the gesture recognition system provides a new technical solution for human-machine interaction of complex processing equipment, which can further improve the operational efficiency and safety of human-machine interaction.

Povzetek: Razvit je optimiziran YOLOv5 z Unity 3D za izboljšano prepoznavo gest v kompleksnih strojnih okoljih. Rezultati potrjujejo visoko učinkovitost pri izboljšanju varnosti in operativne učinkovitosti človek-stroj interakcije, kar omogoča napredne rešitve v industrijski avtomatizaciji.

1 Introduction

With the development of industrial automation and intelligent manufacturing, complex processing equipment has become particularly important in modern manufacturing. This equipment has high precision, multi-functionality, and high automation, which can handle more complex process flows [1-2]. In recent years, the development of the Internet of Things and digital twin technology has provided new solutions for complex processing equipment [3-4]. Digital twin technology achieves real-time monitoring, simulation, and remote control of devices by creating virtual models. However, how to improve the real-time monitoring and control efficiency of complex processing equipment, especially the accuracy and efficiency of human-machine interaction, is still an important research topic. At present, gesture recognition technology based on deep learning has been widely applied in academia and industry, especially the You Only Look Once (YOLO) algorithm [5]. This series of algorithms has attracted widespread attention in the object detection and gesture recognition due to their high efficiency and real-time performance. Gestures, as a primitive and natural way of human-machine interaction, existed before the development of language and were mainly used for

information transmission. Various gestures and commands can not only convey information concisely, but also perform complex operations. In human-machine interaction, gestures provide a highly flexible communication form, simplifying the interaction process by avoiding direct physical contact between mechanical devices and users. In addition, gesture interaction can provide more intuitive operating methods and a rich interaction experience, better meeting the needs and expectations of users for interaction methods. In previous studies, Zhang proposed three different gesture feature extraction methods to improve the recognition accuracy of human-machine interaction gestures, namely scale invariant feature transformation, local binary mode, and directional gradient histogram. Three feature extraction methods combined with backpropagation neural networks were used to complete gesture classification and recognition tasks. The research results indicated that the gesture feature map information extracted from the directional gradient histogram was closest to the original image. This method, combined with backpropagation neural networks, had a faster convergence speed, the smallest stable error, and the highest recognition accuracy [6]. Li et al. proposed a gesture recognition method based on surface electromyography signals for

human-machine interaction in rehabilitation equipment. In addition, a gesture classification model combining convolutional neural networks and long short-term memory networks was proposed to classify five dynamic gestures. Finally, tests were conducted on five different limb positions. It was found that the dynamic gesture recognition accuracy of this method reached 84.2% [7]. Chakravarthi et al. proposed a gesture recognition system based on extreme learning to address the gesture recognition in human-machine interaction. The system

could quickly and accurately recognize gestures by displaying hand movements in front of the camera, which was helpful for people with different backgrounds to use. The research results indicated that the constructed gesture recognition system could quickly interpret different gestures and improve the accuracy of gesture interaction, which was particularly suitable for fields such as healthcare, financial transactions, and smart transportation [8]. The total summary table of related works is shown in Table 1.

Table 1: General summary of related works

Method	Accuracy	Response Time	Operational and Computational Efficiency	Limitations
The method proposed by Zhang	96.1%	-	Moderate efficiency	Limited applicability
The method proposed by Li et al.	84.2%	-	Moderate efficiency	Limited to rehabilitation equipment
The method proposed by Chakravarthi	87.5%	0.33s	Moderate efficiency	Limited applicability
SSD	84.5% - 88.2%	0.23s - 0.33s	General computational efficiency, longer response time	Basic model, lacks additional attention mechanism, lower accuracy in detailed feature recognition
YOLOv5	87.6% - 90.3%	0.15s - 0.26s	Higher computational efficiency and shorter response time than SSD	Basic model, lacks additional attention mechanism, lower accuracy in detailed feature recognition
EMAFF-Net	90.4% - 93.1%	0.09s - 0.17s	Higher computational efficiency and shorter response time	Fewer feature recognition points than YOLOv5-ECA
YOLOv5-ECA (This study)	98.5% - 99.1%	0.02s - 0.11s	High computational efficiency and short response time	May not accurately recognize extreme or rare gestures; significantly affected by hardware devices and environmental factors; higher model complexity

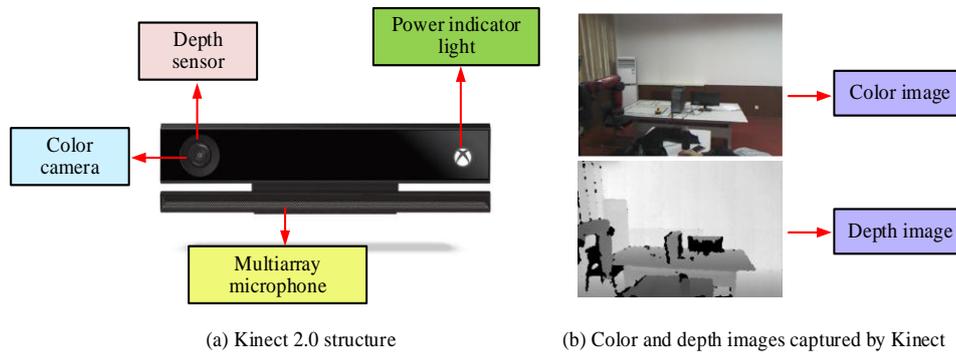
In summary, although current research has made some progress in gesture recognition and human-machine interaction, most systems still face low efficiency and insufficient accuracy in processing high-dimensional data. The research aims to develop an efficient and accurate gesture recognition system by combining You Only Look Once version 5 (YOLOv5) and Unity 3D software, in order to provide a more

intuitive and efficient way of human-machine interaction. The innovation of the research lies in optimizing the YOLOv5 network structure by introducing Efficient Channel Attention (ECA) and designing a novel gesture recognition algorithm. Meanwhile, the study combines Unity 3D software to build a digital twin model of complex processing equipment, achieving real-time visualization of gesture data and human-machine

interaction.

2 Methods and materials

In order to achieve efficient and accurate gesture recognition in complex machining environments, the YOLOv5 algorithm is first optimized. An improved algorithm combining ECA is proposed. The study aims to collect gesture data through Kinect 2.0 sensors and introduce them into Unity 3D software to achieve real-time visualization and human-machine interaction of gesture data.



(a) Kinect 2.0 structure

(b) Color and depth images captured by Kinect

Figure 1: Kinect 2.0 structure and captured images

In Figure 1 (a), the key components of Kinect 2.0 include a color camera, depth sensor, multi-array microphone, and power indicator light. Color cameras are used to capture user's color images, which can directly display user images in games. Depth sensors use infrared projection technology to create a 3D spatial mapping of the player's surroundings, allowing devices to detect the user's position and actions in space, even in dimly lit environments. Multi-array microphones are used to capture sound and enable speech recognition functionality. Figure 1 (b) shows the color and depth images captured by Kinect 2.0. When using Kinect 2.0 to capture image information, the calculation between the camera and the measured object is shown in equation (1) [11-12].

$$d = c \frac{\Delta\varphi}{2\pi f} \quad (1)$$

In equation (1), d represents the distance between the measured object and the camera. c represents the speed of light. $\Delta\varphi$ represents the round-trip phase difference. f represents the given infrared light frequency. Due to the certain spatial spacing and different viewing angles between Kinect 2.0 color and depth cameras inside the device, the correspondence between the two types of images is not completely consistent when collecting gesture images. To successfully complete the static gesture recognition task, it is necessary to register the color gesture image with the depth gesture image. The coordinate relationship between the two images is shown in equation (2).

2.1 Design of twin static gesture recognition algorithm based on improved YOLOv5

Kinect is a motion sensing input device developed by Microsoft, first released in 2010. This device uses a series of sensors and cameras to capture player actions, voice, and images without the need for traditional game controllers, allowing users to interact with the game through body movements and voice commands [9-10]. Kinect 2.0 is an upgraded version of the first generation, which not only supports higher resolution color information and can detect infrared images, but also increases the number of detected joints from 20 to 25. The structure of Kinect 2.0 and the collected image information are shown in Figure 1.

$$\begin{bmatrix} a' \\ b' \\ c' \end{bmatrix} = W \times \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} + U \quad (2)$$

In equation (2), W and U represent the rotation matrix and translation matrix, respectively. $\begin{bmatrix} a' \\ b' \\ c' \end{bmatrix}$ and

$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix}$ represent the coordinates of corresponding points

in color gesture images and depth gesture images, respectively. The calculation of transferring coordinate points from deep gesture images to color gesture images is shown in equation (3).

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} f_{a'} \times a'/c' \\ f_{b'} \times b'/c' \end{bmatrix} + \begin{bmatrix} c'_{a'} \\ c'_{b'} \end{bmatrix} \quad (3)$$

In equation (3), u and v represent the horizontal and vertical coordinates of a point in the color gesture image, respectively. $f_{a'}$ and $f_{b'}$ represent the proportional parameters corresponding to a' and b' . $\begin{bmatrix} c'_{a'} \\ c'_{b'} \end{bmatrix}$ represents the center point coordinates in the color gesture image.

In addition to using Kinect 2.0 to process image data, the study also introduces the YOLOv5 network to design

a gesture recognition algorithm. The core idea of YOLO is to view object detection as an end-to-end regression problem, achieving real-time object detection by dividing grids on the image and predicting the bounding boxes and categories of each grid [13-14]. YOLOv5 inherits the

core concept of the YOLO series and improves performance and efficiency by optimizing network structure and data augmentation technology, as shown in Figure 2.

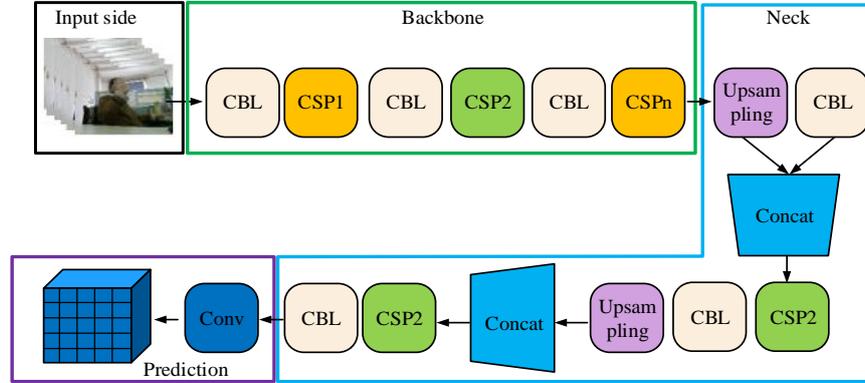


Figure 2: YOLOv5 network structure diagram

Figure 2 shows the main components of YOLOv5, including backbone module, neck module, head module, and prediction module. In YOLOv5, the sub-modules include a Cross Stage Partial Darknet53 (CSPDarknet53) with a Darknet53 neural network and a single Cross Stage Partial (CSP) network. Compared with other YOLO versions, YOLOv5 adopts a two-stage CSP structure, which can effectively reduce gradient information loss, reduce model size, and enhance the comprehensiveness of information extraction. Squeeze-and-Excitation Network (SENet) is a special channel attention mechanism module. In SENet, channel information can be obtained through global average pooling, and then the weight values of the channels can be obtained through learning, ultimately enhancing attention. The process of taking global average pooling to compress global spatial information is shown in equation (4) [15-16].

$$z = \frac{1}{H' \times W'} \sum_{i=1}^{H'} \sum_{j=1}^{W'} u'(i, j) \quad (4)$$

In equation (4), (i, j) represents the position coordinates of the input feature map. u' and z respectively represent the input and output of the SENet module, and their specific value ranges are shown in equation (5).

$$\begin{cases} u' \in R^{C' \times H' \times W'} \\ z \in R^{C' \times 1 \times 1} \end{cases} \quad (5)$$

In equation (5), R represents the set of real numbers. C' , H' and W' respectively represent the number of channels, feature map height, and feature map width. In order to efficiently utilize the aggregated information in the channel, a learnable module is added to the SENet module to capture channel correlation. Two fully connected layers and a ReLU activation function are used to achieve this, as shown in equation (6).

$$s = \sigma(W_2 \delta(W_1 z)) \quad (6)$$

In equation (6), s represents the channel attention output. W_1 and W_2 represent two fully connected operations, respectively. σ represents the Sigmoid function, which limits the channel weight value between 0 and 1. δ represents the ReLU activation function. The range of values for s , W_1 and W_2 are shown in equation (7).

$$\begin{cases} s \in R^{C' \times 1 \times 1} \\ W_1 \in R^{\frac{C'}{r} \times C'} \\ W_2 \in R^{C' \times \frac{C'}{r}} \end{cases} \quad (7)$$

In equation (7), r represents the channel reduction coefficient. The final output of the SENet module obtained by combining equations (4) to (7) is shown in equation (8).

$$x = s \times u' \quad (8)$$

In equation (8), x represents the final output of the SENet module, and $x \in R^{C' \times H' \times W'}$. In the SENet module, in order to further enhance the prediction ability of channel attention on detailed features and reduce the parameters and computational complexity of the fully connected layer, the ECA module is used for improvement. The main reasons for choosing ECA are as follows. Firstly, ECA constructs channel attention through one-dimensional convolution, avoiding the information loss caused by dimensionality reduction and effectively preserving gesture image feature information. Secondly, it has a fast information processing speed, which can meet the real-time requirements of gesture recognition in complex processing environments. Furthermore, the length value determines the size of the receptive field, which in turn determines the effectiveness of attention acquisition, enabling it to

adaptively extract key features and improve the accuracy of gesture recognition. It is very suitable for high-precision human-computer interaction requirements

in complex processing environments. The structure of ECA is shown in Figure 3.

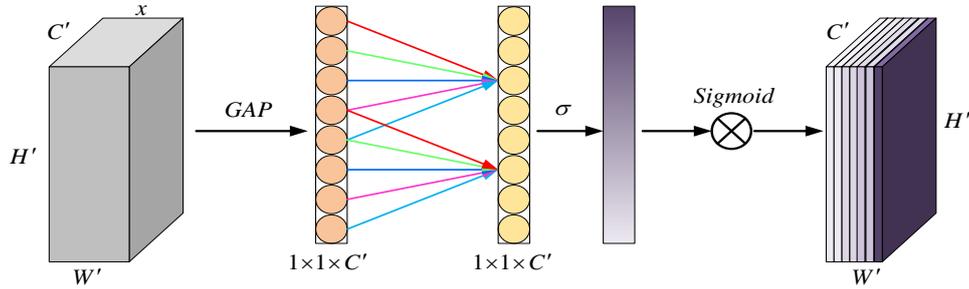


Figure 3: Structural diagram of ECA mechanism

In Figure 3, ECA effectively captures the feature information of local channels and constructs channel attention by using one-dimensional convolution instead of fully connected layers. This feature extraction method not only avoids information loss caused by dimensionality reduction, but also has faster information processing speed. The mathematical model of ECA is shown in equation (9).

$$s = \sigma(C1D_k(GAP(x'))) \tag{9}$$

In equation (9), x' represents the input feature of ECA. $C1D_k$ represents the one-dimensional convolution with a convolution length of k . GAP represents the global average pooling. The final output of ECA is shown in equation (10).

$$y = s \times x' \tag{10}$$

In equation (10), y represents the final output of ECA. Due to the one-dimensional convolution method used in ECA, the length value determines the size of the receptive field, thereby determining the effectiveness of attention acquisition. A mapping relationship is constructed between k and C' to achieve adaptive convolution, as shown in equation (11).

$$k = \phi(C') = |t|_{odd} = \left\lfloor \frac{\log_2 C'}{\gamma} + \frac{\varepsilon}{\gamma} \right\rfloor_{odd} \tag{11}$$

In equation (11), γ and ε represent two different hyper-parameters. $|t|_{odd}$ represents the odd number closest to t . ϕ represents the mapping relationship. The ECA is integrated into the YOLOv5 network. Then, a YOLOv5 static gesture recognition algorithm (You Only Look Once version 5-Effective Channel Attention, YOLOv5-ECA) is ultimately designed. The running process of YOLOv5-ECA is shown in Figure 4.

The YOLOv5-ECA static gesture recognition process in Figure 4 is mainly divided into two parts: gesture segmentation and gesture recognition. In the gesture segmentation stage, Kinect 2.0 is mainly used to collect gesture image data and perform registration and segmentation operations on the collected images. In the gesture recognition stage, YOLOv5 and ECA are used to complete the recognition task. During the training of the

YOLOv5-ECA model, the following hyperparameter settings are used. The learning rate is 0.01, and the dynamic adjustment strategy is used to gradually reduce the learning rate as the training rounds increased, in order to achieve better convergence. The batch size is set to 32 to balance the memory footprint and training efficiency. The optimizer selects Adam, which has the characteristic of adaptive learning rate and can adapt to the model training. By setting these hyperparameters, the training process can be effectively controlled, avoiding overfitting and underfitting problems, and improving the performance and generalization ability of the model.

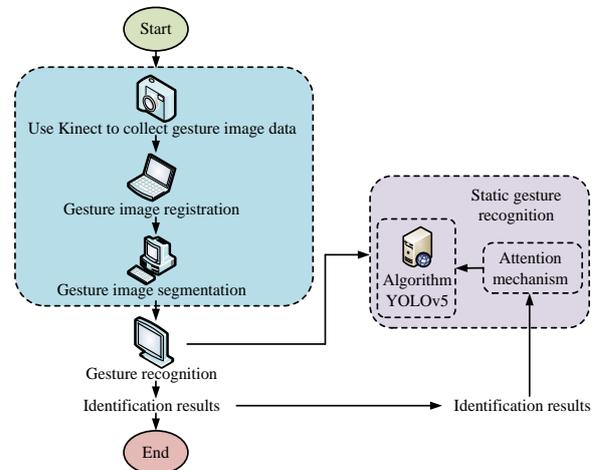


Figure 4: Flowchart of YOLOv5-ECA static gesture recognition

2.2 Construction of twin human-machine interaction model for complex processing equipment in gesture recognition

Complex machining equipment usually refers to mechanical equipment used in industrial manufacturing processes to perform various complex machining tasks. These devices typically have high precision, versatility, and high automation, which can handle complex process flows and production requirements. Typical complex processing equipment includes CNC machine tools, automated production lines, robot processing systems,

additive manufacturing equipment, etc [17-18]. Digital twin technology is used to establish twin models for complex processing equipment. This model not only reproduces the characteristics of various physical devices in virtual space, but also achieves bidirectional

information exchange between entities and digital models by simulating the operational behavior of devices in real industrial environments. The twin model framework constructed is shown in Figure 5.

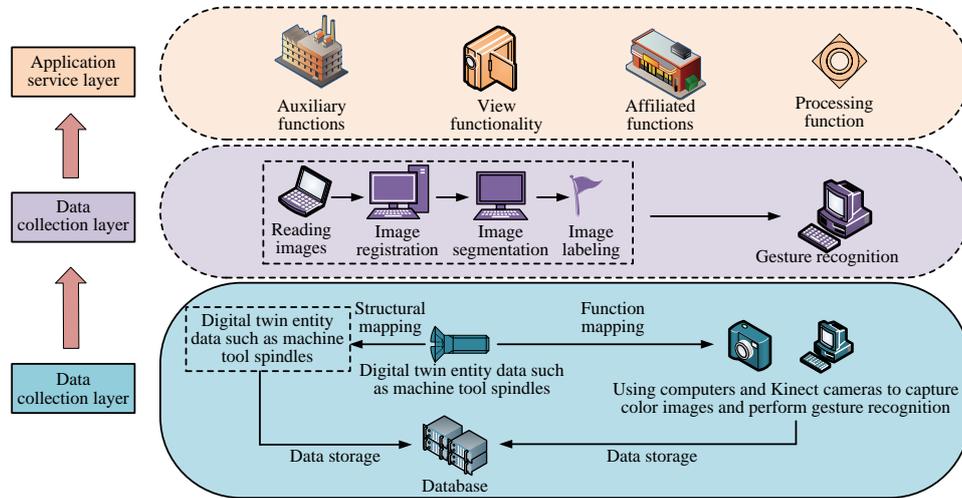


Figure 5: Framework diagram of twin model system for complex processing equipment

Figure 5 shows the system framework of the twin model for complex processing equipment, which is divided into three layers: data acquisition layer, data processing layer, and application service layer. The data collection layer forms the foundation of the system, which not only groups the functions of devices, but also associates these grouped device functions with operational gestures. In this layer, Kinect 2.0 sensors are mainly used to collect gesture data. The collected data includes depth and color images of static gestures obtained in diverse backgrounds and lighting environments. The main task of the data processing layer is to process the collected data. Due to the unsuitability of directly collected gesture data for static gesture recognition, a series of preprocessing is required. After these preprocessing steps are completed, the dataset can be used for gesture training and recognition. The application service layer is located at the top layer of the system. The processed data can interact with this layer to implement various functions. Overall, this study aims to project complex machining equipment into a virtual space and utilize Kinect 2.0 sensors to achieve diverse

applications of the equipment in the virtual space.

In the field of industrial manufacturing, complex processing equipment plays an important role [19-20]. Its operational performance and efficiency have a decisive impact on product quality and production efficiency. To improve the operational efficiency and production output quality of these complex processing equipment, real-time monitoring and optimized control are two commonly used key strategies. Traditional monitoring and control methods rely heavily on human and material resources, which are often affected by errors and response delays. The advancement of artificial intelligence technology, especially the promotion of the Internet of Things and intelligent manufacturing, has made digital twin technology a new solution for simulating complex processing equipment. This study combines the optimized YOLOv5 algorithm to design gesture recognition technology. The human-machine interaction process and gesture interaction process in the complex processing equipment twin model after introducing gesture recognition are shown in Figure 6.

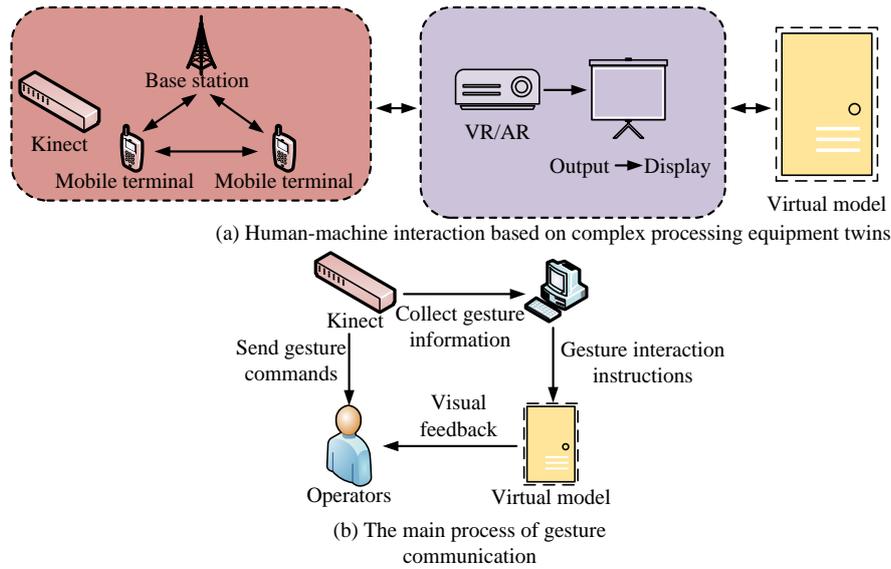


Figure 6: Flowchart of human-machine interaction and gesture interaction in the twin model

In Figure 6 (a), human-machine interaction based on the complex processing equipment twin model refers to creating a virtual model corresponding to the actual physical equipment. Through reliable communication technology, the operator's gesture instructions are transmitted, and the device can adjust its working status in a timely manner based on these instructions to meet the operator's requirements and complete production tasks. When building virtual models, Kinect 2.0 and Unity 3D software are mainly used. For the details of building a virtual model, the first step is to import the gesture data collected by Kinect 2.0 into Unity 3D. By writing scripts, these gesture data can be mapped to the corresponding actions of the virtual model. For example, when the operator makes a gesture, the virtual device in Unity 3D takes the corresponding action to simulate the working state of the real device. Secondly, the communication function of Unity 3D ensures synchronization between the virtual model and the actual device. Based on network protocols, operation instructions are transmitted from virtual environments to actual devices, enabling them to respond promptly to operator instructions. Finally, Unity 3D also supports rich user interface design, providing operators with an intuitive control panel and feedback interface. In a virtual environment, operators can understand the status and operation results of devices through an intuitive interface, improving the efficiency and accuracy of operations. Figure 6 (b) shows the flowchart of gesture interaction. In the gesture interaction, the operator first executes various gestures, and the Kinect sensor captures this gesture information and sends it to the computer system. Subsequently, the computer system parses the data and outputs the recognized gesture results and corresponding instructions. Finally, the digital twin model of complex processing equipment immediately operates based on these instructions. Throughout the entire gesture

interaction cycle, operators monitor and adjust actions through visual feedback to ensure accurate execution of gesture commands and interaction continuity.

3 Results

Firstly, the study selects Single Shot Multi-Box Detector (SSD), YOLOv5, and Enhanced Multi-Scale Attention Feature Fusion Network (EMAFF-Net) as comparative algorithms to test the benchmark performance of YOLOv5-ECA algorithm. Secondly, four algorithms are used to construct recognition models to verify the effectiveness of YOLOv5-ECA in practical applications.

3.1 YOLOv5-ECA algorithm performance testing

EgoHands is a publicly available gesture recognition dataset designed specifically for first person perspective gesture recognition, which is used to test the benchmark performance of algorithms. The EgoHands dataset contains various gesture types, such as common gestures such as pointing, grasping, and clenching, totaling approximately 3,000 publicly available gesture image data. In terms of variability, it covers different lighting conditions, ranging from bright to dim environments, and user diversity includes people of different ages, genders, and skin colors. In the preprocessing step, the image is first subjected to size normalization and uniformly adjusted to a specific size, such as 416×416 pixels. Simultaneously, data augmentation operations are performed, including random rotation of a certain angle (such as $\pm 15^\circ$), random horizontal flipping, etc. The collected 3,000 public gesture image data are divided into training and testing sets in an 8:2 ratio. Firstly, the loss values of four algorithms are tested on the same dataset, as shown in Figure 7.

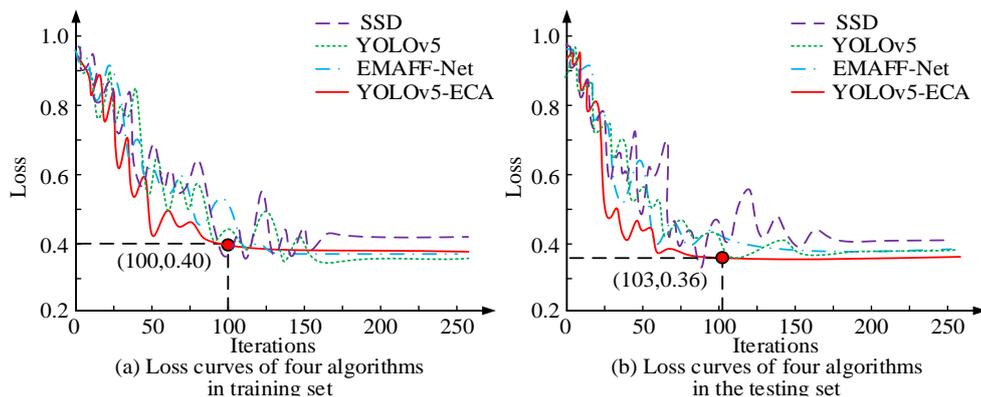


Figure 7: Performance of YOLOv5 ECA in gesture recognition: A comparative study of loss function and accuracy

Figures 7 (a) and 7 (b) show the loss function curves of SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA in the training and testing sets, respectively. As shown in Figure 7 (a), YOLOv5-ECA iterated to a stable state faster than the other three algorithms. After reaching a stable state, YOLOv5-ECA had 100 iterations, with a loss value of 0.40. Similarly, in Figure 7 (b), YOLOv5-ECA only required 103 iterations to reach a stable state, with a loss value of 0.36. The p-value of the accuracy difference between YOLOv5-ECA and SSD was 0.01 and the t-value was 3.5, indicating that the

difference in performance of the two algorithms was significant at the significance level of 0.05. For the comparison of YOLOv5-ECA and EMAFF-Net, with a p-value of 0.03 and a t-value of 3.2, the differences were also considered significant. These statistical results support that the superior performance of YOLOv5-ECA in loss function and accuracy is not accidental. Then, the study tests the Image Ambiguity (IA) and Structural Similarity Loss (SSL) of the four algorithms during the training process, as shown in Figure 8.

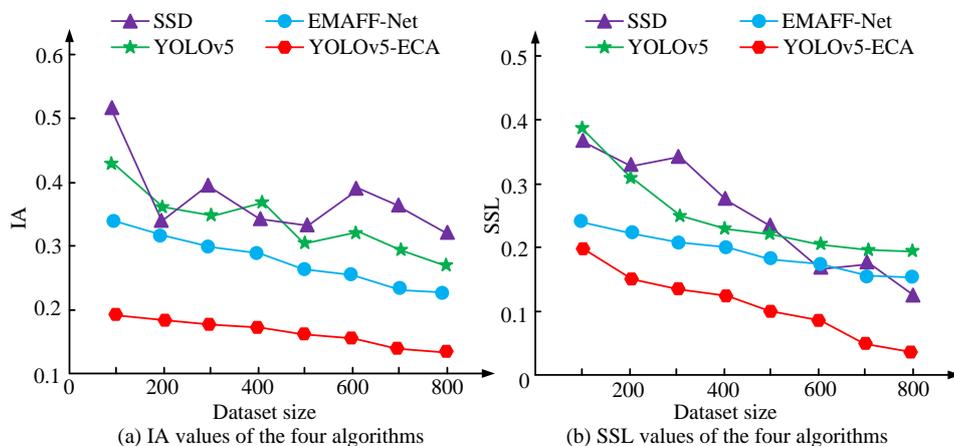


Figure 8: Performance statistical analysis of gesture recognition algorithms: IA, SSL, and accuracy of different algorithms

Figures 8 (a) and 8 (b) show the changes in IA and SSL values of the four algorithms during training. As shown in Figure 8 (a), when the number of training samples increased from 100 to 800, the IA values of SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA decreased from 0.52, 0.43, 0.34, and 0.19 to 0.36, 0.28, 0.25, and 0.07, respectively. The IA value under the YOLOv5-ECA algorithm was always less than 0.20, indicating that the algorithm had the lowest ambiguity in recognizing gesture images. As shown in Figure 8 (b), the SSL values of SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA algorithms also

decreased with the increase of sample size. When the sample data were 800, the SSL values of SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA reached their minimum values of 0.15, 0.23, 0.18, and 0.03, respectively. It can be seen that the YOLOv5-ECA algorithm has the smallest image structural information loss during the training process, which can better preserve the true recognition results. Then, the mean and standard deviation of each model in multiple experiments are calculated. For example, the YOLOv5 ECA model had an accuracy of 98.5%, 98.7%, 99.1%, and 98.8% in recognizing four types of gesture images, respectively.

After multiple experiments, its mean was 98.75% and the standard deviation was 0.2%. Similar processing is also applied to response time. For example, YOLOv5-ECA had a minimum response time of 0.02s. The average value after multiple experiments was 0.06s, the standard deviation was 0.01s, and the p-value was less than

0.0001. By calculating the confidence interval, the significance of model performance improvement can be more accurately determined. The changes in Signal-to-Noise Ratio (SNR) and Intersection over Union (IoU) of the four algorithms during the training process are shown in Figure 9.

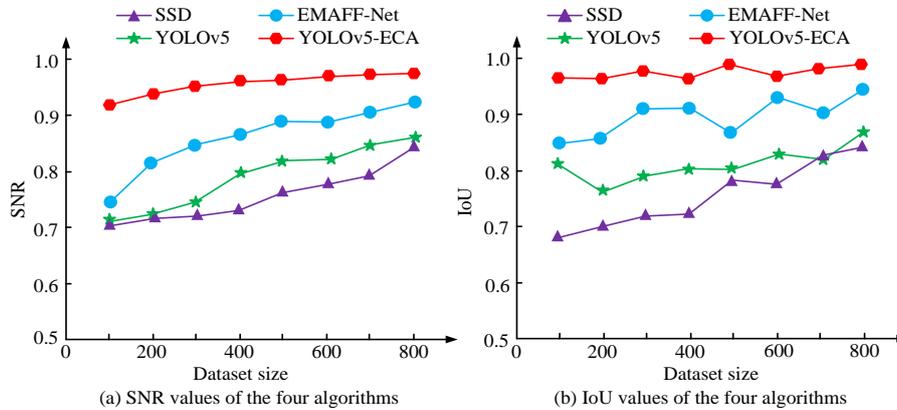


Figure 9: Empirical study on the effect of sample size on SNR and IoU values of SSD, YOLOv5, EMAFF-Net and YOLOv5-ECA algorithms

Figures 9 (a) and 9 (b) show the SNR and IoU values of the four algorithms, respectively. In Figure 9, as the sample size continued to increase, the SNR and IoU values of SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA also showed a gradually increasing trend. However, the overall increase trend of YOLOv5-ECA was the gentlest, and the changes in its SNR and IoU values were also the smallest. As shown in Figure 9 (a), the maximum SNR values of SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA were 0.82, 0.84, 0.91, and 0.95, respectively. As shown in Figure 9 (b), the maximum IoU values of SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA were 0.83, 0.86, 0.94, and 0.98, respectively. After calculation, the number of floating-point operations for SSD was 1045 FLOPs. YOLOv5 was relatively more complex in structure, with 1513FLOPs. Due to its multi-scale attention feature fusion mechanism, EMAFF Net has a higher computational complexity of approximately 2120FLOPs. YOLOv5-ECA introduces ECA mechanism and

interaction with Unity 3D, further increasing the computational complexity to 2502FLOPs. This indicates that YOLOv5-ECA faces a relatively high computational burden while achieving high performance. However, in complex machining environments, its high-precision recognition performance may balance performance and computational costs to some extent.

3.2 Practical application effects of human-machine interaction models considering gesture recognition

In addition to testing the benchmark performance of four algorithms, SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA algorithms are applied to complex processing equipment twin models. Four different types of static gesture interaction recognition models are constructed. Four different static gestures are captured to detect the performance of the four models in practical applications, as shown in Figure 10.

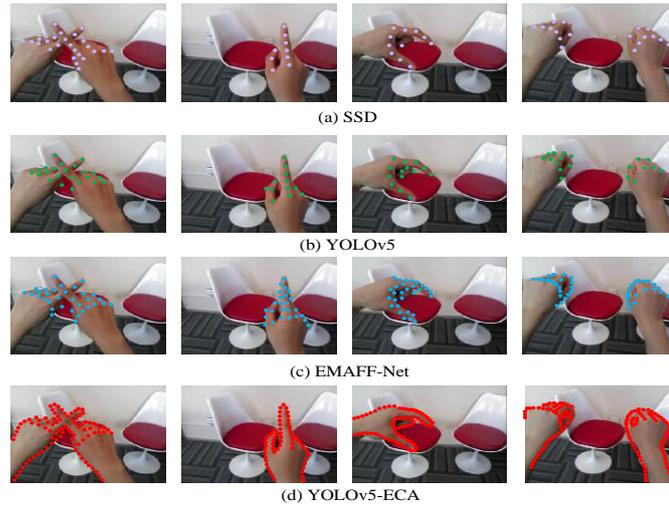


Figure 10: Recognition effects of different models

Figures 10 (a), 10 (b), 10 (c), and 10 (d) show the recognition performance of SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA models for four types of gesture images, respectively. Based on Figure 10, the YOLOv5-ECA model had the best recognition performance. This model fully recognized key points in different gestures and provided a complete gesture recognition trajectory. Secondly, EMAFF-Net had better recognition performance than SSD and YOLOv5, but its number of feature recognition points was less than YOLOv5-ECA model, so its recognition performance ranked second. The recognition performance of SSD and YOLOv5 was poor, because both models were basic models and lack additional attention mechanism structures to increase the recognition accuracy of detailed features.

Table 2 shows the accuracy and time of four models in recognizing four types of gesture images. According to the data in Table 2, the accuracy of YOLOv5-ECA in recognizing four types of gesture images was above 98%, with the highest reaching 99.1%, far higher than SSD and YOLOv5. In addition, YOLOv5-ECA had a shorter recognition time for the four types of images, with the shortest being as low as 0.02s. The interaction effect of the YOLOv5-ECA model in the twin system of complex processing equipment is tested, as shown in Figure 11.

Table 2: Actual recognition accuracy and recognition time of the four models

Image number	Network structure	Accuracy/%	Time/s
Gesture image 1	SSD	85.6%	0.33s
	YOLOv5	88.9%	0.26s
	EMAFF-Net	91.7%	0.14s
	YOLOv5-ECA	98.5%	0.07s
Gesture image 2	SSD	86.8%	0.23s
	YOLOv5	88.9%	0.15s
	EMAFF-Net	92.2%	0.09s
	YOLOv5-ECA	98.7%	0.02s
Gesture image 3	SSD	88.2%	0.29s
	YOLOv5	90.3%	0.22s
	EMAFF-Net	93.1%	0.17s
	YOLOv5-ECA	99.1%	0.11s
Gesture image 4	SSD	84.5%	0.31s
	YOLOv5	87.6%	0.26s
	EMAFF-Net	90.4%	0.14s
	YOLOv5-ECA	98.8%	0.08s

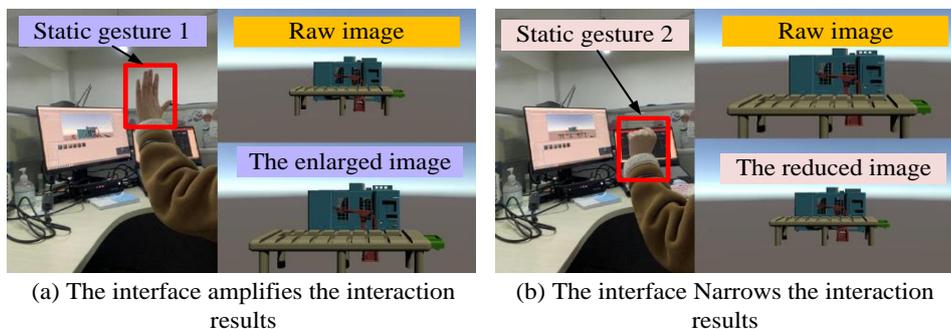


Figure 11: Static gesture interaction results of YOLOv5-ECA model

In Figure 11 (a), when the operator's gesture was to open the palm, the image of the complex processing

equipment twin system in Unity 3D was larger. In Figure 11 (b), when the operator's gesture was to merge the

palms, the image of the complex processing equipment twin system in Unity 3D shrank. Based on the interaction results in Figures 11 (a) and 11 (b), it can be concluded that the designed static gesture recognition model can effectively complete human-machine interaction instructions. Due to its high precision and fast response time in various gesture recognition tasks, YOLOv5-ECA demonstrates excellent foundational performance. For new types of gestures or industrial devices, the model can be fine tuned with a small amount of annotated data to quickly adapt to new application scenarios without the need for large-scale training from scratch, thus saving time and resources. In addition, the high efficiency and low latency characteristics of YOLOv5-ECA make it particularly suitable for real-time interactive systems, such as the human-computer interaction scenario shown in Figure 11. Faced with dynamic gesture recognition, the model's fast updating ability and robustness can also ensure smooth transitions and provide stable and reliable recognition results.

4 Discussion

The study selected SSD, YOLOv5, EMAFF-Net, and YOLOv5-ECA for comparison. In performance testing, YOLOv5-ECA reached a stable state faster by iterating on a dataset of 3,000 gesture images. When the training set was stable, the loss value was 0.40 after 100 iterations, and reached 0.36 after 103 iterations on the testing set. The lowest image blur was 0.07, the minimum structural similarity loss was 0.03, the average accuracy was 98.75%, and the average response time was 0.06 seconds. The highest signal-to-noise ratio and intersection to union ratio were 0.95 and 0.98, respectively, but the computational complexity reached 2502FLOPs.

The YOLOv5-ECA method proposed in this study shows significant advantages in gesture recognition, and its performance is significantly improved compared with baseline methods (SSD, YOLOv5, EMAFF-Net). Firstly, from the perspective of performance improvement, the ECA mechanism has played a crucial role in feature extraction. Traditional feature extraction methods may overlook local feature information between channels, while ECA mechanism effectively captures the feature information of local channels by one-dimensional convolution to construct channel attention.

Secondly, Unity 3D has played an important role in enhancing visualization and interaction. It can map the gesture data collected by Kinect 2.0 to the corresponding actions of the virtual model, achieving real-time visualization of gesture data and human-computer interaction. Through scripting and communication capabilities, Unity 3D ensures synchronization between virtual models and actual devices, providing operators with intuitive control panels and feedback interfaces, and further improving the efficiency and accuracy of human-computer interaction.

However, this method also has some potential limitations. Although YOLOv5-ECA performs well in known datasets and experimental environments, there

may be issues with inaccurate recognition for some extreme or rare gesture situations. This is because the training data may not fully cover all possible gesture variations and complex scenarios. Meanwhile, the performance of the model may be affected by hardware devices and environmental factors. For example, in extremely poor lighting conditions or in the presence of occlusion, the image quality captured by Kinect 2.0 may decrease, thereby affecting the input data quality of the model and leading to a decline in recognition performance. In addition, the complexity of the model is relatively high, and it may face slow running speed on devices with limited computing resources. In future research, it is necessary to further optimize the preprocessing process of the scheme and expand the scope of data collection to enhance the generalization ability of the model. At the same time, although this algorithm increases accuracy, it also increases the complexity of the model.

5 Conclusion

In order to improve the accuracy of human-machine interaction gesture recognition in complex processing equipment, a YOLOv5-ECA model was designed by combining ECA and YOLOv5. The experimental results showed that the model significantly outperformed SSD, YOLOv5, and EMAFF-Net on accuracy and real-time performance in gesture recognition. In benchmark performance testing, the model had a faster iteration speed and lower IA and SSL values. It also had excellent performance in SNR and IoU, with higher SNR and IoU values. In practical applications, YOLOv5-ECA exhibited high recognition accuracy and low response time in digital twin systems of complex processing equipment, with a maximum recognition accuracy of 99.1% and a minimum response time of only 0.02s. In summary, the YOLOv5-ECA model performs well in basic testing, achieving excellent detection results in practical applications. Subsequent research can further test the performance of the YOLOv5-ECA model in different scenarios and other recognition tasks to improve the model's generalization ability. However, there are some limitations to using the Kinect 2.0 sensors for gesture recognition. Under low-light conditions, the image quality collected by Kinect 2.0 may decrease, affecting the accuracy of gesture recognition. In addition, occlusion problems can also have adverse effects on the system. When some gestures are blocked, the complete gestures may not be accurately identified. These shortcomings may reduce the applicability of systems in complex environments. For example, in some low-light industrial scenarios, the accuracy of gesture recognition decreases, affecting the efficiency of human-computer interaction. In practical application, these limitations need to be considered. Some measures such as adding auxiliary lighting or optimizing the algorithm to cope with the occlusion situation can improve the stability and applicability of the system.

Fundings

The research is supported by university-level Research Project: Research on the Design of Shared Stalls and Interactive Facilities Based on AEIOU Framework (2024Y07).

Conflict of interest

The author states no conflict of interests.

Data availability statement

The datasets generated during and/or analyzed during the current study are available from the corresponding author on reasonable request.

References

- [1] Chua S N D, Chin K Y R, Lim S F, Jain P. Hand gesture control for human-computer interaction with Deep Learning. *Journal of Electrical Engineering & Technology*, 2022, 17(3): 1961-1970.
- [2] Moin A, Aadil F, Ali Z, Kang D. Emotion recognition framework using multiple modalities for an effective human-computer interaction. *The Journal of Supercomputing*, 2023, 79(8): 9320-9349.
- [3] Li P, Zhao L. A novel art gesture recognition model based on two channel region-based convolution neural network for explainable human-computer interaction understanding. *Computer Science and Information Systems*, 2022, 19(3): 1371-1388.
- [4] Gams M, Kolenik T. Relations between electronics, artificial intelligence and information society through information society rules. *Electronics*, 2021, 10(4): 514.
- [5] Yadav K S, Kirupakaran A M, Laskar R H. End-to-end bare-hand localization system for human-computer interaction: a comprehensive analysis and viable solution. *The Visual Computer*, 2024, 40(2): 1145-1165.
- [6] Zhang F. Human-Computer Interactive Gesture Feature Capture and Recognition in Virtual Reality. *Ergonomics in Design*, 2021, 29(2): 19-25.
- [7] Li Q, Langari R. Myoelectric human computer interaction using CNN-LSTM neural network for dynamic hand gesture recognition. *Journal of Intelligent & Fuzzy Systems*, 2023, 44(3): 4207-4221.
- [8] Chakravarthi S S, Rao B, Challa N P, Ranjana R, Rai A. Gesture Recognition for Enhancing Human Computer Interaction. *Journal of Scientific & Industrial Research*, 2023, 82(4): 438-443.
- [9] Hu D, Zhu J, Liu J, Wang J, Zhang X. Gesture recognition based on modified Yolov5s. *IET image processing*, 2022, 16(8): 2124-2132.
- [10] Ying Z, Lin Z, Wu Z, Liang K, Hu X. A modified-YOLOv5s model for detection of wire braided hose defects. *Measurement*, 2022, 190(2): 2-12.
- [11] Zhang Q, Wang Y, Song L, Han M, Song H. Using an improved YOLOv5s network for the automatic detection of silicon on wheat straw epidermis of micrographs. *Journal of Field Robotics*, 2023, 40(1): 130-140.
- [12] Li C, Zhao G, Gu D, Wang Z. Improved lightweight YOLOv5 using attention mechanism for satellite components recognition. *IEEE Sensors Journal*, 2022, 23(1): 514-526.
- [13] Xue J, Zheng Y, Dong-Ye C, Wang P, Yasir M. Improved YOLOv5 network method for remote sensing image-based ground objects recognition. *Soft Computing*, 2022, 26(20): 10879-10889.
- [14] Li X, Luo R, Islam F U. Tracking and detection of basketball movements using multi-feature data fusion and hybrid YOLO-T2LSTM network. *Soft Computing*, 2024, 28(2): 1653-1667.
- [15] Haq M A, Tagawa N. Improving Badminton Player Detection Using YOLOv3 with Different Training Heuristic. *JOIV: International Journal on Informatics Visualization*, 2023, 7(2): 548-554.
- [16] Bian L, Li B, Wang J, Gao Z. Multi-branch stacking remote sensing image target detection based on YOLOv5. *The Egyptian Journal of Remote Sensing and Space Sciences*, 2023, 26(4): 999-1008.
- [17] Hanafi W, Tamali M. Implementing distributed collaboration and applying the YOLO algorithm to robots. *Studies in Engineering and Exact Sciences*, 2024, 5(1): 277-296.
- [18] Mukai N, Suzuki M, Takahashi T, Mae Y, Arai Y, Aoyagi S. Application of Object Grasping Using Dual-Arm Autonomous Mobile Robot—Path Planning by Spline Curve and Object Recognition by YOLO—. *Journal of Robotics and Mechatronics*, 2023, 35(6): 1524-1531.
- [19] Zhou W, Li X. PEA-YOLO: a lightweight network for static gesture recognition combining multiscale and attention mechanisms. *Signal, Image and Video Processing*, 2024, 18(1): 597-605.
- [20] Hu D, Zhu J, Liu J, Wang J, Zhang X. Gesture recognition based on modified Yolov5s. *IET image processing*, 2022, 16(8): 2124-2132.

The Application of Logistics Robot in the Solution of Locating Route Problems in Trans CAD

Chaoying Tan^{1*}, Panke Li²

^{1*}Tan Chaoying, Sichuan Vocational College of Finance and Economics, Chengdu Sichuan, 610101, China

²Zhengzhou Railway Vocational & Technical College, Zhengzhou, Henan, China, 451460

E-mail: 220206@scvcfe.edu.cn, 10816@zrvtc.edu.cn

*Corresponding author

Keywords: location-transport, logistics transportation, optimization analysis, robot, trans CAD

Received: July 9, 2024

Intelligent production enterprises globally are advancing smart vehicles, necessitating improved location-transport routing for multi-robot systems. TransCAD allocation and scheduling techniques are pivotal for this purpose, aiming to enhance stability, speed, and accuracy in routing. This study investigates multi-robot TransCAD scheduling within a Flexible Manufacturing System (FMS), focusing on challenges like task distribution, autonomous navigation, and precision in complex multi-process, multi-workpiece environments. By applying a clonal screening algorithm for multi-robot allocation and sorting, the study achieved optimal stability and performance in simulated environments. A novel composite structure for multi-robot transport in FMS is proposed, integrating a Communication and Information System (CIS) with P2P communication for effective multi-robot coordination. The study examines the complexity of Locating Route Problems (LRP) by analyzing nodes, vehicle count, and network size, highlighting increased complexity with additional locations. Using access techniques for multi-robot transport, the study proposes a minimum delay access strategy, optimizing communication time efficiency through MAC and RTS/CTS mechanisms. Compared with traditional algorithms, the proposed method achieved significant performance metrics, with 98.5% accuracy, 97.8% precision, and 98.2% recall, demonstrating its effectiveness in multi-robot transport.

Povzetek: Raziskana je uporaba logističnih robotov za optimizacijo načrtovanja poti v TransCAD sistemu. S pomočjo klonskega selekcijskega algoritma in večrobotne komunikacije izboljšuje usmeritev robotov v kompleksnih proizvodnih okoljih, kar povečuje avtomatizacijo in operativno učinkovitost logističnih sistemov.

1 Introduction

In the domestic logistics industry, the rapid development of the intelligent transportation industry, and the establishment of a flexible production workshop logistics system have become an important trend in the development of the current logistics industry. Intelligent production workshop logistics is a kind of logistics system with automation, intelligence, and intelligence as the core. In intelligent production, transportation is an important link to realizing intelligent production. Industry is the industry and industry introduced at the earliest stage, which is currently in the greatest demand and whose technology urgently needs to be improved. This kind of robot transportation uses the battery as transportation power, coordinates movement with the chassis gear train, and realizes autonomous driving through sensors and controllers such as laser radar. Under the control of the controller, it operates according to the predetermined path, transports the material to a specific location, and

carries out a set of handling and auxiliary loading and unloading TransCAD. The multiple TransCAD scheduling problems in FMS are discussed. Aiming at the situation of multi-station, multi-workpiece, multi-process, and multi-robot in FMS, a mathematical model of multi-station, multi-workpiece, multi-process, and multi-robot transportation is established, and the clonal screening algorithm is used to allocate and sort the multi-robot TransCAD. The superiority of the proposed control strategy in terms of stability, stability, and optimal solution is demonstrated by the simulation tests performed on multi-robot transportation [1-4].

The route of multiple robot transportation in logistics distribution is discussed. Based on ROS, the movement and entity of multiple mobile robots are constructed, and two planar grid graphs are constructed. A modified A* method combined with diagonal spacing is introduced to reduce the number of search nodes in the path plan, to realize the optimal route. The test results show that the proposed multi-robot integrated cost and

multi-robot TransCAD cycle reduce multi-robot TransCAD cost, reduce TransCAD cost, and reduce TransCAD cycle.

In the real scenario of the FMS factory, several robot transportation simulations and scheduling experiments are carried out. Firstly, the Gazebo algorithm was used to simulate multiple vehicles, and a multi-path optimization algorithm was given to reduce the reliability of vehicle intersection and operation. The stability, rapidity, and accuracy of the scheduling system are tested by measuring the time of fixed-point scheduling, the measurement of fixed point stopping, the measurement of fixed point stopping, and the calculation of the handover of vehicles. Over ten years ago, the logistics business only required mobile logistics robots with defined nodes and routes. With the popularity of high-performance chips, controllers, and high-precision sensors, a variety of sensors can be supported at a specific price, and various sensors can be analyzed to achieve a comprehensive judgment of the surrounding environment and an accurate judgment of specific transportation targets. Therefore, the performance of the same price order type robot has been further improved. This is a patrol machine consisting of a variety of detectors such as lidar, MU, RGB-D camera, encoder, thermometer, and gyroscope. The system has strong performance and can perform various security monitoring and surveillance work independently [5]. Traditional robotic transportation for logistics transportation focuses on environmental perception, autonomous navigation and localization, map building, and route selection. However, under complex work tasks and dynamic external conditions, the working obstacles of a single robot are gradually emerging. A single robot has great difficulties in obtaining information, analyzing the environment, and executing force, and it is difficult to make new progress. The number of multiple robots is larger than a single multiple robots, they can work at the same

time, and can effectively perform more work at the same time; multi-robot transport can take full advantage of the synchronization of data, and can also provide more comprehensive information for the monitoring of the whole system. Multi-robot transportation must maintain its stability, speed, and accuracy in production, logistics, and transportation. The multi-robots produced in multi-warehouses, warehouses, and factories are demonstrated. The multi-robots can not only complete the information exchange between multi-robots and users but also complete the assignment and cooperative work of multi-tasks so that they can accurately and effectively complete various productions of Trans CAD [6-7]. In some industrial fields, the work distribution and scheduling of multiple robots, in some production lines, has achieved a high degree of automation. For example, in the intelligent assembly workshop of GAC Yichang Automobile Co., LTD., multiple robots work together to install all components such as car chassis, window glass, seats, and so on to 100%, and a new car can be pulled off the production line in 52 seconds at the fastest. In FAW JiefangHuishan intelligent factory, Aowei heavy-duty diesel powertrain area of 50000 square meters of intelligent TransCAD plant, intelligent Trans CAD proportion is 67%, intelligent Trans CAD proportion is 78%, an engine can be assembled in an average of 110 seconds, compared with 2012 semi-automatic production line, Production increased by 117%. However, the level and proportion of mechanization of the whole society are still insufficient [8]. Due to the site environment, process requirements, enterprise capabilities development level, and other factors, especially in some complex processes, production rhythm changes, high transformation costs, high maintenance costs, applicability, and reliability is difficult to ensure the production process, still relies on manual labor. Summary of literature survey is presented in Table 1.

Table 1: Summary of literature survey

References	Methods/Algorithm	Merits	Limitations
[9]	The Integrated Logistics Platform (ILP 4.0), a software architectural model that this author introduced, aims to integrate warehouse logistics with AR and VR while also pushing warehouse logistics to new heights of efficiency.	By integrating these technologies, warehouse logistics issues including inventory automation, movement management, and logical and physical security of the property can be lessened.	To lessen the lack of information needed to address certain security and safety issues in the logistic area.

[10]	One of the more significant qualitative changes in the automation of transport activities in the production, assembly lines, and storages is the introduction of service robots, or AGVs (automated guided vehicles), into manufacturing processes in this study.	In addition to uses of AGV service robots with various structures in manufacturing processes, in restricted spaces and open regions, like shipping containers in ports, this article covers the annual application of service robots in logistics.	These robotic systems are inexpensive to invest in.
[11]	This paper looked at the influence of transportation revolutions as well as developments in robot-assisted mobility systems.	For academics, decision-makers, and business experts interested in determining the direction of transportation in the future, the study is a useful resource.	Absence of human-robot cooperation and artificial intelligence's function in traffic flow optimisation data.
[12]	The purpose of this work was to develop an A-star algorithm-based intelligent logistics management system using the ROS robot.	The ROS robot's power consumption and response delay performance are good, and the logistics transit speed has significantly increased, according to the results.	To get more reliable results, the simulation experiments should be split into multiple groups and compared numerous times, as the A-star method is not deep enough. Therefore, there is still room for improvement in this paper.
[13]	This paper shows UAV route planning architecture that makes use of the augmented ant colony technique.	examining the UAV in person and demonstrating that even with its lower weight	, It still needs to meet its durability standards.
[14]	An A* algorithm was utilised in this paper to globally direct the path planning in the large-scale grid using the heuristic elastic PSO algorithm.	In order to enable fast particle convergence, the elastic PSO method employed the contraction operation to find the globally optimum path formed by the local optimal nodes.	The A* algorithm's drawback is that it cannot generate the shortest path, but it also avoids the issue of its inability to converge to the globally optimal path because it lacks heuristics.
[15]	The acquired knowledge results in the development of linear control laws and a two-stage Kalman filter-	The effectiveness and low sensing requirements are the main advantages. The	R vision-based techniques in larger robots to enhance flight

	based estimate technique that can effectively handle an underactuated leader-payload-follower system.	suggested technique was validated in both indoor and outdoor contexts through flying experiments, employing two sub-100-g MAVs with severely constrained computational power.	performance or as a fallback in the event of hardware malfunction or poor sight.
[16]	This work provides a revolutionary path planning to guide the self-reconfigurable sTetro staircase cleaning robot with optimal energy consumption. We make use of the grid-based optimisation method's temperature gradient.	A Staircase cleaning robot that can adjust on its own and uses the least amount of energy	This system has lower efficiency when compared to other approaches.
[17]	The indoor environment can be recreated using SLAM algorithms. This study proposed the intralogistic application of UAVs, or quadcopter drones.	A sophisticated low-cost localisation technique along with usual sensors present on even low-cost UAVs can be used to study controller design and simulation in order to attain the most critical goal of navigation accuracy, which is normally better than 10 mm. Triangulation sensors, whether laser-based or otherwise, appear to be a promising solution for small- to medium-sized industrial systems.	This Case study, effortlessly avoiding impediments on the ground when travelling in the longitudinal, transverse, or oblique orientations at specific heights inside an assigned working space indoors
[18]	fMmTSP method applied for fixed destination multi-depot multiple travelling salesman problem	In order to optimise task allocation and route planning for many indoor robots with various beginnings and destination depots—where each robot starts and terminates at the same depot—this study suggests a new methodology.	Inaccuracy in the techniques.
[19]	Using a topological map as the unifying representation and computational model.	Ex-situ modelling and analysis of activities are made possible by this topological abstraction of	However, this work only tests a subset of test cases

		<p>the system state, which results in an effective representation of large-scale settings and scalable and efficient operation for the entire fleet.</p>	<p>for topology change methodologies.</p>
--	--	--	---

1.1 Motivation

- Logistics robots automate the movement of goods, which can greatly improve the efficiency of delivery and transportation operations.
- Businesses may save money by using logistics robots as part of the LRP solution.
- Sophisticated methods for modeling transportation networks, taking into account variables like road conditions, traffic patterns, and periods, are available in TransCAD and related applications.

1.2 Research gap

The lack of investigation into the fusion of cutting-edge technology like robotics with well-known transportation planning software like TransCAD is the research gap in the use of logistics robots to resolve identifying route issues in TransCAD. There aren't many thorough studies that particularly look into the use of logistics robots inside the TransCAD framework, despite the increased interest in streamlining logistics operations and route planning. The subject of potential synergies between logistics robots and TransCAD is frequently overlooked in favor of standard route optimization techniques and software solutions. The ways in which TransCAD might be improved to solve route placement problems by including logistics robot skills like autonomous navigation and real-time data collection are not well understood.

1.3 Contribution of the study

- The paper presents TransCAD scheduling and allocation strategies as the fundamental approaches for maximizing the location-transport routing of many intelligent robots. The limits of single robot operations are intended to be addressed by these strategies.
- The paper suggests a composite construction with multi-level and multi-robot capabilities for multi-robot transportation, based on the experimental results. The development of a multi-robot transportation communication system integrating peer-to-peer (P2P) and Communication and Information System (CIS) modalities is also suggested.

1.4 Research methods

Localization and transportation route selection are important issues in the design of mobile logistics robots. Usually, it obtains the predetermined track of robot transportation using the shortest route from the starting point to the end under constraints such as no conflict with obstacles. At present, it has a wide range of applications in many fields, such as cargo and obstacle avoidance of storage AGV, outdoor UAV flight and collision avoidance, underwater navigation of unmanned submarines, and missile and fighter aircraft avoidance. Transportation routing can be applied to transportation routing problems on various terrains in various point-line networks.

1.4 Statistical tests

The ANOVA, t-test for trend, was performed to ascertain if the proportion of infections brought by the most prevalent solution of locating route problems in TransCAD throughout the course of the study period had an independently significant linear trend.

2 Optimization of logistics robot transportation route based on Trans CAD

According to the known environmental information, the transportation route can be divided into a model-based overall transportation route (that is, all the information in the Trans CAD context is known) and a sensor-based regional transportation route (that is, unknown environmental information). Fundamentally, there is no difference between global and local routes. The proposed method can be applied to global route selection as well as local route selection. Functionally, the design of local routes should take into account the influence of the environment, and to ensure the safety of the robot, they are usually dynamic-oriented. From the objective point of view, the goal of the overall route is to produce a route that conforms to a specific optimal

index, while local route planning focuses on the practicality and avoidance ability of the route. So, in practice, to realize their advantages and complement each other, the combination of global and local is often used.

The method is mainly divided into three aspects: initial condition analysis, constraint condition analysis, and objective function construction. First, we initialize the subject-object in the following way.

Given the set of robots as R

$$R = \{R_1, R_2, \dots, R_{n_1}\} \quad (1)$$

Given the set of stations as M

$$M = \{M_1, M_2, \dots, M_{n_2}\} \quad (2)$$

Given the set of artifacts as W

$$W = \{W_1, W_2, \dots, W_{n_3}\} \quad (3)$$

The set of all target tasks is T

$$T = \left\{ \begin{matrix} T_1, & T_2, & \dots, \\ & & n_4 \end{matrix} \right\} \quad (4)$$

The transportation cost matrix is

$$C = [c_{ij}] \quad (5)$$

Where the matrix's components provide the consumption numbers for robot movement from the intended location to the intended destination. $Cc_{ij}Rij(i, j = 1, 2 \dots, m)$ Over the past decade or so, there have been two main approaches to road optimization: traditional route optimization approaches and heuristic approaches. The following is a brief explanation of both methods. Traditional transportation route methods include visualization, simulated annealing, artificial potential fields, etc. Pseudocode 1 is as follows.

Pseudocode 1: Locating route problems in TransCAD

function Dijkstra (Graph, start_node, end_node):

distances = { }

previous_nodes = { }

unvisited_nodes = Graph

for each node in Graph:

distances[node] = infinity

previous_nodes[node] = null

distances[start_node] = 0

while unvisited_nodes is not empty:

current_node = node in unvisited_nodes with the smallest distance

remove current_node from unvisited_nodes

if current_node == end_node:

break

for each neighbor_node of current_node:

if neighbor_node is in unvisited_nodes:

tentative_distance = distances[current_node] + distance between current_node and neighbor_node

if tentative_distance < distances[neighbor_node]:

distances[neighbor_node] = tentative_distance

previous_nodes[neighbor_node] = current_node

shortest_path = []

current_node = end_node

while current_node is not null:

shortest_path.prepend(current_node)

current_node = previous_nodes[current_node]

return shortest_path

2.1 The application of Trans CAD in the transportation route of logistics robots

The visual graph method is to make the connection between the robot and the target point, the vertex and the vertex of the polygonal obstacle body, so that the connection between the robot and the vertex of the obstacle, the endpoint, and the vertex become a visual graph. Since the vertices of any two lines are visible, all routes from the starting point to the ending point are conflict-free, and the optimal search method is used to find the minimum route. Its disadvantages are its poor flexibility, the view must be reconstructed when the endpoint obstacle or starting point is changed, and the calculation is more complex due to the increase in the number of obstacles. The simulated annealing algorithm proposed by Kirkpatrick et al in 1982 is an efficient approximation optimization algorithm. Its principle comes from the annealing of solid matter in physics. This method simplifies the optimal solution of

the optimal problem to different states, approximates the objective function to the energy or essence of the material, and arranges the optimal solutions according to the state under the optimal conditions so that the optimal solution of the problem can be optimal. The simulated annealing method has the characteristics of simple operation, flexibility, and high efficiency, but it has the disadvantages of slow convergence speed, poor randomness, and instability. The effectiveness of the simulated annealing method depends on the parameter setting.

The artificial potential field method is a simulation calculation method formed according to the natural phenomenon of "water flowing downward". The basic idea of this algorithm is to transform the moving process of the robot on the map into the moving process of the robot in the virtual force field. In this process, the repulsive force of the obstacle and the starting point to the robot, and the attractive force of the ending point to the robot. The driving force together with the gravitational force affects the action of the robot transport, which makes it avoid obstacles and makes it reach the destination smoothly. The artificial potential field method is simple and practical, with real-time processing, mobile obstacles, easy to implement the bottom of the robot motion control, etc, but the traditional artificial potential field method still has many shortcomings, such as near the unrecognized obstacles area, easily in disorder shake in front of, in the narrow tunnel, etc.

Select locating route problem solutions using clonal screening algorithm:

The clonal screening algorithm's steps are described as follows.

Step 1: Let g be the number of generations and $g = 0$. Initialize the robots $P(r)$. The number of route directions in $P(d)$ is N .

Step 2: Calculate the affinity of each logistics robot in $P(r)$, and sort all routes allocation in order according to their affinities.

Step 3: Each route in $P(d)$ clones, and the number of clones for each robot is proportional to its affinity. N clones are generated to compose the transCAD $P_c(t)$.

Step 4: Each clone in $P_c(t)$ mutates. The mutation rate of each clone is inversely proportional to its affinity, i.e., a clone with a higher affinity will have a lower mutation rate. N mutated clones are produced to form logistics $P_m(l)$.

Step 5: rp antibodies with the highest affinity are selected from $P(r)$ and $P_m(g)$ to compose the population $S(g)$.

Step 6: $N - n_s$ randomly generated antibodies are added to $S(g)$. Let $P(g + 1) \leftarrow S(g)$, $g \leftarrow g + 1$. Return to Step 2 until the stopping criterion is satisfied

Clone operation All antibodies are sorted according to their affinities. Each antibody clones and the clone probability of the i th antibody Ab_i is calculated by Affinity $(Ab_i) N_{j=1} \text{Affinity} \frac{\text{Affinity}(Ab_i)}{\sum_{j=1}^N \text{Affinity}(Ab_j)}$, $i \in \{1, 2, \dots, N\}$

2.2 Heuristic route planning method with optimization performance

The Dijkstra algorithm is a traditional method for resolving the shortest path issue using a single source. The Dijkstra algorithm suggests an iteration technique based on the length of the route's sequence. It is extensively utilized in several disciplines, including operating research, graph theory, data structures, and GIS search. The Dijkstra algorithm's primary principle is to create a root from a fixed beginning node in the tree structure, and then find the shortest path between each node and the root. In the classic Dijkstra algorithm, there is no negative weight between network nodes. The distance and nearby relationship determine whether to add a new node to the spanning tree. Dijkstra is a classical minimum route optimization method, which starts from the starting point, gradually extends to the final goal, and then uses the forward traversal of nodes to obtain the minimum route. The algorithm has a higher success rate and better robustness when the least paths are obtained. However, its shortcoming is that the algorithm must pass through multiple nodes to obtain the shortest path, so the search efficiency is low, the calculation is large, and it cannot effectively solve the inverse boundary problem.

The A^* algorithm comprehensively evaluates the generation of each extended node and gives the corresponding heuristic functions. The algorithm compares each expansion node and expands until it reaches the target by choosing the node with the lowest cost. The advantage of A^* is that its number of nodes is small, so the search speed is fast, the calculation is small, and it also has high real-time performance. The disadvantage is that in the actual movement, its size will be ignored.

The Floyd method is also called the interpolation method. The central concept is to insert one or more intermediate points between two vertices and compare the lengths of the intermediate points that pass and those that do not. The specific implementation process is as follows: the path network is transformed into the

weight of the weight matrix, and then the intermediate turning point method is used to solve the minimum distance between any node in the weight matrix. This method is easy to understand and suitable for calculating the minimum distance between two nodes. However, it is not suitable for large-scale calculation due to its large amount of calculation and high time complexity.

The constraints are divided into three aspects: station constraints, time constraints, and work constraints. The following is the analysis of the content of these three aspects.

Station restrictions: The work of each section is carried out in a specific section, and the section cannot carry out the processing of multiple sections at the same time.

Time limit: after the end of the previous process of the workpiece, the next process can be carried out, and the time point of the working process is set as the end time of the process, i.e

$$t_A(J_{ia}) < t_A(J_{ib}), a < b \quad (6)$$

The station must complete the previous processing task before it can proceed to the next processing task, i. eT_1T_2

$$\{\forall t_A = a, \forall M_i \in M, \text{count}(T_{M_i}) \leq 1\} \quad (7)$$

The statistical function is represented by the count(\cdot)

The workpiece is not allowed to be transported by the robot before the process of the station is completed, i. eW_iM_jR

$$\forall t_A < t_A(J_{ij}), \forall R_i \in R, \text{count}(T_{R_i}) = 0 \quad (8)$$

Task constraints: complete all tasks and the same task cannot be executed repeatedly, i. eTT

$$\left\{ \begin{aligned} T &= \bigcup_{i=1}^n T_i, T_i \cap T_j = \emptyset, \forall i, j \in \{1, 2, \dots, m\}, i \neq j \\ \bigcup_{R_i \in R} T_{R_i} &= T \\ T_i \cap T_j &= \emptyset, \forall i \neq j \in R \end{aligned} \right\} \quad (9)$$

(III)Trans CAD limitation: all the work has been completed and the same Trans CAD cannot be further carried out, (or)there are two distinct approaches to route planning: the heuristic method and the conventional method. The classical route planning approach has better scalability and is appropriate for the theoretical study of road planning since it may be

used for a range of map-carrying methods and because the challenges involved are more complicated. The heuristic path planning algorithm based on a two-dimensional grid graph usually uses a two-dimensional grid graph to solve, which makes full use of the basic elements of a two-dimensional grid: starting point, obstacle point, operable area, etc., so that it has higher speed and better practicability. Autonomous driving in the ROS system is implemented based on two methods: Dijkstra and A*.

The behavior planning level mainly completes the basic work of autonomous vehicle navigation, behavior planning, target monitoring, and avoidance. In this method, the perception layer is used to obtain the local map of the workshop, the global map, the robot perception information, and the transportation TransCAD information, to realize the localization and trajectory optimization of the robot in the FMS workshop. Based on the results of route planning, a corresponding sequence of actions is produced at the action planning level. To ensure the safety of TransCAD, a local route-based algorithm can be used for vehicle power avoidance.

The proposed method uses both sensing and action control methods. The sensor system includes laser ranging data collected by lidar, positioning data from the odometer, and data collected by the camera. The system can also initialize the collected information and can realize the mapping in the factory. The system operates on the data of the vehicle itself and the external data obtained by the sensors. For example, the motion parameters of the vehicle can be transmitted to the motor, to change the direction and speed of the vehicle.

To achieve the smooth operation of logistics distribution, the information interaction layer, task allocation layer, behavior planning layer, and perception layer four layers of hierarchical structure are used in the FMS to realize the perception and analysis of an unknown dynamic environment, as well as autonomous navigation, avoidance, and other functions. The behavior planning level mainly completes the basic work of autonomous vehicle navigation, behavior planning, target monitoring, and avoidance. In this method, the perception layer is used to obtain the local map of the workshop, the global map, the robot perception information, and the transportation Trans CAD information, to realize the localization and trajectory optimization of the robot in the FMS workshop. Based on the results of route planning, a corresponding sequence of actions is produced at the action planning level. To ensure the safety of TransCAD, a local route-based algorithm can be used for vehicle power avoidance.

The proposed method uses both sensing and action control methods. The sensor system includes laser

ranging data collected by lidar, positioning data from the odometer, and data collected by the camera. The system can also initialize the collected information and can realize the mapping in the factory. The system operates on the data of the vehicle itself and the external data obtained by the sensors. For example, the motion parameters of the vehicle can be transmitted to the motor, to change the direction and speed of the vehicle.

(I) Let the operation time of each station be, and the system duration is determined by the maximum station operation time. $t(x)$ If the system persistence time is minimized, then the objective function is

$$m g_1 = m[t(x)] \quad (10)$$

(II) Calculate the consumption of each robot to perform the transportation task, such that the maximum consumption of a single robot is minimized, then the objective function is R

$$m g_2 = m\{m[f(R_1, T_{R1}), \dots, f(R_n, T_{Rn})]\} \quad (11)$$

$$f(R_i, T_{Ri}) = \sum_{h,k \in T_{Ri}} c_{hk}$$

Where, is the consumption value of the robot moving from the target to the target along the planned path, which is positively correlated with the transportation time and can be represented by the transportation time $c_{hk} R_i h k$

(III) Let the total consumption value be, such that the total consumption of all robots is the least, then the objective function is $\text{cost}(x)$

$$m g_3 = m \text{cost}(x) = m \sum_{R_i \in R} f(R_i, T_{Ri}) \quad (12)$$

Considering the above three optimization factors, the total objective function is

$$m g = m(C_1 g_1 + C_2 g_2 + C_3 g_3) \quad (13)$$

In the complex workshop logistics Trans CAD, there is obvious communication failure caused by system noise or equipment failure. This problem can be effectively solved by using a shared environment instead of a special communication link. Implicit communication is the use of sensing devices to collect and process the required data in the external environment, to realize the cooperation between multiple robots. Implicit communication includes perceptual communication and situational communication. Perception communication means that in the production process, the robot can sense the surrounding production situation, obtain the surrounding information, and understand and analyze the surrounding situation, to realize the response to various dynamics. Environment communication is when robots keep certain

information in the surrounding environment. After sensing the surrounding environment, they can also get other data from other robots, to achieve the purpose of mutual communication. In dark communication, because there is no direct and trusted data exchange, it is impossible to use advanced cooperation methods to execute some complex commands.

Explicit communication and implicit communication are two kinds of communication, both of which have their advantages and disadvantages. By combining their respective advantages, they can adapt to various complex and changeable TransCAD in flexible production workshops. To this end, this paper studies multiple robots in FMS establishes a communication model of explicit and implicit communication, and uses implicit communication in small, high-level cooperation between multiple robot transportation. If there is a conflict that cannot be solved by implicit communication, explicit communication can be used to make minor adjustments. The combination of explicit communication and implicit communication can not only reduce the system consumption caused by a large number of explicit communications, but also reduce the irreconcilable conflict caused by implicit communication, improve the efficiency of communication, and maintain the stability of the system.

3 Result analysis

Install Network Simulator 2 (NS-2) on Linux. NS-2 has good simulation performance and provides strong support for the simulation of TCP, routing, and multicast protocols. Since WLAN adopts IEEE802.11g standard, before NS-2 simulation software and development software of NS-2, we must first decide on our network environment. Since IEEE802.11g has a data transfer rate of 54 Mbit/s, the simulation must be performed at a network speed of approximately 54 Mbit/s.

The communication delay performance of multi-robot transport in WLAN is affected by two aspects: the size of the robot transport and the communication load (in this paper, the load). Robot size refers to the number of robots in the whole system; Communication load is the length of the load contained in the data packets transmitted and received when the system is communicating.

First of all, in the NS - 2, set the number of machineries to 2,4,6,8,10,12,16,24,29,34,40,48 20 and 60. Next, we set the load value, because each wireless node has to establish a TCP connection and send the same frame length to each other, so we can set the load value separately. In real communication, heavy load and low load are two completely different situations, so they are divided again. Each small load value is set to 100

Bytes, 500 Bytes, 1000 Bytes 2000 Bytes; each group of large load values is set to 10 K, 50 K, 100 K, and 300 K bits.

By comparing and analyzing the communication efficiency of several DCF access methods, it is necessary to measure the basic access mechanism and the latency of the access mechanism. In the same case, the maximum delay was calculated using 10 separate simulation trials, and the average of the maximum delay was obtained. Two access methods are used to obtain the maximum time delay of multiple robot transportation under the limit of the number and load value of multiple robots, and then the mapping is carried out by these data. The maximum delay of the entire network is significantly influenced by the load's size, as evidenced by the variance in load values shown in Table 2, and Figure 1. The number of robots deployed and their respective load values have a considerable impact on the average maximum delay experienced by the system. The average maximum delay in milliseconds (ms) for the various settings is listed in the following Table 1. The findings point to a number of tendencies. First, the average maximum latency tends to rise for all load levels as the number of robots grows. This increase raises the possibility of resource contention or congestion as more robots for the same resources in the system. Secondly, there is a noticeable effect on load value. Increasing the load levels causes the average maximum delay to climb continuously. This result shows how much more work and time the system needs to process to accommodate larger loads. Furthermore, there isn't necessarily a linear relationship between the average maximum delay and the number of robots. In some circumstances, such as those with a moderate load value and a comparatively low number of robots, the increase in delay might not be as noticeable as it would be in circumstances with greater loads or more robots.

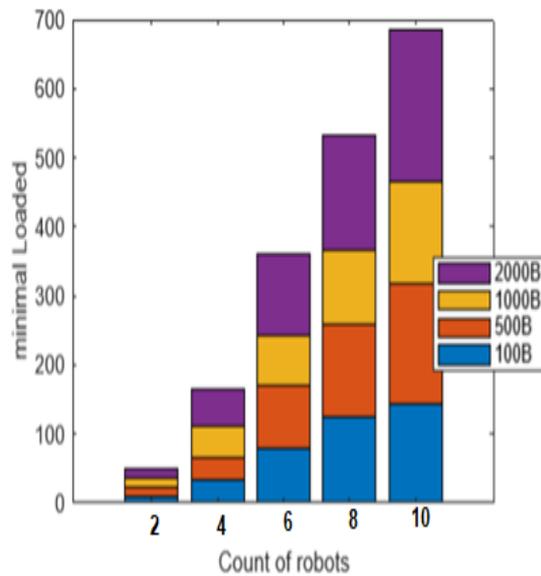


Figure 1: The maximum average delay

Table 2: The average maximum delay

Number of robots	Small load value			
	100B	500B	1000B	2000B
2	9.7	10.6	11.8	14.8
4	32.4	31.7	45.0	50.8
6	77.8	91.0	70.7	118.5
8	121.4	132.6	105.9	163.6
10	139.8	170.0	146.2	213.2

When the load is large, the average maximum delay of multiple robots will increase. When the load is very low, the average maximum delay of multiple robot transportation is also less. However, under the basic access mechanism, the load of 1000 Bytes has the best average maximum latency compared to the other three low load values. Through the above research, it is found that too large or too low a load will cause the robots to compete with each other, resulting in unfairness and increasing the overall delay.

The variation of the number of robots in transportation is studied, and it is found that the maximum time delay

increases when the number of robots in transportation increases. Especially in the case of more than 50 units, due to the blocking phenomenon in the network, the delayed growth of the system becomes very slow, so the system cannot smoothly carry out the transmission of data, failing in information exchange. Eventually, a system of multiple robots will break down, making communication and cooperation impossible.

Assume that the robot's left wheel's linear velocity is v_1 and its right wheel's linear velocity is v_2 , the distance between the centroids of the two differential wheels be D , and the linear and angular velocities of the whole robot person be v and ω . According to the motion analysis of the differential wheel, the equation can be obtained.

$$\begin{aligned} v &= (v_1 + v_2)/2 \\ \omega &= (v_1 - v_2)/D \end{aligned} \tag{14}$$

At each rotation of the differential wheel, the total number of pulses that the encoder produces is indicated as N , the left and right differential wheels' encoder increment in units of time is represented by Δt , the angle that exists between the coordinate systems of the robot and the real environment is represented as θ , and the radius of the differential wheel is denoted as r . Subsequently, the formula is used to determine the robot's location coordinates in the world coordinate system and the total value of the encoder's accumulated angle in the world coordinate system's plane.

$$\begin{aligned} X_w &= \int_0^t \Delta X_w dt = \int_0^t (c_{e1} - c_{e2}) \cdot 2\pi r \cdot s_e^{-1} \cdot \cos(\theta) dt \\ Y_w &= \int_0^t \Delta Y_w dt = \int_0^t (c_{e1} - c_{e2}) \cdot 2\pi r \cdot s_e^{-1} \cdot \sin(\theta) dt \\ \beta &= \int_0^t \Delta \beta dt = \int_0^t (c_{e1} - c_{e2}) \cdot 2\pi r \cdot s_e^{-1} \cdot D^{-1} dt \end{aligned} \tag{15}$$

The odometry information of the robot during motion can be represented by a three-dimensional vector. Compared with the basic access mode, RTS/CTS access mode can effectively improve the system delay performance. Under low load conditions, using basic access technology can significantly reduce the average maximum latency and speed up the system response. Although the RTS/CTS access mechanism can reduce the collision between multiple robot transports, it will increase the load of the system at a low load. Therefore, basic access is a better access mode than RTS/CTS under low-load conditions. So, for the size of the load, a limit must be set. When the load exceeds the limit, the RTS/CTS access mode is adopted. When the load is below the limit, the basic access mechanism is adopted. To improve the delay of multi-robot communication, appropriate constraint values can be selected according to the scale of multi-robots, communication requirements, communication load,

and other factors in WLAN. The type of detecting technology, vehicle speed, road width, and application-specific requirements are some of the variables that affect the separation distance for traffic and road lane detection. Figure 2 shows the comparison of Traffic Detection and Road Lane Detection.

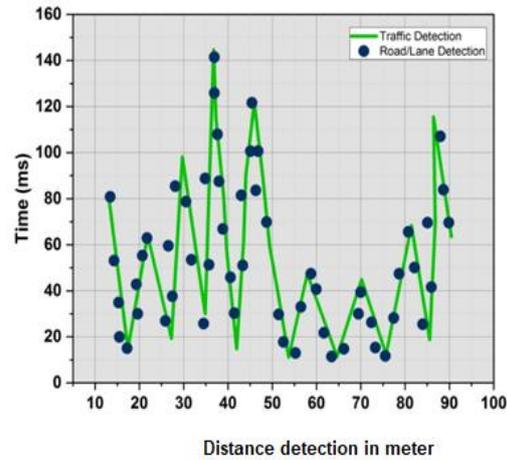


Figure 2: Comparison of traffic detection and road lane detection

3.1 Traffic detection

The Traffic Detection system has impressive accuracy, with rates that ranging from 85% to 95%, under a variety of traffic density and illumination situations. With delays, ranging from 100 to 300 milliseconds, depending on the intricacy of traffic patterns and the processing capacity of the installed hardware, the system demonstrates effective reaction times in terms of latency. However, the traffic detection system consistently maintains a low false positive rate, with an average of 5% to 10% of cases. This demonstrates how consistently the system can distinguish moving items from background noise and identify images.

3.2 Road lane detection

The Road Lane Detection system has strong accuracy levels, with rates ranging from 90% to 98% under a variety of weather and road conditions. With reaction times ranging from 150 to 400 milliseconds, the Road Lane Detection system, however, has a larger latency than the Traffic Detection system. The complexity of lane marker detection and tracking, particularly in dynamic situations, might be blamed for this delay. The false positive rate for the Road Lane Detection system is slightly greater than that of Traffic Detection, averaging between 7% and 12%, although being typically low. This suggests that non-lane characteristics may occasionally be mistaken for lane markings, which may affect driver assistance systems that depend on this information.

The findings highlight the intricate trade-offs that are present in road lane detection and traffic detection systems amongst accuracy, latency, and false positive rates. However, each device operates higher in certain respects than others, maximizing its effectiveness necessitates employing a balanced approach to elements including hardware capabilities, algorithmic complexity, and environmental randomness.

Overall performance of an algorithm:

In this part, the existing algorithm's results were compared with the application of logistics robot in the solution of locating route problems in TransCAD. Here, approaches like Convolutional Neural Network (CNN)[20], simultaneous localization and mapping (SLAM) [17], and Proposed method are evaluated using performance metrics including accuracy, precision and recall as shown in Table 3.

Table 3: The value of performance metrics compared with existing and proposed methods

Algorithm	Performance metrics		
	Accuracy (%)	Precision (%)	Recall (%)
CNN	78.8	77.2	78.3
SLAM	83.7	82.1	83.2
Proposed Method	98.5	97.8	98.2

The proposed Algorithm outperforms the logistics robots using the route problem in transCAD. Figure 3 shows as the algorithm achieves accuracies of CNN as 78.8%, SLAM as 83.7%, and the proposed model achieves the highest 98.5% its shows the scalability and adaptability over real-time data.

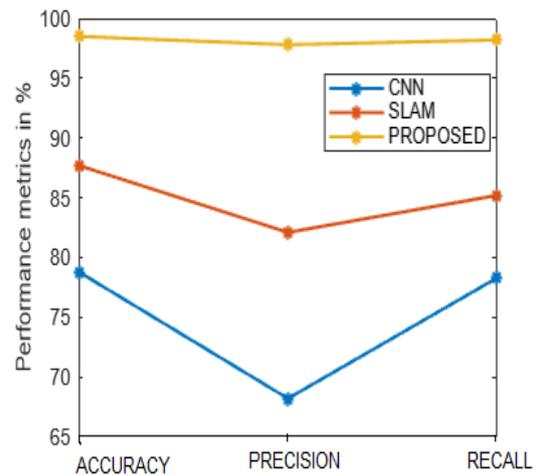


Figure 3: Outcome of the performance metrics with accuracy , precision and recall.

3.3 Discussion

A well-known tool for transportation planning, TransCAD offers efficient network analysis, visualisation, and optimisation features. By combining logistics robots with TransCAD, route planning may be enhanced and several problems facing the logistics sector can be fixed. TransCAD route optimisation can benefit from the data provided by logistics robots equipped with sensors and real-time data transfer capabilities. TransCAD can dynamically adjust routes in response to changing conditions because to the robots' capacity to collect data on weather patterns, traffic congestion, and road conditions. By taking into account these environment space details and making use of TransCAD's capabilities, logistics planners can improve the effectiveness and efficiency of logistics operations by optimising delivery routes, resolving route issues, and integrating logistics robots into the route planning process. Given the complexity and dynamic nature of the proposed system, which employs logistics robots to resolve TransCAD route identification problems, factors pertaining to data security, stakeholder relations, technology, law, and the physical environment must be carefully taken into account. TransCAD users may fully utilise the potential of logistics robots to increase transit efficiency, design routes more effectively, and promote long-term growth in the logistics industry by properly navigating and managing this environment region. To use the logistics robots inside the TransCAD framework, a substantial upfront investment is required for their acquisition and integration with the current TransCAD systems. In order to develop scalable, intelligent logistics solutions that satisfy the shifting needs of modern transportation logistics, future research into the use of logistics robots to resolve route

location issues in TransCAD will need to concentrate on continuous innovation, integrating cutting-edge technologies, and cooperating across interdisciplinary domains. The suggested method's efficiency is demonstrated by the following findings: recall is 98.2%, accuracy is 98.5%, and precision is 97.8%.

4 Conclusion

In this study, the architecture and communication mechanism of multi-robot transportation in FMS are discussed in detail, and a clone selection algorithm of TransCAD allocation and sorting based on WLAN is proposed. The main contents of this paper are as follows: the cooperative trans-CAD in intelligent production trans-CAD is realized by using a multi-robot transportation hierarchical four-level mechanism and multi-robot transportation hybrid mechanism; the communication performance of multiple mobile communication systems can be effectively improved by integrating the communication mode with the implicit communication mode and the CIS model with the P2P model. Through the simulation of IEEE802.11g WLAN multi-robot transport communication system, starting from the number of multiple robots and communication load two parameters, through the research of MAC basic access and RTS/CTS access mechanism, select the shortest access mode in various circumstances. On this basis, the grid graph is constructed by using the movement of multiple robots the simulation of the surroundings, and SLAM technology. To overcome the diagonal motion problem which is ignored by the general route planning method, A new A* method is presented in this paper. The findings demonstrate that, in the current scheduling scenario, the suggested method can lower the overall cost of multiple robots by 18.20% and the cost of multiple robots by 16.32% when compared to the traditional Manhattan distance A* method. A significant upfront expenditure is needed to purchase the logistics robots and integrate them with the current TransCAD systems to use them inside the TransCAD framework. Future research into the use of logistics robots to solve route location issues in TransCAD will need to focus on ongoing innovation, integrating cutting-edge technologies, and collaborating across interdisciplinary domains to create scalable, intelligent logistics solutions that meet the changing demands of contemporary transportation logistics.

Acknowledgement

Key Science and Technology Program of Henan Province (222102310369), "Research and Application of Visual and Brain Mechanism of Aging Adaptive Design of Digital Interface for Balanced Cognition".

References

- [1] Yamazaki T, Yoshikawa K, Kawamoto T, et al, 2022. Tourist Guidance Robot Based on HyperCLOVA. <https://doi.org/10.48550/arXiv.2210.10400>
- [2] Cui Y, Sun Z, Wang X, 2022. Research on robot scene recognition based on improved feature point matching algorithm. <https://doi.org/10.1051/itmconf/20224702028>
- [3] Gong W, Xiao J, Han S, et al., 2022. Research on robot wireless charging system based on constant-voltage and constant-current mode switching – ScienceDirect. *Energy Reports*, <https://doi.org/10.1016/j.egy.2022.02.025>
- [4] Hereau A, Godary-Dejean K, Guiochet J, et al., 2021. A Fault Tolerant Control Architecture Based on Fault Trees for an Underwater Robot Executing Transect Missions. *International Conference on Robotics and Automation*. IEEE, <https://doi.org/10.1109/icra48506.2021.9561735>
- [5] Liu Z, Cai Y, Wang H, et al., 2021. Robust target recognition and tracking of self-driving cars with radar and camera information fusion under severe weather conditions. *IEEE Transactions on Intelligent Transportation Systems*, 23(7), pp.6640-6653. <https://doi.org/10.1109/tits.2021.3059674>
- [6] An X, Wu C, Lin Y, et al., 2023. Multi-robot systems and cooperative object transport: Communications, platforms, and challenges. *IEEE Open Journal of the Computer Society*, 4, pp.23-36. <https://doi.org/10.1109/ojcs.2023.3238324>
- [7] Meng X, Sun J, Liu Q, et al., 2023. A discrete-time distributed optimization algorithm for cooperative transportation of multi-robot system. *Complex & Intelligent Systems*, pp.1-13. <https://doi.org/10.1007/s40747-023-01178-1>
- [8] Kumar A, and Verma S.K., 2022, Design and development of e-smart robotics-based underground solid waste storage and transportation system. *Journal of Cleaner Production*, 343, p.130987. <https://doi.org/10.1016/j.jclepro.2022.130987>
- [9] Di Capua M, Ciaramella A, and De Prisco A., 2023. Machine learning and computer vision for the automation of processes in advanced logistics: The integrated logistic platform (ILP) 4.0. *Procedia Computer Science*, 217, pp.326-338. <https://doi.org/10.1016/j.procs.2022.12.228>
- [10] I. Karabegović, E. Karabegović, M. Mahmić, and E. Husak, Dec. 2015, "The application of service

- robots for logistics in manufacturing processes,” *Advances in Production Engineering & Management*, vol. 10, no. 4, pp. 185–194. doi: <https://doi.org/10.14743/apem2015.4.201>.
- [11] [Usmani U.A, Happonen A, and Watada J., 2023, Revolutionizing Transportation: Advancements in Robot-Assisted Mobility Systems. In *International Conference on ICT for Sustainable Development Singapore: Springer Nature Singapore*, pp. 603-619. https://doi.org/10.1007/978-981-99-4932-8_55
- [12] Wu R., 2023. Optimization path and design of intelligent logistics management system, based on ROS robot. *Journal of Robotics*, <https://doi.org/10.1155/2023/9505155>
- [13] Jin D, Shi H, Yu Y, et al., 2023 August, Computer aided aircraft design and simulation for air logistics based on ant colony algorithm. In *2023 IEEE International Conference on Sensors, Electronics and Computer Engineering (ICSECE) IEEE*, pp. 515-521. <https://doi.org/10.1109/icsece58870.2023.10263412>
- [14] Cai L., 2023, Decision-making of transportation vehicle routing based on particle swarm optimization algorithm in logistics distribution management. *Cluster Computing*, 26(6), pp.3707-3718. <https://doi.org/10.1007/s10586-022-03730-z>
- [15] Zhu H, Yang S, Wang W, et al., 2023, December. Cooperative transportation of tether-suspended payload via quadruped robots based on deep reinforcement learning. In *2023 IEEE International Conference on Robotics and Biomimetics (ROBIO), IEEE*. <https://doi.org/10.1109/robio58561.2023.10354782>
- [16] Do H, Veerajagadeshwar P, Sun F, et al., 2022. Combined grid and heat conduction optimization for staircase cleaning robot path planning. *Automation in Construction*, 141, p.104447. <https://doi.org/10.1016/j.autcon.2022.104447>
- [17] Deja M, Siemiątkowski M.S, Vosniakos, G.C, et al. 2020. Opportunities and challenges for exploiting drones in agile manufacturing systems. *Procedia Manufacturing*, 51, pp.527-534. <https://doi.org/10.1016/j.promfg.2020.10.074>
- [18] Mantha B.R, Jung M.K, de Soto B.G, et al., 2020, Generalized task allocation and route planning for robots with multiple depots in indoor building environments. *Automation in Construction*, 119, p.103359. <https://doi.org/10.1016/j.autcon.2020.103359>
- [19] Das G, Cielniak G, Heselden J, et al., 2023. A Unified Topological Representation for Robotic Fleets in Agricultural Applications. *Authorea Preprints*. <https://doi.org/10.22541/au.169357512.24867804/v1>
- [20] Z. Han, Oct. 2023, “Multimodal intelligent logistics robot combining 3D CNN, LSTM, and visual SLAM for path planning and control,” *Frontiers in Neurobotics*, vol. 17. doi: <https://doi.org/10.3389/fnbot.2023.1285673>

Strengthening Accounting Information Systems with Advanced Big Data Mining Algorithms: Innovative Exploration of Data Cleaning and Conversion Automation

Shu Yang

Department of International Trade, Shanxi Vocational College of Tourism, Taiyuan, 030031, Shanxi, China

Email: yangshu_sxcj@163.com

Keywords: accounting information system, big data mining algorithms, decision support, financial performance

Received: October 11, 2024

With the rapid development of big data technology, accounting information systems are facing unprecedented challenges in processing and analyzing massive financial data. In order to improve the efficiency and accuracy of data processing, this article deeply explores the application of big data mining algorithms in optimizing accounting information systems. By introducing advanced big data mining algorithms, accounting information systems have achieved automation in data cleaning, transformation, and analysis, significantly reducing manual intervention and improving data processing efficiency. This article compares the performance of different big data mining algorithms in the analysis of accounting informationization risk transactions. Through practical verification, we found that the selected algorithm performs well in terms of accuracy, reaching over 95%, which is a significant improvement compared to traditional methods. Meanwhile, in terms of computation time, the algorithm has also demonstrated significant advantages, reducing computation time by over 30% when processing datasets of the same size. These performance improvements not only improve the operational efficiency of accounting information systems, but also provide enterprises with more accurate and timely financial information. In addition, this article also conducted a survey on the intelligent management of accounting information systems. We collected valuable opinions on the current status of accounting information intelligent system management by distributing survey questionnaires to in-service MBA and MPAcc students. The survey results show that over 76% of accounting personnel and almost all management personnel (91.18%) agree with the intelligent features of accounting information systems and believe that establishing a separate accounting knowledge base or knowledge management system is necessary. This discovery further emphasizes the importance of optimizing accounting information systems and provides direction for future research.

Povzetek: Analizirani so algoritmi rudarjenja velikih podatkov za izboljšanje računovodskih informacijskih sistemov. Raziskava poudarja potrebo po inteligentnih sistemih v računovodstvu ter potrjuje pozitivne učinke na finančno odločanje in obvladovanje tveganj v podjetjih.

1 Introduction

The intelligent accounting information system can display the financial and operational status of enterprises in different regions through technologies such as geographic information systems (GIS), providing support for the global layout and risk management of enterprises. The system can penetrate into various business areas of the enterprise, provide targeted financial management and decision support, and help the enterprise achieve refined management [1]. By utilizing big data analysis and artificial intelligence technology, the system can deeply mine and analyze the integrated data, discover patterns and trends in the data, and provide scientific basis for enterprise decision-making. On this basis, researching and constructing an intelligent agent dynamic accounting information platform has important practical significance and future value. In the stage of computerized accounting, computers are widely used in daily accounting operations. Such as setting accounting subjects, filling out accounting vouchers, registering accounting books, cost accounting,

and preparing accounting statements. This not only greatly reduces the workload of accounting personnel, but also improves the accuracy and efficiency of accounting [2]. It is not like computerized accounting, where manual accounting is simply simulated through computer technology. At present, the existing accounting information systems in China are mainly used to process transactions that have already occurred [3].

Data warehouse technology is a special data storage technology that can extract data from numerous databases and convert it into a special new format, providing decision analysis for decision-makers. It is a collection of data that reflects continuous historical changes. The data source of a data warehouse is not unique and often includes multiple sources, including internal and external data of the enterprise (such as survey reports, documents, etc.) [4]. It reorganizes, arranges, and stores a large amount of historical or current comprehensive data as needed, providing random queries, comprehensive data, and trend analysis information over time. In the intelligent

interactive visualization accounting information system, to establish a big data analysis platform, Hadoop, a big data processing architecture, can be considered [5].

The emergence of big data mining algorithms provides strong technical support for the optimization of accounting information systems. These algorithms are based on artificial intelligence and machine learning technologies, capable of automatically extracting valuable information from massive data, discovering correlations and patterns between data, and providing decision support and risk management for enterprises. Through the application of big data mining algorithms, accounting information systems can achieve automated and intelligent data processing and analysis, greatly improving work efficiency and accuracy, reducing labor costs and error rates. The application of big data mining algorithms in optimizing accounting information systems has broad prospects and profound significance [6]. Firstly, through automated processing and analysis of data, big data mining algorithms can significantly improve the efficiency of accounting information systems. Traditional accounting information systems require a significant amount of manpower and time to process data, while big data mining algorithms can automate these tasks, greatly shortening the data processing cycle. Secondly, big data mining algorithms can improve the accuracy of accounting information systems [7]. Traditional data processing methods often have errors and omissions, while big data mining algorithms can ensure the accuracy and reliability of data through precise algorithms and models [8]. Therefore, this article designs a big data mining algorithm to provide predictive analysis and decision support models for accounting information systems. By mining and analyzing historical data, algorithms can predict future financial trends and market changes, provide valuable information and suggestions for enterprises, and help them make wiser decisions.

The contribution points of research innovation are as follows:

1) The innovation proposed in this article lies in the successful application of big data mining algorithms to accounting information systems, achieving automation and intelligence in data processing.

2) This article uses advanced model recognition technology and data integrity verification algorithms to verify accounting information data files item by item. This method not only effectively identifies anomalies and errors in the data, but also ensures the integrity and consistency of the data. Meanwhile, by comparing with other algorithms, the identity verification technology proposed in this article has been significantly improved, further enhancing the security of accounting information systems.

In addition to innovations in data processing and validation, this article also proposes the application of big data mining algorithms in predictive analysis and decision support. By mining patterns and trends in data, big data mining algorithms can provide valuable predictive information and decision-making recommendations for enterprises.

2 Related work

In the field of tax management, the application of artificial intelligence technology is gradually demonstrating its enormous potential and value. Some scholars have explored how artificial intelligence can help modernize tax management systems, and through deep learning and data mining, artificial intelligence can identify potential tax risks. And provide warnings to help the tax department take timely prevention and control measures [9]. This can not only reduce the occurrence of tax violations, but also ensure the security and stability of national taxation. In terms of auditing, artificial intelligence can quickly identify potential financial fraud and violations through data analysis and pattern recognition. Artificial intelligence can also automate audit procedures, reduce manual intervention, and improve the accuracy and efficiency of audits [10]. In terms of accounting, artificial intelligence can automate daily accounting and reporting work, reducing the workload of accounting personnel. Artificial intelligence can also provide accurate financial forecasts and decision-making recommendations for enterprises based on historical data and predictive models [11].

The application of financial technology based on artificial Internet of Things, especially the development of big data management algorithms, has brought unprecedented opportunities and challenges to the financial industry. It explores the application and development of big data management algorithms based on artificial Internet of Things in financial technology [12]. In the field of financial technology, the application of AIoT enables financial institutions to collect and analyze large amounts of financial data in real-time, improving the efficiency and accuracy of financial services. In the credit field, big data management algorithms are used to help financial institutions quickly assess borrowers' credit status, achieve accurate risk pricing and risk control [13]. In the investment field, big data management algorithms can provide personalized investment advice and risk management solutions for investors by analyzing and predicting historical data. The application of artificial intelligence technology in the field of accounting is becoming increasingly widespread. From automated accounting processing, intelligent auditing to predictive analysis, artificial intelligence technology is gradually changing traditional accounting methods. Automated accounting processing is one of the earliest applications of artificial intelligence technology in the field of accounting. Through machine learning and natural language processing techniques, artificial intelligence can automatically recognize, classify, and input financial data, greatly improving work efficiency and accuracy [14]. Intelligent auditing utilizes artificial intelligence technology to deeply analyze and mine large amounts of financial data, helping auditors quickly identify potential financial fraud and violations, and improve audit efficiency and accuracy [15].

The management accounting information system should be a comprehensive system that integrates budgeting, performance evaluation, analysis and

forecasting, and decision support. Its core lies in not only acquiring and processing traditional financial accounting information, but also incorporating more non-financial information to meet the diverse decision-making needs of management [16]. By deeply integrating artificial intelligence and big data technology, this system can efficiently extract valuable information from massive data and provide accurate and timely decision-making basis for management. Empowered by big data and artificial intelligence, management accounting information systems can achieve intelligent analysis of data, reveal hidden patterns and trends behind the data, and help enterprises make more informed decisions. At the same time, research achievements from companies such as Informatica in data management and integration (such as the application of platforms like Informatica PowerCenter in data cleaning,

transformation, and integration). This provides strong technical support for the management accounting information system, ensuring data quality and consistency, and further enhancing the system's decision support capabilities [17]. In order to effectively solve a series of large and complex data problems currently existing in enterprises, and to effectively handle various current and historical data distributed inside and outside the enterprise, it is necessary to establish various themed database management systems. In this way, accounting personnel and decision-makers can access the integrated database system through various front-end analysis software tools. And make various decisions based on accurate and comprehensive historical information, quickly putting the overall solution of the business plan into practice [18].

Table 1: Key results gap in reference research

Reference research	Key Results	Technological gap	The necessity of carrying out this work
[9]	Artificial intelligence helps modernize tax management systems and identify tax risks	The scalability and accuracy of tax risk identification models need to be improved	Develop more efficient and accurate tax risk identification algorithms to improve tax management efficiency
[10]	Artificial intelligence quickly identifies financial fraud and violations, automates audit procedures	The automation level of audit procedures is limited, and real-time performance needs to be improved	Enhance the intelligence and real-time performance of audit procedures, reduce manual intervention, and improve audit efficiency
[11]	Artificial intelligence automates accounting and reporting work, providing financial forecasting and decision-making recommendations	The accuracy and real-time performance of the prediction model need to be further optimized	Develop more accurate predictive models to provide valuable decision support for enterprises
[12]	The application of AIoT in financial technology improves the efficiency and accuracy of financial services	The adaptability of big data management algorithms in the financial field needs to be strengthened	Exploring big data management algorithms applicable to the financial sector to enhance the level of financial services
[13]	Big data management algorithms are used for credit risk assessment, achieving risk pricing and risk control	The real-time and accuracy of credit risk assessment models need to be improved	Develop more efficient credit risk assessment algorithms to enhance the risk management capabilities of financial institutions
[14]	Artificial intelligence automatically recognizes, classifies, and inputs financial data to improve accounting efficiency	The accuracy and efficiency of accounting automation processing need to be further improved	Optimize accounting automation processing flow, improve work efficiency and accuracy
[15]	Intelligent auditing utilizes artificial intelligence technology to deeply analyze financial data and improve audit efficiency	The depth and breadth of audit data analysis need to be expanded	Strengthen the ability to analyze audit data, improve audit efficiency and accuracy
[16]	The management accounting information system should obtain relevant information beyond financial accounting information and provide decision support	The integration and intelligence level of information systems need to be improved	Building a more intelligent management accounting information system to provide comprehensive decision support
[17]	Blockchain ensures the authenticity and integrity of accounting records, improves transparency and audit efficiency	The application of blockchain technology in the field of accounting needs to be further expanded and optimized	Explore more application scenarios of blockchain technology in the accounting field to enhance the transparency and audit efficiency of accounting work
[18]	Establish a thematic database management system, provide an integrated database system, and support decision-making	The integration and real-time performance of database management systems need to be strengthened	Build a more integrated and real-time database management system to provide comprehensive information

			support for accounting personnel and decision-makers
--	--	--	--

From the above table, it can be seen that although artificial intelligence technology has made significant progress in fields such as tax management, auditing, accounting, and financial technology, there are still many technological gaps, such as shortcomings in scalability, accuracy, and real-time performance. It is particularly important to carry out this work in order to effectively solve these problems. By conducting in-depth research on the application of big data mining algorithms in accounting information systems, we can further optimize algorithms for tax risk identification, automation of audit procedures, financial forecasting, and decision recommendations. It improves the accuracy and efficiency of accounting automation processing, strengthens audit data analysis capabilities, builds a more intelligent management accounting information system, and explores more application scenarios of blockchain technology in the accounting field. These tasks will help improve the overall performance of accounting information systems and provide more comprehensive, accurate, and real-time information support for enterprises.

incorporating it into the management accounting information system is extremely necessary. Based on this, this article believes that an intelligent modern management accounting information system should include four subsystems (Figure 1).

3 System model construction

3.1. Design of an intelligent accounting analysis management system

At present, accounting informatization in China has been around for nearly 40 years, but there has not been a significant breakthrough in the field of accounting informatization. This article believes that there are two main reasons: firstly, the market focuses more on the construction of accounting business level systems, while neglecting research on accounting information systems related to management decision-making; Secondly, due to the fact that Chinese enterprises often rely solely on experience to make decisions in the risk management process, without the support of corresponding data, it often leads to major mistakes. It can be inferred that establishing a risk management information system and

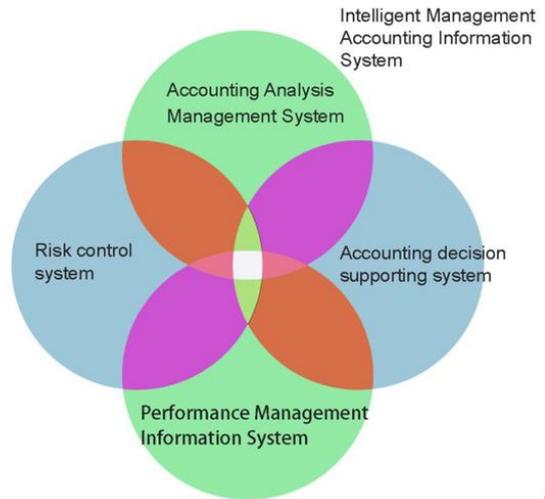


Figure 1: Intelligent management accounting information system

From its functional perspective, it is equivalent to a traditional accounting information system, with the goal of improving data processing capabilities. After introducing data mining technology, this subsystem can effectively enhance data processing capabilities and obtain a large amount of accounting information. Enterprise managers can use performance management information systems to set effective performance goals for each employee and connect the enterprise strategy with each employee.

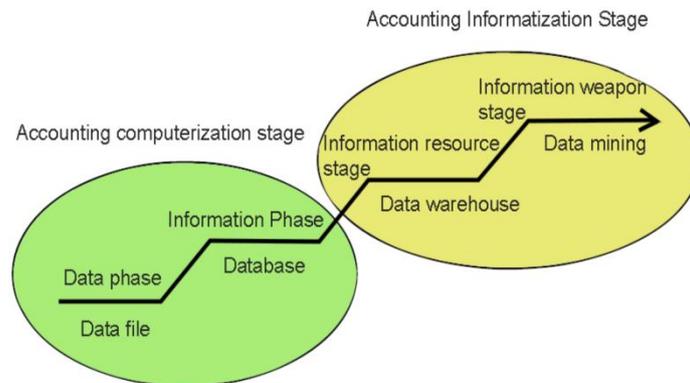


Figure 2: Different stages of accounting informatization

Figure 2 shows the different stages of accounting informatization. Data mining techniques can help analyze historical data, determine which indicators have the

greatest impact on organizational performance, and adjust the indicators and weights in the balanced scorecard accordingly. By mining patterns in historical data, future

performance trends can be predicted, providing a basis for organizations to set reasonable performance goals. Data mining technology can monitor organizational performance data in real-time, detect abnormal situations and issue warnings, helping managers take timely measures to make adjustments. Through in-depth mining and analysis of performance data, problems and deficiencies in the organization's operations and management can be identified, providing support for formulating improvement measures and making scientific decisions. Accounting decision support systems can help decision-makers in enterprises better utilize their financial information to make effective decisions. It is based on modern management science and information technology, utilizing techniques such as quantitative economics, operations research, and control theory to establish relevant models, while utilizing computer technology to solve semi-structured and unstructured accounting problems.

3.2. Overall architecture of intelligent accounting analysis and management system

The accounting analysis management system can be divided into four subsystems: data extraction, data warehouse storage, information processing, and information visualization display (Figure3). With the help of these four subsystems, the goal of integrating, collaborating, sharing, controlling, intelligentizing, and integrating accounting business management can be achieved. The task of the accounting data warehouse storage subsystem is to store accounting data not only in the database according to the requirements of accounting transaction processing, but also synchronously in a data warehouse that is convenient for data mining according to the theme; The task of the accounting information processing subsystem is to process accounting data into accounting information; The task of the accounting information visualization display subsystem is to provide the processing results to information users in various visual ways through human-computer interaction.

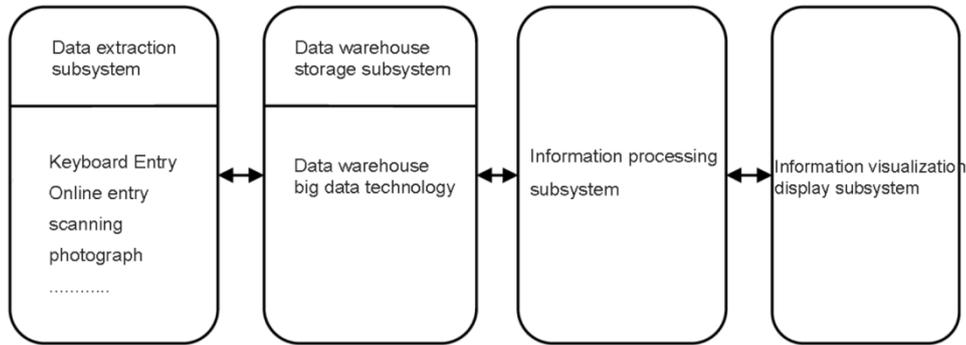


Figure 3: Composition of accounting analysis management system

The data is stored in a standardized format in a unified Data Warehouse (DW) to achieve effective data sharing (Figure 4). The data in DW is theme oriented, such as sales, production, or customers. The data is organized around a specific theme, and when targeted towards theme users, it can determine how the business is conducted and its reasons. Based on DW technology, enterprise managers

can discover the relationships between accounting information through OLAP technology. With the help of DM technology, they can understand the hidden value behind data, discover potential favorable information, and provide favorable support for enterprise decision-making.

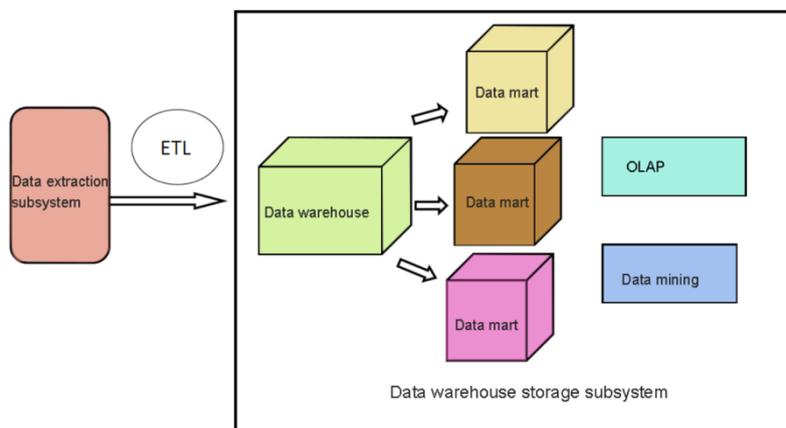


Figure 4: Accounting data retrieval system

3.3. Construction of accounting information system model under big data mining algorithms

In this article, we propose and validate an optimization scheme for accounting information systems based on big data mining algorithms, which has achieved significant results in feature selection and model accuracy. By comparing the data accuracy and training time before and after feature selection, we found that the model after feature selection significantly reduced data complexity and training time while maintaining high accuracy. This discovery emphasizes the importance of feature selection in improving model performance and efficiency. Compared with classical discrete models and project response theory models, our model performs well in accuracy prediction and testing results on three open datasets. Although classical discrete models may still be applicable in some cases, our model demonstrates higher accuracy in most cases. This may be due to our model adopting more advanced algorithms and a more comprehensive feature set, which can better capture

complex patterns and relationships in the data. It should be noted that the accuracy of the theoretical model reflected in the project is the lowest, which may be due to certain issues in model design or implementation. In order to improve the accuracy of the model, it may be necessary to re-examine its assumptions, parameter settings, or data processing methods in the future to ensure that it better adapts to practical application scenarios. In the preprocessing step, the first step is to remove duplicate data, process missing values, correct erroneous data, etc. Convert data of different dimensions to the same dimension to improve the training effectiveness of the model. Traditional accounting information systems often require a significant amount of time and manpower when processing large amounts of data. And big data mining algorithms can automatically process and analyze data, greatly improving data processing efficiency. This enables accounting personnel to obtain accurate information faster and provide more timely data support for enterprise decision-making. Figure 5 shows the data mining process for accounting and financial management systems.

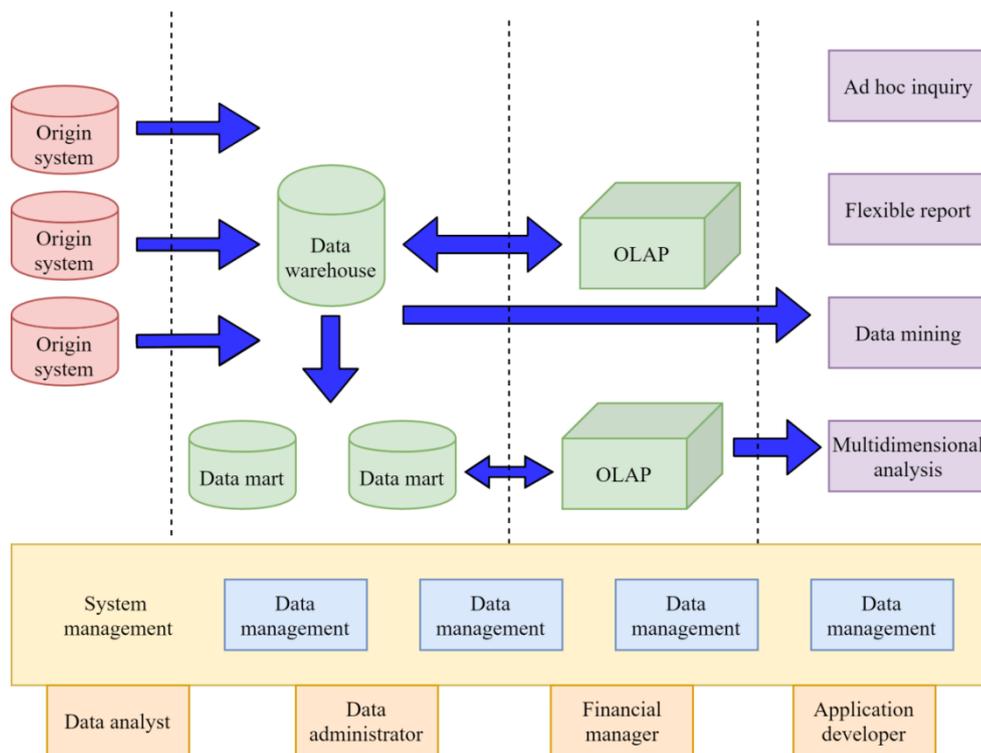


Figure 5: Data mining process for accounting and financial management systems

In each sample $X_i \in S_{\min}$ in the minority class sample S_{\min} , we used the Nearest Neighbors class of the sklearn library to find the nearest neighbors of a given sample. Then, we randomly selected a point from these nearest neighbors and calculated the difference in feature vectors between the sample and this point. Synthesize a minority class sample X_{new} as:

$$X_{new} = X_i + (X_i - Y_j) \times \delta \quad (1)$$

If it is assumed that the upsampling rate is n , then n points should be randomly selected from the k nearest neighbor points found, and then this method is used to synthesize artificial samples for each minority class sample.

Real time segmentation of time series data has broad application prospects in the field of accounting. Through reasonable testing and validation methods, we can ensure the accuracy of segmentation and the effectiveness of the model, thereby providing strong support for financial management and decision-making in enterprises. The

application of real-time segmentation of time series data in the field of accounting is crucial, as it can help us better understand the changing trends of financial data and make accurate predictions and analyses based on them. Real time segmentation is based on certain statistical feature indicators for time-series data. Time series data is a collection of data points arranged in chronological order, usually used to represent the trend of a certain phenomenon over time. In the field of accounting, time series data may include daily transaction volume, monthly revenue, annual profit, etc. The goal of real-time segmentation is to divide time series data into multiple data segments, so that the data sequences in each segment follow the same statistical model. In this way, we can apply corresponding statistical methods to analyze and predict each data segment:

$$X = \{x(t_1), \dots, x(t_i), \dots, x(t_c), \dots\} \quad (2)$$

Where t_c is the current moment. Data stream segmentation is the segmentation of X into a series of consecutive non-empty data segments $\{X_1, \dots, X_j, \dots, X_s, \dots\}$, where:

$$X_j \{x(t_{j,1}), \dots, x(t_{j,l}), \dots, x(t_{j,n_j})\} \subset X, \quad j = 1, 2, \dots, s \quad (3)$$

$$t_{j,l} \in \{t_1, \dots, t_i, \dots, t_c, \dots\}, t_{j,l} < t_{j,l+1}, l = 1, 2, \dots, n_j \quad (4)$$

n_j is the data sequence length of X_j . X_s is the data segment containing the current data $x(t_c)$, which is called the current data segment.

Let the data in X_j be described by a linear regression model:

$$x(t) = f(t, \theta_j) + \varepsilon_j(t), t \in \{t_{j,1}, \dots, t_{j,n_j}\} \quad (5)$$

The linear regression model corresponding to data segment X_j is:

$$f(t, \theta_j) = a_j t + b_j \quad (6)$$

The model parameter vector is:

$$\theta_j = [a_j, b_j]^T \quad (7)$$

$$\sigma_i = (H(i) \cdot u^{m_i})^x \quad (8)$$

$$CHAL = \{(i, v_i)\} \quad (9)$$

Generate a unique serial number for each file block. This serial number can be generated based on a timestamp, file hash, or other unique identifier. Meanwhile, generate a corresponding random number for each file block. This random number can be used for subsequent encryption, signing, or other security operations. After grouping and numbering all user file blocks, merge these file blocks in a certain order into one large file or data stream. This merging process can be based on the sequence number of

file blocks or other sorting rules to ensure that the merged files are ordered.

$$\mu_k = \sum_{i=1}^n v_i m_{k,i} + \mu_r \quad (10)$$

In the formula: μ_r represents the random number generated by the cloud storage server for each user during each verification process. Compute the signature:

$$\sigma = \prod_{k=1}^K \left(\prod_{i=1}^n \sigma_{k,i}^{v_i} \cdot r_k \right) \quad (11)$$

To further improve the reliability of data, error correction techniques can be used, such as solving linear system equations for error correction. This method is usually applied in data communication and storage, and when there may be a small number of errors in the received data, linear equations are solved to recover the original data. In the context of cloud storage, if errors occur during data transmission or storage, linear system equations can be used for error correction through the following steps:

$$s_j = \sum_{k=1}^v Y_k X_k^j, j = 1, \dots, n - k \quad (12)$$

The error value Y_k can be determined where:

$$X_k = \alpha^{i_k} \quad (13)$$

Accompanied by:

$$S_j = R(\alpha^j) = e(\alpha^j) = \sum_{k=1}^v e_k (\alpha^j)^{i_k}, j = 1, \dots, n - k \quad (14)$$

Corresponding α^{i_k} and e_k to X_k and Y_k respectively, then:

$$S_j = R(\alpha^j) = e(\alpha^j) = \sum_{k=1}^v Y_k (X_k)^j, j = 1, \dots, n - k \quad (15)$$

That is, X_k gives the wrong position and Y_k gives the wrong value.

4 Intelligent investigation of accounting information system management

In order to study the necessity of accounting knowledge management and compare the views of management and accounting fields on knowledge management, this article specifically selected in-service MBA (representing management personnel) and MPAcc (representing accounting personnel) students as the objects to conduct a survey on the current situation of accounting information intelligent system management. Hope to obtain the opinions of middle and senior management and finance department personnel on accounting knowledge management in the enterprise. A total of 70 survey questionnaires were distributed, and 55

valid questionnaires were collected, with an effectiveness rate of 78.57%, including 21 MPAcc questionnaires and 34 MBA questionnaires. Among the MPAcc survey respondents, 57.14% had more than 5 years of work experience, and 61.90% in the finance department; Among the MBA survey respondents, 55.88% have more than 5 years of work experience, and 41.18% are middle-level managers, which can basically reflect the actual situation of the enterprise. The main results are as follows:

More than 76% of MPAcc and MBA employees believe that personal experience accumulation contributes to the development of departmental business and the company. This indicates that implicit knowledge related to personal experience is of high importance to departmental business, and there is a clear need for intelligent information management.

Table 1: Intelligent attributes of accounting information systems

Intelligent attributes of accounting information systems	MPAcc	MBA
Existence	76.19%	91.18%
Not existence	23.81%	8.82%
Total	100%	100%

Table 1 addresses the intelligent management issues inherent in accounting knowledge compared to other knowledge. The survey results show that 76.19% of accounting personnel agree with the intelligent characteristics of accounting information systems, while almost all managers (91.18%) agree with the professionalism and specialization of accounting knowledge. managers believe that even if a company already has an information management system, it is still necessary to establish a separate knowledge base or knowledge management system that is suitable for the specialization of accounting knowledge. This also indicates that although knowledge management has been promoted for many years, accounting knowledge management in the financial field is still necessary. It is also not difficult to see that management personnel have a higher demand for accounting knowledge management than accounting personnel, indicating that accounting personnel have a high level of professionalism in providing useful decision-making information, and it is necessary to make up for it through intelligent accounting information management methods.

5 Result analysis

By selecting features from different training datasets, this article ultimately used 70 financial indicators out of 8 features for data analysis

Table 2: Accuracy before and after feature selection

	Quantity of features	Accuracy	Training time
Data before feature selection	70	75.91%	0.0469
Data after feature selection	8	82.16%	0.0411

According to the results in Table 2, it can be clearly observed that the accuracy of the feature set after feature selection has increased by 6.25%, and there are 8 feature subsets after feature selection, which is much lower than the original feature set. This reduces the complexity of the data, improves the learning efficiency of learners, and thus reduces training time. The significant reduction from 70 original financial indicators to 8 feature subsets not only simplifies the data, but also reduces the complexity of the model. Fewer features mean that the model needs to learn and process less information, which usually leads to faster training speed and better generalization ability. The adjustment parameters of the model were tested for motion accuracy prediction on the recommendation model, item reflection theory model, and classical discrete model in three open datasets. From Figure 6, it can be seen that although the results of the classical discrete model are not significantly different from the model proposed in this paper, this does not mean that the classical model has no advantages or value. In practical applications, classical models may still have certain applicability, especially in specific contexts or datasets. Therefore, when choosing a model, it is necessary to weigh specific requirements and data characteristics. The theoretical model reflected in the project has the lowest accuracy, which may be due to issues in model design or implementation. In order to improve the accuracy of theoretical models, it may be necessary to re-examine the assumptions, parameter settings, or data processing methods of the model to ensure that it better adapts to practical application scenarios.

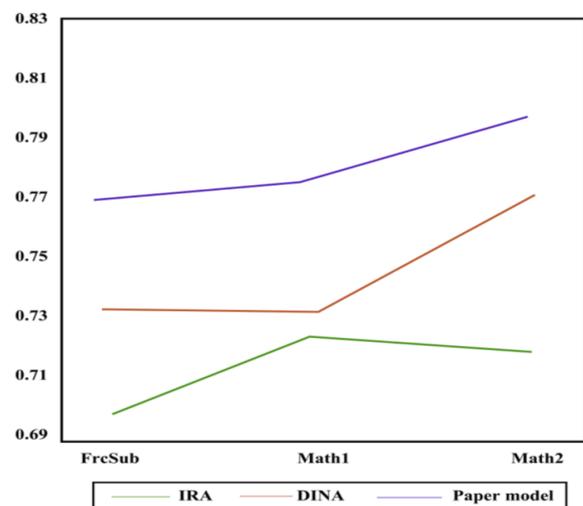


Figure 6 Accuracy prediction and testing results of models in three open datasets

In the experimental validation section of this article, we not only emphasized the significant improvement in model accuracy after feature selection, but also conducted more in-depth statistical analysis and validation to support our conclusions. Firstly, we compared the performance of the models before and after feature selection using the data in Table 2. From the table, it can be seen that the feature set after feature selection has improved accuracy by 6.25%, while the training time has also been reduced. This result not only indicates that feature selection can reduce data complexity and improve model learning efficiency, but also validates our effectiveness and accuracy in the feature selection process. To further enhance the credibility of the conclusion, we conducted a statistical significance test. Specifically, we used paired sample t-test to compare the difference in model accuracy before and after feature selection. The results indicate that this difference is statistically significant ($p < 0.05$), further supporting our conclusion. In addition, we conducted a comprehensive comparison and analysis of the performance of different algorithms on three open datasets. As shown in Figure 6, although there is no significant difference between the results of the classical discrete model and the model proposed in this paper in some cases, we still emphasize the applicability of the classical model in specific environments or datasets. This viewpoint not only reflects our comprehensiveness and objectivity in model selection, but also provides readers with richer information and references.

As shown in Figures 7 and 8. As the sample size increases, the training time also increases. This is because larger datasets require the model to perform more calculations and iterations to fit complex patterns and relationships in the data. Therefore, when dealing with large-scale datasets, it is necessary to consider the training efficiency of the model and the required computational resources. It is valuable to compare the accounting risk identification model in this article with the models in general traditional methods in terms of accuracy and average absolute error. This helps to evaluate the advantages and disadvantages of the model in practical applications, as well as its performance on different indicators.

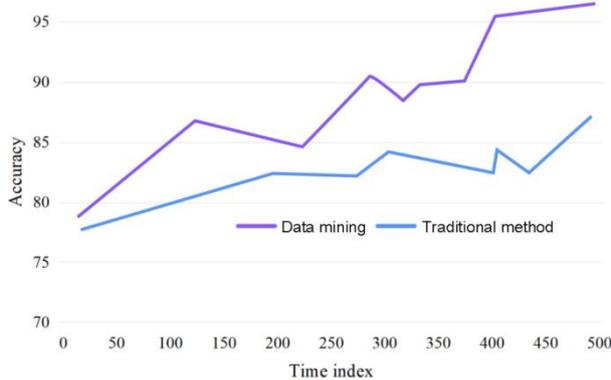


Figure 7: Accuracy comparison

Through the case analysis of introducing the algorithm in this article into enterprise management accounting, it can be understood that intelligent accounting information systems can save costs and create more value for enterprises. Through continuous exploration in the field of management accounting informatization, multiple major businesses of enterprises have achieved first place in the industry, and the employee development index is far higher than the average of Baosteel Group. According to relevant professionals, an intelligent management accounting information platform is the future development trend of enterprises. In order to win in competition, enterprises must combine various IT technologies (big data, BI, artificial intelligence, cloud computing, etc.) to establish a management accounting information system that is in line with their own development.

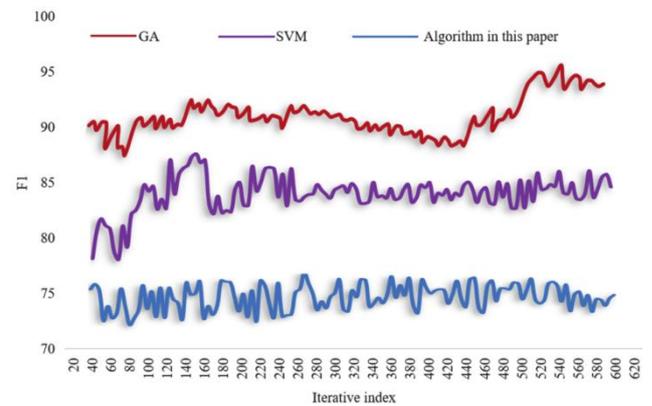


Figure 8: Comparison of F1 values for different algorithms

According to the data in Figure 8, the recall rate and F1 value of our algorithm are higher than the other two algorithms. This result verifies that the algorithm has certain advantages compared to other algorithms. In the input factors of the model, factors such as actual situation are related to the accuracy of the prediction effect. Analyzing the prediction effect from multiple factors can improve prediction accuracy, and higher prediction accuracy can enable enterprises to effectively control financial risks.

In this article, we propose and validate an optimization scheme for accounting information systems based on big data mining algorithms, which has achieved significant results in feature selection and model accuracy. By comparing the data accuracy and training time before and after feature selection, we found that the model after feature selection significantly reduced data complexity and training time while maintaining high accuracy. This discovery emphasizes the importance of feature selection in improving model performance and efficiency. Compared with classical discrete models and project response theory models, our model performs well in accuracy prediction and testing results on three open datasets. Although classical discrete models may still be applicable in some cases, our model demonstrates higher accuracy in most cases. This may be due to our model

adopting more advanced algorithms and a more comprehensive feature set, which can better capture complex patterns and relationships in the data. It should be noted that the accuracy of the theoretical model reflected in the project is the lowest, which may be due to certain issues in model design or implementation. In order to improve the accuracy of the model, it may be necessary to re-examine its assumptions, parameter settings, or data processing methods in the future to ensure that it better adapts to practical application scenarios.

6 Conclusion

This study delves into the optimization of accounting information systems based on artificial intelligence algorithms, especially in the context of the big data era. By introducing big data mining algorithms, accounting information systems have been significantly improved, effectively addressing the challenges of processing massive amounts of data. These algorithms not only improve data processing efficiency and reduce manual intervention, but also achieve automation in data cleaning, transformation, and analysis. These advances not only optimize the workflow of accounting practice, but also enhance the accuracy and efficiency of data processing. Experimental verification shows that compared with traditional single user authentication techniques. This means that in practical applications, accounting personnel can obtain and process data faster, improving work efficiency. Although this study conducted model recognition and data integrity verification in terms of data security and privacy protection, it did not delve into how to implement higher-level data protection and privacy encryption technologies at the algorithmic level. With the continuous development of big data and artificial intelligence technology, data security and privacy protection will become increasingly important issues. In response to the above limitations, future research can further explore the application of other artificial intelligence technologies (such as deep learning, reinforcement learning, etc.) in accounting information systems to find more suitable combinations of algorithms and technologies for specific scenarios, thereby further improving the performance and efficiency of the system.

References

- [1] Tang, W., Yang, S., & Khishe, M. (2023). *Profit prediction optimization using financial accounting information system by optimized DLSTM*. *Heliyon*, 9(9), 1. <https://doi.org/10.1016/j.heliyon.2023.e19431>
- [2] Bhima, B., Zahra, A. R. A., Nurtino, T., & Firli, M. Z. (2023). *Enhancing organizational efficiency through the integration of artificial intelligence in management information systems*. *APTISI Transactions on Management*, 7(3), 282-289. <https://doi.org/10.33050/atm.v7i3.2146>
- [3] Hu, K. H., Chen, F. H., Hsu, M. F., & Tzeng, G. H. (2021). *Identifying key factors for adopting artificial intelligence-enabled auditing techniques by joint utilization of fuzzy-rough set theory and MRDM technique*. *Technological and Economic Development of Economy*, 27(2), 459-492. <https://doi.org/10.3846/TEDE.2020.13181>
- [4] Berdiyeva, O., Islam, M. U., & Saedi, M. (2021). *Artificial intelligence in accounting and finance: Meta-analysis*. *International Business Review*, 3(1), 56-79. <https://doi.org/10.37435/NBR21032502>
- [5] Xia, W. H., Zhou, D., Xia, Q. Y., & Zhang, L. R. (2020). *Design and implementation path of intelligent transportation information system based on artificial intelligence technology*. *The Journal of Engineering*, 2020(13), 482-485. <https://doi.org/10.1049/joe.2019.1196>
- [7] Hua, H., Wei, Z., Qin, Y., Wang, T., Li, L., & Cao, J. (2021). *Review of distributed control and optimization in energy internet: From traditional methods to artificial intelligence-based methods*. *IET Cyber-Physical Systems: Theory & Applications*, 6(2), 63-79. <https://doi.org/10.1049/cps2.12007>
- [9] Saragih, A. H., Reyhani, Q., Setyowati, M. S., & Hendrawan, A. (2023). *The potential of an artificial intelligence (AI) application for the tax administration system's modernization: the case of Indonesia*. *Artificial Intelligence and Law*, 31(3), 491-514. <https://doi.org/10.1007/s10506-022-09321-y>
- [10] Faccia, A., & Petratos, P. (2021). *Blockchain, enterprise resource planning (ERP) and accounting information systems (AIS): Research on e-procurement and system integration*. *Applied Sciences*, 11(15), 6792. <https://doi.org/10.3390/app11156792>
- [12] Andronic, M., Iatagan, M., Uță, C., Hurloiu, I., Dijmărescu, A., & Dijmărescu, I. (2023). *Big data management algorithms in artificial Internet of Things-based fintech*. *Oeconomia Copernicana*, 14(3), 769-793. <https://doi.org/10.1109/ISAIEE57420.2022.00032>
- [13] Ahmad, A. (2024). *Ethical implications of artificial intelligence in accounting: A framework for responsible ai adoption in multinational corporations in Jordan*. *International Journal of Data and Network Science*, 8(1), 401-414. <https://doi.org/10.5267/j.ijdns.2023.9.014>
- [14] Peng, Y., Ahmad, S. F., Ahmad, A. Y. B., Al Shaikh, M. S., Daoud, M. K., & Alhamdi, F. M. H. (2023). *Riding the waves of artificial intelligence in advancing accounting and its implications for sustainable development goals*. *Sustainability*, 15(19), 14165. <https://doi.org/10.3390/su151914165>
- [15] Ahmad, A. (2024). *Ethical implications of artificial intelligence in accounting: A framework for responsible ai adoption in multinational corporations in Jordan*. *International Journal of Data and Network Science*, 8(1), 401-414. <https://doi.org/10.5267/j.ijdns.2023.9.014>
- [16] Peng, Y., Ahmad, S. F., Ahmad, A. Y. B., Al Shaikh, M. S., Daoud, M. K., & Alhamdi, F. M. H. (2023). *Riding the waves of artificial intelligence in advancing accounting and its implications for*

- sustainable development goals*. Sustainability, 15(19), 14165. <https://doi.org/10.3390/su151914165>
- [17] Zhang, Y., Xiong, F., Xie, Y., Fan, X., & Gu, H. (2020). *The impact of artificial intelligence and blockchain on the accounting profession*. IEEE Access, 8, 110461-110477. <https://doi.org/10.1109/ACCESS.2020.3000505>
- [18] Lee, C. S., & Tajudeen, F. P. (2020). *Usage and impact of artificial intelligence on accounting: Evidence from Malaysian organisations*. Asian Journal of Business and Accounting, 13(1), <https://doi.org/10.22452/ajba.vol13no1.8>

Advances in Machine Learning Framework for Near-Infrared Spectroscopy: A Taxonomic Review on Food Quality Assessment

Nguyen Thi Hoang Phuong¹, Hieu Nguyen Van², Xuan Nguyen Thi Thanh², Phien Nguyen Ngoc^{3,4,*}

¹Faculty of Information Technology, Pham Van Dong University, Quang Ngai, Vietnam

²The University of Danang, University of Science and Technology, Viet Nam

³Center for Applied Information Technology, Ton Duc Thang University, Ho Chi Minh City, Vietnam

⁴Faculty of Information Technology, Ton Duc Thang University, Ho Chi Minh City, Vietnam

E-mail: nthphuong@pdu.edu.vn, nvhieuqt@dut.udn.vn, nttxuan@dut.udn.vn, nguyennngocphien@tdtu.edu.vn

*Corresponding author

Keywords: Machine learning, deep learning, near-infrared spectroscopy, advances in framework, food quality

Received: November 01, 2024

This review taxonomically analyzes and evaluates recent advances in machine learning (ML) frameworks applied to near-infrared spectroscopy (NIRS) for food quality assessment. Through a comprehensive literature search across IEEE Explore, ScienceDirect, and Springer (2021-2024), we examine key framework components: data acquisition, public datasets, preprocessing, wavelength selection, and advanced ML architectures. Our analysis reveals the current state: miniaturized devices and multi-device data collection are expanding spectral coverage, while public datasets focus mainly on nutritional indices, lacking safety-related data. Framework-wide challenges persist in device compatibility, dataset comprehensiveness, and model interpretability. Recent advances show promising developments through: specialized deep learning architectures achieving 97-100% accuracy, data transformation techniques (2D-COS, GAFD) enhancing interpretability, hybrid traditional-deep learning models, and effective transfer learning for cross-device applications. Based on these insights, we propose three critical research directions: expanding food safety datasets through regulatory partnerships, developing multi-level fusion for heterogeneous device data, and creating automated techniques for model optimization and interpretability. These directions are vital for advancing ML-NIRS applications in food quality assessment, improving both efficiency and reliability.

Povzetek: Analizirani so napredki v strojnih učnih modelih za spektroskopijo bližnjega infrardečega spektra (NIRS) pri oceni kakovosti hrane. Pregled obsega ključne komponente, kot so zbiranje podatkov, predprocesiranje in izbira valovnih dolžin. Predlagane so tri raziskovalne smeri: širitev podatkovnih zbirk, razvoj fuzije večnivojskih podatkov in avtomatizacija optimizacije modelov za boljše zanesljivost ocenjevanja kakovosti hrane.

1 Introduction

Food quality and safety have emerged as critical concerns for both the food industry and global consumers [1]. The burden of foodborne diseases and economic losses due to poor quality or spoiled food at the production and distribution stages is enormous [2], requiring careful monitoring of food composition regularly. NIRS, as an analytical technique that can provide complex “chemical fingerprints” of food samples related to their composition, quality, and safety [3–5], has been combined with classical statistical methods and advanced ML to address this issue.

ML techniques and IoT development have brought about considerable changes in many fields [6–8]. Applying ML, mainly supervised learning, to multivariate NIRS spectral analysis has significantly changed food quality assessment and assurance. These studies have been diverse across various data types and increased rapidly in the past two years [9]. From 2022 to 2024, along with traditional ML meth-

ods such as Principal Component Analysis (PCA), Partial Least Squares (PLS), Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), and Deep Learning (DL) such as Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), Autoencoders (AE), as in Figure 1, there are four major trends in NIRS applications. This includes (1) detecting contaminated and adulterated food, identifying adulterants, and determining the level and residual concentration of chemicals in agricultural and livestock products [10, 11]; (2) developing sustainable agriculture through monitoring crop growth, soil nutrients, and various components of crops to improve care and early detection and treatment of crop diseases [12]; (3) determining the optimal harvest time to achieve maximum economic yield [13]; and (4) evaluating product quality, particularly for high-value economic items [14–19], etc.

However, ML for NIRS is still inceptive compared to other fields for several reasons. First, ML on NIRS spectra requires specialized data, which is difficult to collect due

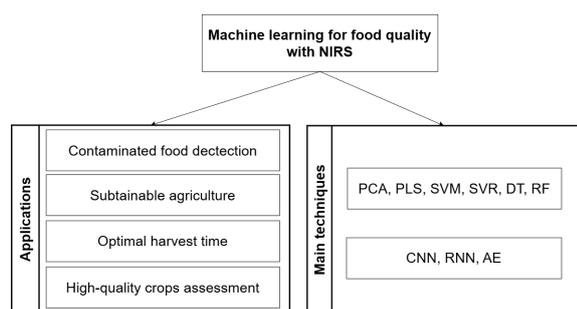


Figure 1: NIRS applications and ML techniques for NIRS

to expensive spectrometers and reference chemical data, in chemical content determination problems [9]. Second, studies are often published in agricultural or interdisciplinary chemometrics journals that combine data science and chemistry, as in Figure 2, so computer scientists' access to technical developments is more limited. With the potential of developing ML to solve social food quality problems, this needs to be further promoted by supporting a technical overview.

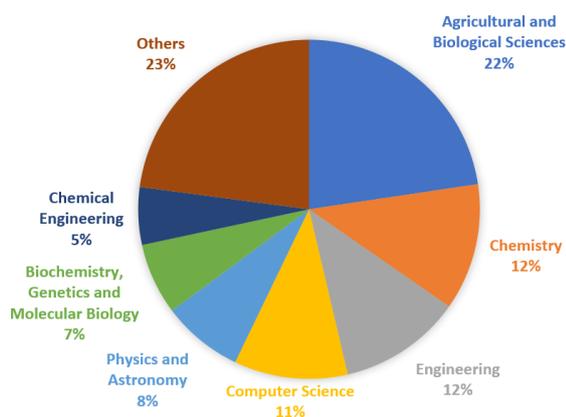


Figure 2: Subject area from 2021 to 2024 according to Scopus analysis in ML for NIRS food quality

Therefore, this study aims to shift the focus to technical surveys, technological advances, availability of public datasets, and development potentials not mentioned in previous review studies. We conduct two main review bases:

1. First, we summarize and discuss the contributions and limitations of recent review articles on ML in NIRS and food analysis in particular. From there, we find gaps that need to be exploited and further evaluated on techniques and data.
2. Second, we synthesize and evaluate recent new research articles, classifying and highlighting information on data (NIRS tools, data fusion techniques, public datasets), as well as ML techniques (preprocessing, wavelength selection, and advanced ML

architectures) that have not been covered in existing reviews.

From this background, gaps were identified from a computer science perspective to conduct future research in food inspection.

2 Methodology

In the fourth quarter of 2023, a thorough investigation was conducted through IEEE Explore, ScienceDirect, and Springer, employing controlled vocabulary in ML, NIRS, and food quality analysis, as in Figure 3. The search focused on emerging ML techniques for NIRS in food quality assessment, utilizing specific terms such as “deep learning”, “chemometrics”, “NIR spectroscopy”, and “food”. In addition to these above primary keywords, we conducted deeper searches focusing on specific components of our ML framework. Each framework component served as secondary keywords - notably “preprocessing” and “wavelength selection” - to thoroughly identify recent studies focusing on improvements in these critical areas. For machine learning algorithms, we specifically searched for both traditional methods (PCA, PLS, SVM) and emerging deep learning architectures (CNN, RNN, AE, GAN) to track their evolution and applications in NIR analysis. Additionally, the NIR dataset was also searched extensively on Mendeley and Zenodo. This hierarchical search approach, structured according to our ML framework taxonomy, enabled us to systematically evaluate recent advances in specific methodological aspects rather than just general applications. Through this focused search strategy, we could better assess how recent research has contributed to advancing different components of the ML framework in NIR spectral analysis.

The literature review methodology branches into two distinct paths: (1) comprehensive review papers and (2) original research articles. The first branch focuses on conducting rigorous analyses of existing literature to identify key challenges, significant contributions, and unexplored territories within the machine learning domain for NIR analysis. The second branch encompasses original research papers, systematically categorized according to their contributions to the ML framework - spanning from data acquisition and public datasets to preprocessing/wavelength selection and advanced architectural innovations. This dual-branch approach ensures both a broad understanding of the field's current state through synthesized reviews and a detailed examination of specific technical advancements through original research contributions.

Stringent filters were applied, including English language restriction and consideration of peer-reviewed articles, reviews, books, book chapters, and conference papers from the four years (2021-2024). The research database has been updated to include publications up to September 2024 to ensure the most informed and up-to-date discussion. The

strategy aimed to capture the latest innovations at the intersection of machine learning and NIRS for non-destructive and rapid evaluation of diverse food quality traits, excluding older publications beyond the scope of emerging techniques.

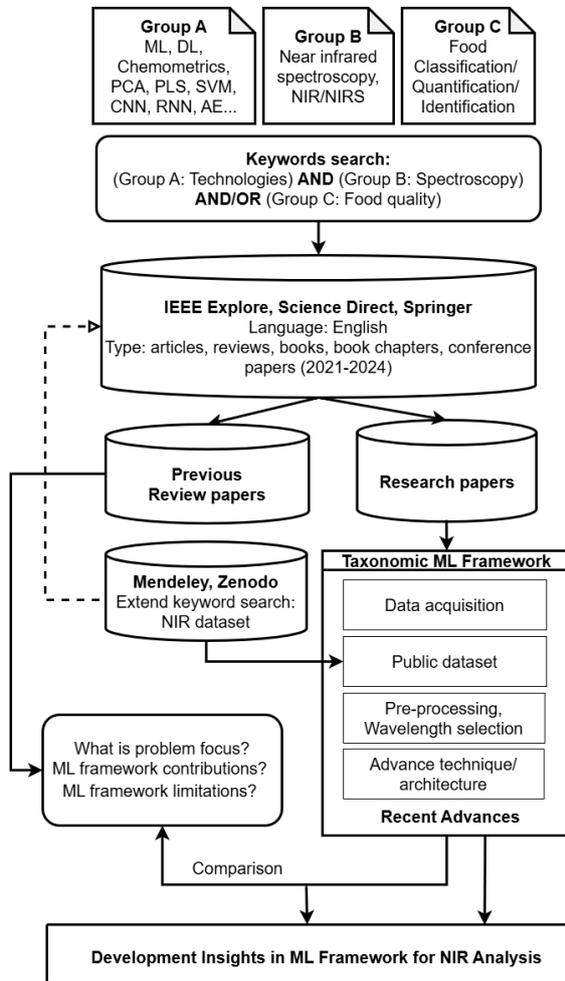


Figure 3: Research process flow chart

3 Previous review studies: an overview

Our analysis of previous ML in NIR spectroscopy for food quality and safety draws from 18 review articles published between 2021 and 2024. This comprehensive overview reveals significant progress in applying advanced ML techniques to NIR data and exposes critical challenges, as in Table 1, 2. While studies demonstrate the potential of methods ranging from traditional multivariate analysis to sophisticated deep learning algorithms, they also underscore persistent limitations in datasets, model optimization, and real-world applicability.

Building upon this overview, we conduct an additional review in the next section to explore unaddressed aspects that significantly impact machine learning trends in food spectroscopy. This supplementary analysis aims to fill crucial gaps and provide insights into emerging directions that could shape the future of ML-driven NIRS in food analysis.

These studies are mainly related to ML for food quality and safety and cover a range of applications or techniques. Recent review articles have examined hyperspectral imaging and NIR spectroscopy combined with advanced algorithms for non-invasive assessment of parameters, including nutritional composition (e.g., protein, moisture, fatty acids), adulteration/defect detection, and geographical origin discrimination in various food products. Both traditional multivariate analysis methods, like PCA and PLS, and increasingly sophisticated ML algorithms, including SVM, ANN, and CNNs, have been explored for relevant tasks such as multi-class food classification and quality prediction. This demonstrates the general feasibility of data-driven modeling approaches on spectroscopic data for food evaluation. However, significant limitations persist regarding dataset availability and model optimization, transferability, and interpretability, specifically for NIR food applications using advanced machine learning.

Despite the widespread application of NIRS in food quality and safety assessment, significant challenges remain in developing and sharing suitable datasets for machine learning research in this field. Firstly, current technique reviews tend to cover NIR datasets broadly without an in-depth analysis tailored to the food domain. Secondly, the mentioned datasets primarily consist of Vis-NIR spectral range (< 1000 nm) stored in MATLAB data files, which poses challenges for developing ML research applications (typically developed with Python). However, this spectral range is often considered less informative in chemical information than the 1000–2500 nm [4]. Thirdly, other current shared datasets with broader coverage (400 - 2500 nm) are just suitable for simple classification tasks, lacking the detailed chemical information required for laboratory-based quality assessment regression problems. Besides, researchers also highlighted significant limitations in collecting valuable data, labeling, data enrichment, and practical deployment due to high costs. Last but not least, many previous studies utilizing NIRS for food quality and safety inspection have relied on small datasets, often with fewer than 200 samples, limiting model robustness and generalizability. Therefore, building NIRS spectral datasets with appropriate wavelength bands, relevant chemical parameterization, and proper labeling would be more practical when addressing real-world problems.

Additionally, reviews focusing specifically on deep learning also need more details regarding optimal network architectures, data requirements regarding sample size and variability, and quantitative benchmarking on relevant food NIR datasets. There is no in-depth discussion of deep learning or other advanced machine learning methods, nor is

Table 1: Previous Review Studies (1)

The issues	Review focus	Related Contributions	Limitations
Quantification of food bioactives by NIR spectroscopy: Current insights, long-lasting challenges, and future trends [3]	Factors affecting model performance. Algorithm used: Mostly PLS; SVM, MLR, BP-ANN, CNN	Effects of sample prep, analyte concentration, instrument features on performance. Compares benchtop/portable NIR. Proposes FAIR data management. Suggests theoretical calculations for interpretation.	Limited datasets (< 200 samples). Difficulty in choosing pre-processing/regression methods. Interpretability/transferability issues. Lacks DL focus.
Food quality 4.0: From traditional approaches to digitalized automated analysis [5]	Traditional vs emerging techniques. Algorithm used: Mostly PLSR, PLS-DA; SVM	Industry 4.0 innovations (AI, DL, sensors) in spectroscopic. Portable/miniatuized NIR-AI for evaluation. HSI as non-destructive quality technique.	Brief NIRS-food analysis mention. No in-depth NIRS/DL applications with food datasets. Mostly Vis-NIR datasets.
DL for NIRS data modeling: Hypes and benefits [9]	Potential benefits and pitfalls of using DL for modeling NIRS	DL auto-transforms spectral data without preprocessing. Shallow DL success with small datasets (< 1000). DL efficiency for complex food analysis tasks (multi-class, multi-response).	Small, under-optimized datasets in DL-chemometrics comparisons. Limited food NIR spectra DL modeling. No large food quality/safety datasets for DL benchmarking.
Multivariate analysis of food fraud: A review of NIR based instruments in tandem with chemometrics [10]	Chemometrics with NIRS, HSI. Algorithm used: Mostly PLSR, PLS-DA; SVM, PCA, SIMCA, ANN, KNN	Overview of NIR/HSI principles. Summarizes chemometrics for spectral processing. Reviews NIR/HSI with chemometrics for adulterant detection. Compares classification/regression models. Discusses advantages/limitations for food authenticity.	No DL discussion. Focuses on PCA, PLS, and SVM. No specific NIR datasets for food. Lab-prepared samples, not real-world fraud. Limited assessment of method robustness and applicability.
AI-based techniques for adulteration and defect detections in food and agricultural industry: A review [11]	AI techniques combined with sensors. Algorithm used (NIR-specific): Mostly SVM, PLSR, ANN, PCA, CNN, Random Forest.	AI for food authentication/quality. CNN (2015-2022). Challenges: technique standardization, algorithm selection, data fusion, fast detection, severity quantification, framework development.	Lack of sensor device, data acquisition, preprocessing details. Impact on model performance not discussed.
Computer vision and DL in insects for food and feed production [20]	Applications	CV and NIRS for non-invasive assessment of nutritional composition, moisture, protein, fat, fatty acids in live insects.	Limited NIRS applications mentioned. No technical aspects discussed.
Application of NIRS for the nondestructive analysis of wheat flour: A review [21]	Application. Algorithm used: MPLS, PLSR, RF, RBF, LDA.	NIR fundamentals, recent developments for wheat flour quality/safety assessment. Four development areas: data quality, chemometrics, affordable tools, data integration.	Focuses on traditional ML for classification/regression.
Quality analysis and authentication of nutraceuticals using near IR (NIR) spectroscopy: trends and applications [22]	Novel analytical trends and applications from a chemistry/metabolomics perspective	NIR trends for nutraceutical quality control (HSI, portable devices). Targeted/untargeted metabolomics applications. Geographical classification.	No DL discussion for NIR data. No specific NIR nutraceutical datasets were analyzed.

Table 2: Previous Review Studies (2)

The issues	Review focus	Related Contributions	Limitations
A research review on DL combined with HSI in multiscale agricultural sensing [23]	Applications and limitations. Algorithm used (specific to NIR): Mostly CNN, AE	DL models for food quality, ripeness, moisture, nitrogen, chlorophyll, sugar prediction. HSI range (250-2500 nm) is broader than NIRS.	Data collection/real-world application challenges. Limited food authentication focus. CNN, SAE, and RNN without detailed evaluation.
Efficient extraction of deep image features using CNN for applications in detecting and analyzing complex food matrices [24]	Principle, architectures, applications of feature extraction methods, CNNs	1D CNN feasibility for NIRS food classification/defect detection. CNN features outperform traditional ML. HSI-CNN examples for cereal variety/quality classification.	Not NIRS-specific. Limited NIRS-based food analysis datasets. Lacks DL challenges for food NIRS data.
A Review of ML for NIRS [25]	ML, especially DL. Algorithm used: Mostly PLS, ELM, SVR, SVM, SLFN, DT, RF, AE, CNN, RNN, LSTM, GRU, GAN.	Summarizes NIR modes, instruments, preprocessing, datasets, feature selection. Covers traditional ML (PLS, SVM, ELM) and DL (CNN, RNN, autoencoders) for NIR food data.	No in-depth ML-NIR food analysis. Limited food spectroscopy dataset details. No model performance comparison. No data augmentation/transfer learning discussion.
Are standard sample measurements still needed to transfer multivariate calibration models between NIR spectrometers? [26]	Recent developments in calibration transfer (CT) methods	Mentions DL for multivariate calibration and transfer learning in NIRS model updating.	No in-depth DL-NIR food analysis discussion. No public DL-NIR food datasets were mentioned. Brief food applications (temperature/form adaptation).
Recent advances and application of ML in food flavor prediction and regulation [27]	Algorithm used: SVM, DT, RF, KNN, ELM, ANN	Principles, advantages, application, challenges of ML for food flavor prediction/regulation.	Traditional ML focus. Few NIRS flavor prediction studies. Small sample sizes. Limited public NIRS flavor datasets.
ML applications for multi-source data of edible crops[28]	Fusion of multi-source data with ML techniques	CNN and ResNet for edible crop classification using 2D spectral/HSI.	Not NIRS-specific. No NIRS algorithm performance
AI in sensory and consumer studies of food products [29]	Applications. Algorithm used: ANN, SVM, CNN	ML, particularly NIRS, for predicting sensory responses from physicochemical data.	ANN common, but basic supervised learning prevalent. More DL research needed for complex spectroscopic data. Lack of public NIRS-sensory datasets limits validation/research.
DL in analytical chemistry [30]	Applications. Algorithm used: CNN, DNN, LSTM, GAN, AE	DL applications in analytical chemistry. AEs for cereals, CNNs for vibrational spectral data classification/quantification, chemical component determination, geographical discrimination.	No in-depth NIRS-food DL analysis. Lacks food application datasets discussion. No focus on DL techniques/datasets for food NIRS.
Recent advances in assessing qualitative and quantitative aspects of cereals using nondestructive techniques[31]	Chemometrics based AI & ML. Algorithm used: PLS, PCR, LDA, PLSR, PCA, KNN	Chemometrics for cereal analysis. NIRS for amylose, moisture, texture, authentication. Preprocessing, models, performance parameters for NIRS.	Lacks DL/large dataset focus for NIRS. No public datasets. Limited portable NIRS discussion for on-site cereal analysis.

there any analysis provided comparing the performance of different algorithms specifically for NIRS data. The reviews mainly focus on conventional chemometrics techniques like PCA, PLS, SVM, etc. Even where public reference datasets exist, few studies thoroughly validate and compare deep learning techniques using these resources. Based on the mentioned contributions and limitations, in these review studies, we focus on exploring key aspects related to measurement devices, data collection, publicly available NIR food datasets, and new ML architectures applied to NIRS spectral data. This review aims to contribute to expanding multidisciplinary progress at the intersection of NIRS, data science, food science, and industrial applications.

4 Recent advances of machine learning for NIRS

Based on the comparison with previous research results in section 3, this section classifies new contributions to machine learning frameworks in processing NIR spectral data in classification, pattern recognition, and component content regression problems, as in Figure 4.

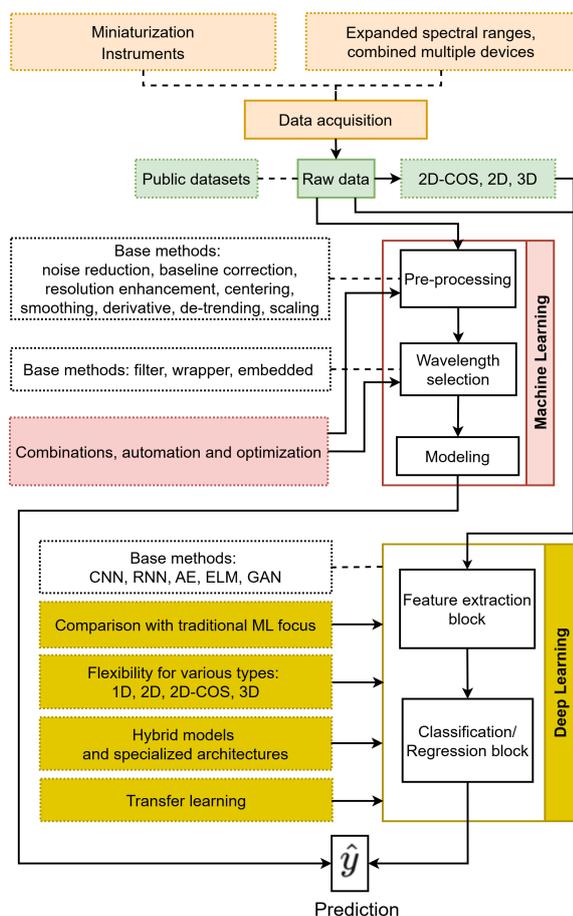


Figure 4: Advances in ML framework for NIRS

4.1 Data acquisition

Effective and accurate data acquisition using near-infrared spectrometers is crucial for food quality assurance. This relies on appropriate spectrometer selection, compatible measurement modes, and relevant reference components.

4.1.1 Miniaturization NIR spectral instruments

In recent years, NIR spectral instruments have undergone a notable transformation in size and portability, progressing from traditional benchtops to various compact, handheld, pocket-sized, miniaturized, and real-time versions [4, 32]. Traditionally, benchtop systems provide standardized design, broad spectral range, and high performance. For instance, the NIRSystems 6500 covers the full NIR range from 400 nm to 2500 nm, and the FT-NIR Frontier operates between 900-2500 nm [33]. The NIR spectral dataset covering the entire spectral range is typically measured using this instrument.

On the other hand, portable NIR spectrometers prioritize practical applications and offer greater flexibility, revolutionizing modern measurement techniques. However, this comes at the cost of design uniformity and spectral range. The diverse designs and technologies used in portable systems lead to variations in spectral coverage and resolution compared to their benchtop counterparts. For example, the SCiO operates within a narrow 740-1070 nm range [34], while the IAS-3120 [35], NIR-S-G1 [36], and NIR-M-R2 [37] function within the more common 900-1700 nm range. Specialized instruments like the Spectromètre portable NIR [38] cover a distinct 1750-2150 nm range. Many other devices from various manufacturers in [32, 34, 39–41] also have differing spectral ranges. This diversity challenges the development of multivariate analysis models with heterogeneous data. Despite these challenges, the ongoing miniaturization of NIR spectral instruments continues to expand their applicability across diverse fields, driving innovation in portable spectroscopic analysis.

4.1.2 Expanded spectral ranges and combined multiple devices

NIR spectrometers are capable of operating in various measurement modes: diffuse reflectance for solid surfaces; transmittance for gases, liquids, or cuvette-contained semi-solids; transreflectance combining reflectance and transmittance for semi-solids without cuvettes; interactance using a specialized probe for enhanced solid sample analysis; and a transmittance mode accounting for scattering in dense samples [25, 42]. Each mode varies in cost, signal quality, and speed, and the selection depends on the specific needs of sample material and analysis requirements. These modes are determined by the detector, wavelength selector, and light source. For instance, the portable SCiO spectrometer uses a silicon array detector, a bandpass filter as the wavelength selector, and an LED light source, making it suitable for transmittance measurements of liquid or thin solid sam-

ples like vanilla solutions to determine vanillin content [34]. The DLPR NIRscan Nano employs an InGaAs detector, reflective diffraction grating, and tungsten lamps, enabling diffuse reflectance analysis of solid samples such as adulterated almond flour [43].

NIR datasets are typically collected by scanning multiple times and averaging values for each sample, potentially including 1500-1700 variables for the full NIR range [4]. However, achieving such comprehensive datasets is uncommon, due to challenges in sensor sensitivity across the broad spectral range and the high cost of instruments for full spectrum. NIRS data in various types of foods typically collected exclusively from either the short-wave (800-1100 nm, often together with the Vis-NIR region) or long-wave (1100-2500 nm) regions or the middle region [44]. The selection of NIR spectral range is crucial for analysis effectiveness. The key difference lies in signal strength across wavelengths [4]. In short-wave regions (700-1000 nm), signals become progressively weaker with 3rd and 4th overtones and show extensive band overlap, making detection and analysis more challenging. In contrast, long-wave regions (1000-2500 nm) contain stronger first and second overtones with better peak separation. This spectral range is especially effective at detecting absorption from specific molecular bonds - primarily those containing hydrogen atoms (O-H, C-H, N-H) and certain strong bonds (C=O, C≡N) - which produce characteristic absorption patterns useful for chemical analysis and machine learning model development. While long-wave regions typically provide better analytical performance, the choice of spectral range often requires balancing between analysis accuracy and economic considerations, as instruments for shorter wavelengths can use simple glass components, making them significantly more cost-effective than specialized equipment needed for longer wavelengths. The sub-1700 nm range is the most common due to its sufficiency in capturing key chemical components of the sample, as demonstrated by numerous studies [4, 44].

Recently, obtaining comprehensive and valuable information across the entire NIR spectral range often requires employing at least two instruments to capture data from both short-wave and long-wave regions. For instance, the NIR spectral of almond flour [43] were obtained from 900-2500 nm using three portable spectrometers: DLPR NanoNIRscan (900-1700 nm), MicroNIR 1700 (950-1650 nm), and NeoSpectra FT-NIR (1350-2500 nm) with 16 nm resolution. The classification models using portable near-infrared devices achieved 100% sensitivity and over 95% specificity in identifying almond flour adulterations above 5% (w/w), while the PLS regression models obtained coefficients of determination above 0.90 and RMSEP values between 3.2-4.8% for quantifying almond flour purity. Moreover, many real-time datasets are also collected with additional white reference and black reference spectra inside the instrument. These are two reference environments with no light at all, and light is reflected at 99.99%, which can be designed as supplementary components within standard

measurement devices [45]. Compared to datasets measured in laboratory settings, as done previously, this trend of real-world data collection is more practical and is being targeted by food companies.

Therefore, the trend of using handheld or ultracompact devices for direct measurements at the production site, collecting the wide NIR spectrum by many devices, and supplementing reference spectra within the devices, are notable highlights in current food authentication data collection devices. This convenience creates significant opportunities for generating valuable datasets and improving food safety control.

Discussion: Development insights

The using multiple NIR spectral devices trend in food analysis offers opportunities for more comprehensive data collection but also presents challenges due to data heterogeneity. To address this issue and create larger datasets, developing effective data fusion techniques becomes crucial.

In NIRS, these multilevel fusion techniques enable integration and standardization of data from diverse instruments, effectively leveraging data resources from multiple institutions and organizations to build large, valuable datasets without depending on homogeneous equipment. Data fusion, as in [46–48], occurs at three levels: low-level fusion directly combines raw data from multiple sources, mid-level fusion integrates extracted features from different data sources, and high-level fusion combines decisions or interpretations made from separate data sources to conclude. Promising approaches for data fusion include data normalization, selection of important variables, application of advanced machine learning methods such as neural networks or transfer learning, utilization of multi-block data analysis techniques, and development of advanced spectral correction methods. These multilevel fusion techniques enable integration and standardization of data from diverse NIR instruments, effectively leveraging data resources from multiple institutions and organizations to build valuable large datasets without depending on homogeneous equipment, thereby significantly improving the effectiveness of multivariate analysis models in food quality control.

Moreover, compared to using a single handheld device with a narrow spectral range, fusing data from multiple devices to expand the spectral range offers advantages similar to laboratory benchtop with full-range spectra. Specifically, there are more additional spectral information from different spectral regions, facilitating better discrimination and quantification of components in complex samples with high spectral overlap while maintaining the mobility and convenience of handheld devices in practical applications.

4.2 Public datasets

This section reviews publicly accessible NIR spectroscopy food datasets, addressing the limitations of proprietary datasets in developing robust, generalizable machine learning models. These datasets typically include spectral data (often in .csv, .xlsx, or .unsc formats) spanning vari-

ous wavelength ranges, predominantly 900–2500 nm, with some datasets covering both visible and near-infrared regions (Vis-NIR). In NIR spectral archiving, the .unsc file format is a standard format similar to .csv and .xlsx. This is a proprietary file format of The Unscrambler software. To open and process this file, it is necessary to use The Unscrambler or convert it to a more familiar format for machine learning. Corresponding reference values for key food quality parameters are usually provided. While hyperspectral imaging (HSI) offers both spatial and spectral information, its large file sizes and complex processing limit widespread use in portable devices. Therefore, this review focuses on spectral data from NIRS, balancing information richness and practicality for rapid, non-destructive food quality assessment. The datasets cover diverse food products, highlighting the need for larger, more diverse, and standardized spectral databases to advance NIRS applications in food science.

4.2.1 Milk composition in transmittance mode

Collected over eight weeks on a dairy farm, the milk dataset in [45] includes transmittance mode spectra spanning the 960–1690 nm range for 1224 raw milk samples from 41 cows, each with a 2.86 nm/pixel resolution. The dataset incorporates raw milk spectra and white and dark reference spectra used for calibration. With accompanying laboratory reference values for essential milk components such as fat, protein, lactose, urea, and somatic cell count, this dataset goes beyond spectral information by including details like cow ID, milk yield, and time intervals between milkings. Formatted as a .csv file with comprehensive variable descriptions, the dataset aims to facilitate chemometric analysis and the development of multivariate calibration models for predicting milk parameters.

4.2.2 Handheld NIR for chicken breast filets

In a non-destructive manner, portable miniaturized NIR spectrometry captured diffuse reflectance data (908–1676 nm, with an evenly distributed spectral resolution, resulting in 125 variables/measurement) from chicken breast filets in [49]. This data helped differentiate fresh and thawed filets and assess bird growth conditions. NIR measurements were taken from 153 commercial chicken file samples in three modes: direct contact with meat and through the top foil (with or without an air pocket). Thawed samples were generated by freezing and thawing. Multivariate statistics were applied to the 4590 raw NIR spectra.

4.2.3 SpectroFood dataset

The SpectroFood dataset in [50] is a comprehensive hyperspectral meta-dataset aimed at non-destructive estimation of dry matter content across multiple crops. It comprises visible/near-infrared (VIS/NIR) hyperspectral data coupled with corresponding dry matter measurements for

four crops - apples, broccoli, leeks, and mushrooms. In total, 1028 samples were measured using four different calibrated hyperspectral imaging cameras across the spectral range of 398–1717 nm, with all measurements capturing the VIS/NIR range of 470–900 nm. Specifically, 240 apple samples were measured in the 430–990 nm range with 141 bands, 250 broccoli samples in the 470–900 nm range with 150 bands, 288 leek samples in the 398–1717 nm range with 421 bands, and 250 mushroom samples in the 400–998 nm range with 204 bands. The dataset provides the mean reflectance spectrum extracted for each sample in a tabular (.csv) format, along with the corresponding dry matter percentage which ranges from 8.1% to 87%. Additionally, the raw hyperspectral image data for each crop is also provided as .mat files. This multi-crop, multi-sensor dataset aims to facilitate the development of generalized AI/ML models for dry matter estimation that can robustly handle data from different imaging systems and crops.

4.2.4 Reflectance spectral dataset of pre-cooked pasta

This dataset in [51] comprises 1200 Vis-SWIR reflectance spectra (350–2500 nm, with 2151 variables/measurement) of 6 Pennette 72 and 6 Mezze Penne pre-cooked pasta samples with varying salt levels, measured in both frozen and thawed states. The spectra were non-destructively acquired using a portable ASD FieldSpec 4 Standard-Res spectrophotometer, with 50 spectra collected per sample. The data is provided as a .mat file containing a dataset object with rows labeled for sample ID, dry matter content (42.8%, 46.7%, 47.5%), pasta type, and physical state. The averaged spectra highlight differences in the visible region based on salt content, while frozen and thawed samples differed in reflectance intensities across most wavelength ranges, especially 350–1450 nm, 1600–1850 nm, 2100–2400 nm. This annotated Vis-SWIR dataset has valuable reuse potential for developing multivariate classification and regression models to rapidly inspect pre-cooked pasta quality by combining portable spectroscopy and chemometric techniques.

4.2.5 Sugar content measurements of grapes berries in various maturity stage

This dataset in [52] involves 274 samples, each composed of 100 grapes, representing three grape varieties: Syrah, Fer, and Mauzac. The dataset is structured as a CSV file, where rows represent samples and columns include variables such as tray keys, grape varieties, sugar content, and reflectance spectra. Sugar content ranges from 100 to 300 g/L across varieties. Grape sorting was performed using NaCl densimetric baths, followed by hyperspectral acquisition. A total of 274 reflectance spectra were obtained, covering red (Syrah, Fer Servadou) and white (Mauzac) grape varieties.

4.2.6 Mango fruits

The dataset in [53] comprises 186 NIR spectra (1000-2500 nm, $\log(1/R)$ absorbance) of intact mangoes from 4 cultivars, acquired using FT-NIR with 64 coadded scans/sample. Raw spectral data in .xls and .unsc formats. Reference data includes vitamin C (mg/100g), soluble solids ($^{\circ}$ Brix), and total acidity (mg/100g). This dataset enables the development of prediction models for rapid, non-destructive quality evaluation of whole mangoes using NIR spectroscopy.

4.2.7 Enhanced NIR spectra of intact mangoes

This dataset in [54] provides original and enhanced near-infrared (NIR) spectral data (1000-2500 nm, 1557 wavelength variables) of 58 intact Kent mango samples. The spectra were acquired using a Fourier transform NIR spectrometer, with 32 coadded scans per sample. The raw absorbance spectra were enhanced using algorithms like multiplicative scatter correction (MSC), baseline linear correction (BLC), and their combination MSC+BLC. The original and enhanced spectral data in .unsc and .xlsx formats for predicting two key internal quality traits - total acidity (TA) and vitamin C content. Model performances were evaluated against reference TA and vitamin C values measured by standard methods, using metrics such as coefficient of determination (R^2), correlation (r), root mean square error (RMSE), and residual predictive deviation (RPD).

4.2.8 Cocoa beans

The dataset in [55] contains NIR absorbance spectra (1000-2500 nm) with 32 co-added scans at 0.2 nm resolution for a total of 72 bulk samples of intact cocoa beans, with each sample amounting to 50g. The spectra data is provided in both .xlsx and .unsc file formats. For each sample, the actual moisture content (%) and fat content (%) were measured using standard laboratory methods like thermogravimetry and Soxhlet extraction, respectively. The measured moisture content ranged from 6.74% to 12.08% with a mean of 9.04%, while the fat content ranged from 35.26% to 45.75% with a mean of 40.32%.

4.2.9 Vis-NIR spectra for sugarcane across multiple spectrometers

This dataset in [56] provides Vis-NIR absorbance spectra and corresponding chemical reference data for 60 sugarcane samples, which were analyzed using 8 different spectrometers. These include one laboratory spectrometer (LabSpec 4) and seven micro-spectrometers (NIRscan Nano, F750, MicroNIR1700, MicroNIR2200, NIRONE 2.2, SCIO, TellSpec). The spectral ranges covered span from 350-2500 nm for the LabSpec 4, while the micro-spectrometers capture narrower ranges, such as 1750-2150 nm for the NIRONE 2.2 device. The reference chemical

data encompasses total sugar content (ranging from 1.1-51% dry matter), crude protein content (0.9-9.6%), acid detergent fiber (26-59.3%), and in-vitro organic matter digestibility (13-66.6%). This open-access dataset facilitates comparing prediction performance across the various spectrometers employed.

4.2.10 Enhanced Vis/NIR spectral dataset of intact Cucurbitaceae fruits

The dataset in [57] comprises Vis/NIR absorbance spectra (381-1065 nm) of 300 samples from 6 Cucurbitaceae fruit types, including zucchini, bitter melon, ridge melon, chayote, and cucumber. The spectra were acquired using a NirVana AG410 portable spectrometer, with each sample scanned 6 times. The data is provided in .xls and .unsc formats. Reference data on soluble solids and water content were determined by standard wet chemistry methods.

Discussion: Development insights

Analysis of existing public NIR datasets in the food industry reveals a clear trend: prioritizing nutritional indices and product quality. This trend reflects economic development goals through enhancing nutritional value and optimizing production processes. However, a notable gap exists - the relative scarcity of data related to food safety factors, particularly in identifying chemical toxins and other hazards. Based on the above survey, to the best of our knowledge, there are currently no relevant public NIR datasets. This could be due to some factors: the high cost of preparing absolutely safe samples for reference, the cost of chemicals for testing and measuring unsafe levels, the sensitivity of data related to food contamination, and the technical challenges of detecting low concentrations of substances using NIR methods.

The first approach to closing this gap is to promote the creation, development, and sharing of NIR spectral datasets through active partnerships with food safety regulatory bodies. Through one project supported by the People's Committee, we are currently collecting NIR spectral data in collaboration with regional food safety authorities to build comprehensive datasets covering 19 common chemical residues found in daily food products, following Ministry of Health standards. This NIR data is validated through independent testing using traditional chemical reference methods. This data collection strategy aligns with regulatory requirements and enables standardized datasets for market surveillance while advancing ML applications in food quality.

The second approach is to improve small food safety datasets, creating richer, larger-scale, and more valuable datasets. Methods such as data enrichment with Generative Adversarial Networks (GANs) and spectral diffusion models that can generate synthetic spectra are untapped potential directions. This approach can increase both the quantity and quality of available data, partly addressing the difficulty of high cost in collecting spectral data related to chemicals and chemical residues in food. This will support the

development of comprehensive, rapid, and non-destructive analytical methods, bringing double benefits: economic development in parallel with consumer health protection.

In these recent datasets containing over 1000 spectra, techniques are mainly in traditional ML, such as PLS for [50] in [58], PLS-DA, SVM, ANN or combine methods Random Subspace Discriminant Ensemble (RSDE) for [49] in [59]. In [59], RSDE demonstrated superior performance with over 95% classification accuracy. Its innovative ensemble architecture combines multiple submodels through random subspace projection and majority voting, delivering enhanced accuracy and reliability while inherently reducing noise sensitivity and overfitting risks. The trend across studies shows a shift from single methods to ensemble and hybrid approaches for handling complex spectral data. Another instance, a milk composition study using NIR transmittance spectra [45] showed that combining SO-PLS with appropriate preprocessing improved prediction accuracy by 5-25%. While these studies demonstrated the effectiveness of hybrid approaches, they examined different food products under varying conditions, making direct comparisons challenging. Notably, DL approaches were not extensively explored, suggesting potential opportunities for investigating modern architectures like CNNs and transformer models that have proven effective in other computer vision and spectral analysis applications. Further research exploring traditional and DL techniques across standard datasets would be valuable to establish their relative effectiveness and generalizability.

4.3 Pre-processing and wavelength selection

While deep learning has developed increasingly, traditional ML models still dominate in NIR tasks. In this context, data preprocessing and feature selection are two commonly exploited factors, improving the quality of input data, reducing data dimensionality, extracting important information from NIR spectra, and enhancing the performance of models. This review focuses on the recent trends and notable developments in data preprocessing and feature selection for NIR spectral analysis in the last few years.

The identification of recent advances in preprocessing and wavelength selection methodologies was conducted systematically and compared with established baseline methods (as in Figure 4).

4.3.1 Pre-processing

Pre-processing of NIR data is important to improve model performance. These common techniques have been shown in [25, 44], mainly including noise reduction, baseline correction, resolution enhancement, centering, smoothing, derivative, de-trending, and scaling methods. Most studies apply only one or two pre-processing methods, chosen based on experience rather than a systematic assessment of optimal methods. This is a highly complex process that requires expert knowledge to select the most appropriate pre-

processing method for each specific dataset. Choosing the wrong pre-processing method can lead to the loss of important information or add more unwanted noise, thereby negatively affecting model performance. Therefore, developing a systematic approach to evaluate and select optimal pre-processing techniques for NIR spectrum data is an important direction for future research.

Recently, the trend in NIR spectral data pre-processing has seen significant advancements, particularly towards automation and optimization of procedures. Methods such as Synergy Adaptive Moving window algorithm based on the Immune Support Vector Machine (SA-MW-ISVM) [62], Automatically generating pre-processing strategy (AgoES) [65], and Sequential preprocessing through ORThogonalization (SPORT) [66] have been developed to automate the selection and combination of optimal pre-processing methods, thereby reducing manual intervention and experimentation time. Although the Self-expansion Full Information Optimization strategy (SFIOS) [57] is not entirely automatic in pre-processing selection, it nevertheless provides a comprehensive optimization strategy that includes pre-processing, as in Figure 5. Simultaneously, the trend of combining multiple pre-processing methods, as in Table 3, such as Savitzky-Golay (SG) with Standard Normal Variate (SNV), has become popular to leverage the advantages of each method. Researchers in [65] developed strategies to find optimal pre-processing pipelines, evaluating 150 different combinations of 16 common techniques. Similarly, another study [64] explored an automated method to combine up to 4 out of 9 pre-processing types for NIR data of coconut milk. Using various models (Multi-Layer Perceptron (MLP), k-Nearest Neighbors (KNN), and Partial Least Squares (PLS)), they achieved the best results with KNN on Micro-NIR data, obtaining a classification accuracy of 0.97 and a regression RPD of 16.108. Both studies [64, 65] emphasize the importance of automated optimization in pre-processing pipelines, as no single method consistently outperforms others across all scenarios.

Notably, there is a close integration between pre-processing and machine learning algorithms, as demonstrated in SA-MW-ISVM, where both pre-processing and SVM parameters are optimized simultaneously. Similarly, ensemble learning techniques are widely applied, as in AGoES and SPORT, to combine various pre-processing models. Additionally, many methods focus on selecting relevant spectral variables, thus contributing to dimensionality reduction and improving model performance.

Not limited to a single data type, new methods have been developed with adaptability for various data types, including both solid and liquid data, as well as spectral types beyond NIR. Moreover, the trend of sharing open-source pre-processing algorithms is increasing, thereby enabling the research community to continue developing and improving these methods. Interestingly, while many methods focus on large datasets, some approaches like Extended Multiplicative Signal Augmentation (EMSA) emphasize improving performance on smaller datasets, meeting the needs of

Table 3: Preprocessing methods and their applications

Preprocessing methods	Single	Combine	Remarks
SG smoothing, MSC, SNV, 1st Der, 2nd Der [60]	x	Multiple combinations of SG with other methods	SG alone or combined with SNV gave the best results for most models; Improved classification accuracies compared to raw spectra; SG+SNV optimal for predicting most physicochemical properties; Preprocessing crucial for developing robust NIR models for melon seed powder authentication
SM, DF, NM, CT, DE [61]	x	Multiple schemes from multiple methods	SFIOS auto-optimization strategy combines methods, and provides statistical info; en-iViSSA ensemble improves model performance; DF highlights spectral info; CT and DE good for solid samples; MSC and SNV minimize scattering effects; Open source
SG smoothing, SG derivatives, MSC, SNV, Autoscale, Normalization [62]	x	Multiple schemes from 6 methods	SA-MW-ISVM algorithm optimizes preprocessing, wavelength selection, and SVM parameters simultaneously; Improves prediction accuracy by up to 44% compared to PLS; Selects relevant wavelengths, reducing variables to ~30% of full spectrum; Chooses appropriate preprocessing combinations; Applicable to NIR and other spectroscopy data
SG, EMSC, EMSA, Batch Aug [63]	x	Multiple schemes from 4 methods	EMSA can replace pre-processing for CNNs; Combination of pre-processing and augmentation improves results for conventional classifiers on small datasets; CNNs benefit from traditional pre-processing; Augmentation especially beneficial for small datasets
BSO3, 1st Der, 2nd Der, SNV, MSC, MS, SG filter, SS [64]	x	Multiple schemes from 9 methods	Automatic strategy combines preprocessing and ML hyperparameter tuning; Improves classification (up to 98% accuracy) and regression (RPD up to 16) for coconut milk adulteration detection using FT-NIR and Micro-NIR
Baseline correction, Scatter correction, Scaling [65]	x	150 combinations from 16 single methods	AGoES automatically generates and evaluates all preprocessing combinations; Different optimal combinations are found for each ML algorithm and property predicted; Improved model performance compared to raw spectra for most cases; SVM with AGoES preprocessing performed best overall
SG smoothing/derivatives, SNV, VSN [66]	x	Multiple combinations via SO-PLS	SPORT method automatically combines and selects optimal preprocessing sequence; Outperformed single preprocessing and stacking approaches; Selected parsimonious combinations of 2-3 methods; Order of combination had minor impact on performance

many practical applications.

Discussion: Development insights

Effective preprocessing depends on the type of data. However, the variety of handheld measurement devices available today, as discussed in Section 4.1, and the rapid development of deep learning raise questions about the future role and necessity of data preprocessing. Possible research directions include: (1) more flexible preprocessing automation, adapting to different types of data, serving complex multi-class classification problems in practice, such as classifying unsafe fruits and vegetables in

food safety inspection, (2) developing preprocessing methods suitable for many different measuring devices, aiming at a model that can transfer technology between localities with asynchronous devices, and (3) evaluating if the need for preprocessing, with fluctuating data in different measuring environments, such as markets and supermarkets, where humidity, light, and temperature can all affect measurement values, thereby affecting the effectiveness of machine learning models.

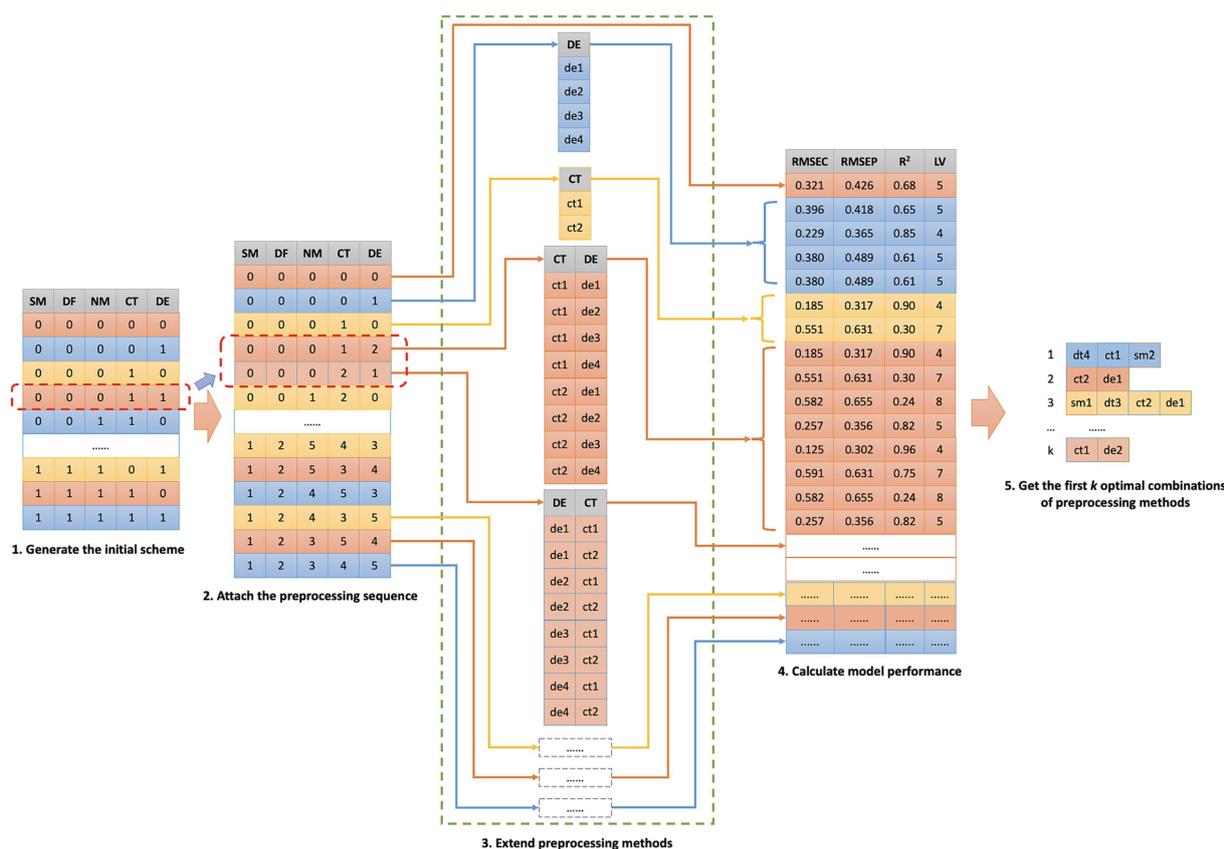


Figure 5: The optimal preprocessing scheme of SFIOS [61]

4.3.2 Wavelength selection

Wavelength selection is critical for NIR analysis. By identifying the most informative spectral regions, it simplifies complex data, improves model performance, and focuses on key information, as in Table 4, 5. These methods are categorized by their approach and can be broadly grouped as follows:

1. Filter methods utilize statistical criteria to evaluate the relevance of individual wavelengths to the target variable, including Variable Importance in Projection (VIP), Correlation-based Feature Selection (CFS), Relief, Fisher score, and Chi-squared test;
2. Wrapper methods assess the quality of wavelength subsets based on the performance of a prediction model, such as Genetic Algorithm (GA), Successive Projections Algorithm (SPA), Recursive Feature Elimination (RFE), Sequential Forward Selection (SFS), and Sequential Backward Selection (SBS);
3. Embedded methods integrate feature selection within the model-building process, including Least Absolute Shrinkage and Selection Operator (LASSO), Elastic Net, Ridge Regression, Random Forest, and Competitive Adaptive Reweighted Sampling (CARS);

4. New hybrid methods combine the strengths of different approaches, such as Multi-Feature Extraction combined with LASSO (MFE-LASSO), Maximal information coefficient - Successive projections algorithm combined with Extreme learning machine - Genetic algorithm (MIC-SPA-GA-ELM), and Beluga whale optimization with Iterative variable subset optimization (BWO-IVSO).

Each method has its strengths and limitations, making it suitable for specific data types and analytical objectives. Filter methods efficiently screen large variable spaces to rank features. For instance, VIP was used to select 13 key wavelengths out of 209 initial variables for adulterant detection in quinoa flour [67], improving the R_p^2 from 0.94 to 0.98 and reducing RMSEP from 3.04% to 1.60%. Wrapper methods, such as GA and SPA, directly optimize subsets based on model performance. In a study on durian fruit quality assessment [19], GA selected 23 wavelengths for dry matter prediction and 19 for total soluble solids, improving the model's accuracy. Hybrid methods combine filter and wrapper approaches for robust performance, as seen in the MIC-SPA-GA-ELM [68] combination used for tobacco and corn samples, which showed the best accuracy and robustness. However, most techniques retain spectroscopic variables related to key functional groups and structural chemistry to develop broadly applicable, physically

interpretable models, as in Table 5.

Wavelength selection can also be divided into interval or peak selection, as in Figure 6. Peak selection involves choosing specific, individual wavelengths that are most informative. For example, four peaks (1428, 1704, 1892, 1912 nm) were identified as crucial in a vineyard water status prediction study [69]. On the other hand, interval selection chooses continuous ranges of wavelengths. In the same vineyard study, three intervals (1402-1508, 1676-1750, 1870-1926 nm) were selected. This approach can be particularly useful when certain regions of the spectrum are known to be associated with specific molecular structures or properties of interest.

The current trend is to combine multiple variable selection methods to optimize results. For instance, the BWO-IVSO approach applied in aflatoxin B1 analysis in peanuts significantly improved model performance compared to using the full spectrum [72]. Another example is the two-step approach using RRelief and MIC, followed by Elastic Net, for azodicarbonamide detection in wheat flour [74]. Besides, recent studies also emphasize the importance of optimizing model parameters. Algorithms like Harris Hawks Optimization (HHO) and Rime Optimization Algorithm (RIME) have been used to fine-tune parameters of models such as Kernel-based Extreme Learning Machine (KELM), significantly improving model performance [78].

The number of selected wavelengths or spectral regions typically varies based on the complexity of the sample and the specific analytical objectives. For instance, 18 wavelengths were selected for sugar in tobacco leaves, while 24, 34, 26, and 16 wavelengths were chosen for moisture, oil, protein, and starch in corn, respectively [68]. In some cases, a few carefully selected wavelengths can provide comparable or even superior results to models using the full spectrum, while significantly reducing computational requirements. **Discussion: Development insights**

Wavelength selection is a popular method in NIR spectral processing. This is based on the interpretability of the results, which are related to the functional groups representing the sample. Although this method is effective, as with preprocessing, whether wavelength selection is necessary or using the raw data itself with deep learning architectures is more effective needs to be further compared.

Combining both methods is also a possible direction for development. This method takes advantage of both the non-linear learning capabilities of the neural network and maintains the interpretability of the final model.

4.4 Advanced NIR food spectral analysis techniques

4.4.1 Traditional machine learning

Traditional ML methods play a crucial role in NIR spectral analysis, focusing on addressing multicollinearity issues and improving generalization capabilities. In a typical pipeline, these techniques involve data pre-processing,

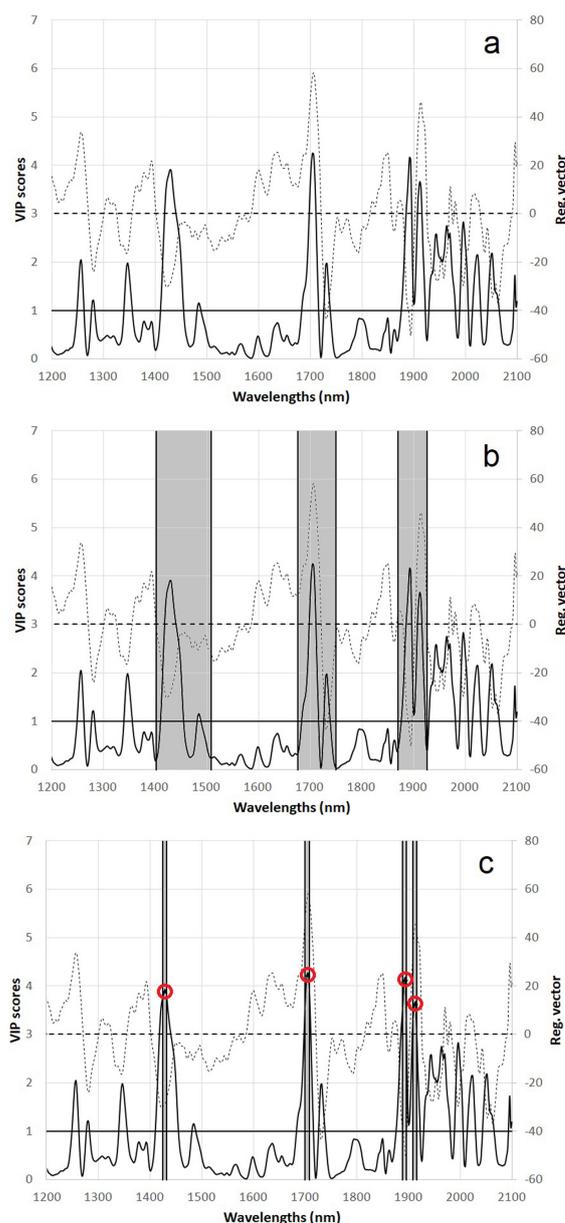


Figure 6: Wavelength selection methods comparison: (a) Manual VIP-based selection after preprocessing, (b) Interval selection approach, and (c) Peak selection method highlighting key wavelengths and bandwidths [69]

feature selection/ extraction, and applying traditional ML algorithms to model the selected features and generate outputs. Popular techniques include PCA, PLS, ELM, SVM, SVR, DT, and RF. These methods concentrate on extracting important features, minimizing data redundancy, and building effective predictive models for various NIR applications. The recent highlights of traditional ML on NIR spectra largely lie in the improvements in pre-processing strategies as well as wavelength selection, and the effective feature extraction methods mentioned earlier. However, traditional methods often face limitations in

Table 4: Wavelength selection methods for NIR spectra (1)

Task	Dataset	Variable Selection Methods	Selected Wavelengths	Results	Remarks
Quinoa flour adulteration [67]	54 samples, 941-1674 nm, 209 vars	VIP	13 wave-lengths	Initial (209 vars): $R_p^2=0.94$, RMSEP=3.04%. Selected (13 vars): $R_p^2=0.98$, RMSEP=1.60%	PLSR model improved for quinoa flour adulteration prediction
Vineyard water status [69]	288 samples, 1200-2100 nm, 501 vars	Interval, Peak, IPLS	3 intervals: 1402-1508, 1676-1750, 1870-1926 nm. 4 peaks: 1428, 1704, 1892, 1912 nm	Initial (501 vars): $R_p^2=0.84$, RMSEP=0.167 MPa. Selected (9-33 vars): $R_p^2=0.77-0.78$, RMSEP=0.186-0.201 MPa	3 methods for key wavelengths, simplified model, comparable accuracy
Durian quality [19]	278 samples, 860-1750 nm	SPA, GA, VIP	GA: DM 23, TSS 19 wave-lengths	Full: $R^2=0.83$, RMSEP=4.96% (DM); $R^2=0.81$, RMSEP=3.71% (TSS). Selected: $R^2=0.85$, RMSEP=4.50% (DM); $R^2=0.66$, RMSEP=5.15% (TSS). Accuracy: 94.20%	3 methods improved DM and TSS prediction, GA best
Robust NIR model [68]	Corn: 80, 1100-2498 nm	LARS, CARS, SPA, UVE, MIC-combined	24 (moisture), 34 (oil), 26 (protein), 16 (starch)	MIC-SPA-GA-ELM: Best accuracy, robustness	Combined methods improve model accuracy, stability
S-ovalbumin in eggs [70]	150 samples, 900-1700 nm, 390 vars	SPA, IRIV	SPA: 16, IRIV: 14	PLSR (16): $R_c^2=0.90$, RMSECV=8.92%, $R_p^2=0.84$, RMSEP=9.98%. PLSR (14): $R_c^2=0.91$, RMSECV=8.44%, $R_p^2=0.86$, RMSEP=9.33%	IRIV-selected (14) PLSR best for S-ovalbumin prediction
Adaptive PLS [71]	Corn: 80, 2498-1100 nm. Wine: 44, 900-5000 cm^{-1}	CARS, MCVUE, LARS, ABUSE	Corn: 4 regions. Wine: 4 wavelengths	ABUSE PLS: Corn (25): RMSECV=0.003, RMSEP=0.004. Corn (26): RMSECV=0.009, RMSEP=0.013. Wine (4): RMSECV=5.12, RMSEP=4.66. Wine (3): RMSECV=5.13, RMSEP=4.63	ABUSE selects key peaks, best performance, fewer variables, improved accuracy, reduced time
Aflatoxin B1 in peanuts [72]	100 samples, 955-1702 nm, 128 vars	IVSO, BWO-IVSO	IVSO: 32, BWO-IVSO: 18	SVM (full): RMSEP=31.4602, $R_p=0.9608$, RPD=3.6799. SVM (IVSO): RMSEP=30.4587, $R_p=0.9633$, RPD=3.8009. SVM (BWO-IVSO): RMSEP=24.6322, $R_p=0.9761$, RPD=4.6999	BWO-IVSO removes redundancy, noise, enhances AFB1 analysis accuracy

learning complex and nonlinear features when dealing with complex and high-dimensional NIR spectral data. This leads to the need for deep learning approaches, which can automatically extract complex features and efficiently pro-

cess high-dimensional data, thereby improving accuracy and generalization capabilities in many NIR spectral analysis tasks.

Table 5: Wavelength selection methods for NIR spectra (2)

Task	Dataset	Variable Selection Methods	Selected Wavelengths	Results	Remarks
Dual-sPLS NIR [73]	Rice: 447, 12481-3595 cm^{-1} , 1153 vars	PLS, iPLS, SiPLS, mw-PLS	149 optimal vars	SiPLS (149 vars): RMSEP reduced 0.2284 to 0.1952	Variable selection improves model, saves computation time
ADA in wheat flour [74]	101 samples, 0-300 mg/kg, 7012 vars	RReliefF, MIC, EN	Two-step: 500, then 40	PLSR (7012): RMSEP=2.53%, $r=0.975$. PLSR (500 MIC): RMSEP=1.32%, $r=0.992$. PLSR (40 MIC+EN): RMSEP=0.78%, $r=0.997$	MIC+EN eliminates irrelevant vars, retains key info, improves accuracy
Talcum in wheat flour [75]	123 samples, 1050 vars	EN, GA	EN+GA: 55	EN+GA (55/1050): GBDT: $R^2=0.9778$, RMSEP=0.8905, RPD=6.8099	Detects low talcum concentrations in wheat flour
GBM-PLS for corn [76]	120 samples, 7 countries, 867-2535 nm, 949 vars	RC, CARS, XGBoost, LightGBM, CatBoost	Moisture: 6 (CatBoost). Protein: 6 (LightGBM)	Best: Moisture - $R_V^2=0.97$, RMSEV=0.45%, RPDV=6.20. Protein - $R_V^2=0.82$, RMSEV=0.51%, RPDV=2.41	GBMs good for wavelength selection. CatBoost best for moisture, LightGBM for protein. SHAP identified key wavelengths
Strawberry SSC [77]	630 samples. Reflectance: 600-1080 nm (949). Transmittance: 600-950 nm (805)	SI, SPA, UVE, CARS	Transmittance (CARS): 33	Best (Transmittance CARS-PLS): $R_p=0.928$, RMSEP=0.412 $^{\circ}\text{Brix}$, RPD=2.670	Transmittance with CARS best. 3 strawberries/sec. More research needed
Zearalenone in wheat (CSA-NIR) [78]	131 samples, 901-1701 nm, 228 vars	VCPA, BOSS, CARS	CARS: 107 (best)	Best (CARS-RIME-KELM): $R_p^2=0.9900$, RMSEP=18.4610 $\mu\text{g/kg}$	CARS best. RIME improved KELM. CSA-NIR effective for zearalenone detection
Zearalenone in wheat (FT-NIR) [79]	116 samples, 10,000-4000 cm^{-1} , 3112 vars	CARS, SVM-RFE, MFE-LASSO	MFE-LASSO: 38 (best)	Best (MFE-LASSO-PLS): $R_p^2=0.9545$, RMSEP=18.6442 $\mu\text{g/kg}$, RPD=4.3198	MFE-LASSO best. FT-NIR+MFE-LASSO-PLS effective for zearalenone detection

4.4.2 Deep learning architecture

In NIRS analysis, DL architectures have demonstrated great potential in enhancing the efficiency and accuracy of analytical processes. Each architecture possesses unique operational mechanisms and advantages suitable for different challenges in NIR analysis, as shown in [25]. Stacked Autoencoders (SAE) excel at learning low-dimensional representations of input data, effectively reducing noise and focusing on essential information. Variational Autoen-

coders (VAE) function as generative models capable of producing new samples, proving valuable for data augmentation. CNN are adept at learning local features of NIR data, while RNN process NIR data as time series, capturing sequential relationships. ELM show high efficiency in scenarios with limited training samples, and GAN can generate new training data, addressing data scarcity issues. A thorough understanding of the mechanisms and advantages of each architecture enables researchers to select the most appropriate method for specific applications in NIR analy-

sis, ultimately leading to more accurate and reliable results in this crucial field of spectroscopy. In this study, we focus only on recent notable DL research and the salient value points not discussed in previous review articles. In particular, papers that have comparisons with traditional ML are prioritized, as in Table 6, 7, 8.

Firstly, DL models achieve superior performance over traditional ML methods in NIRS analysis through their innovative architectures such as automatic hierarchical feature extraction, attention mechanisms, deep temporal learning via recurrent networks, and intelligent feature fusion through dense connections, enabling more accurate classification, regression, and anomaly detection tasks. Previously, many DL studies were done without comparison with traditional methods that have achieved high performance. This has received more attention recently. These results confirm the potential of DL in enhancing performance for regression, classification, and anomaly detection tasks in NIR spectral analysis, opening up possibilities for wide-ranging applications across fields requiring high accuracy and reliability.

In classification tasks, for instance, Convolutional Neural Network-Attention (CNN-ATT) achieved 100% accuracy in categorizing chickpeas into HTC and ETC classes [84], surpassing traditional SVM models. In this study, CNN-ATT enhances performance through its attention block mechanism that dynamically weighs and focuses on the most relevant features in input data, allowing the network to adaptively prioritize important spectral information while filtering out less significant signals. Similarly, the Transformer in [92] achieves remarkable classification accuracy (99.31%) through its three-layer encoder and multi-head attention mechanism that effectively extracts semantic information from both vibration and Vis/NIR spectral data while dynamically adjusting feature weights via an attention feature fusion module to focus on the most relevant information for apple moldy core detection, better than PLS-DA, SVM, ELM. In another study, for corn variety recognition [87], CNN reached 99.2% accuracy, outperforming traditional methods like KNN, SVM, and PLS by 25.78%. In a more complex scenario of identifying adulterated beef and mutton [85], ResNet with 2DCOS, as in Figure 7, achieved 100% accuracy, significantly surpassing PLS-DA's 32–50% accuracy. A total of 1,878 synchronous and asynchronous 2D-COS spectra were obtained from transforming 1D spectra across 23 diverse adulteration patterns (5 pure meats and 18 mixed samples with varying proportions of 25%, 50%, and 75%) into 2D images to enhance resolution and analytical sensitivity while providing multi-dimensional information through auto and cross-correlation peaks, enabling accurate detection of both components and mixing ratios in meat samples, with characteristic markers at different wavelengths. Besides, ResNet achieves superior performance through its innovative skip connection technique that allows direct data flow between layers, effectively preventing gradient vanishing and enabling faster, more accurate training than traditional CNNs when handling this com-

plex multi-class classification problem.

For regression tasks, DL models consistently showed superior performance. For instance, in predicting lead content in oilseed rape [81], the Transfer Stacked Auto-encoder (T-SAE) model achieved R^2 values of 0.9215 and 0.9349 for leaves and roots respectively, outperforming PCA-SVM and SAE. A highlight of this research is that T-SAE achieves high performance through its dual-model transfer mechanism, where network weights are initialized from pre-trained SAE models while allowing deep feature layers to be trained from scratch with random weights, effectively combining information from both leaf and root spectral data to achieve superior classification accuracy (98.75%) in lead stress detection. In another study, for estimating soluble solids content in pears [83], SpectraNet-32 achieved the best results, surpassing classical methods like PLS, MLR, and SVM. In predicting cooking time for chickpeas [84], 1D-CNN also outperformed traditional regression methods. Regarding anomaly detection and complex analysis, DL models also excelled. For pesticide residue recognition on garlic chive leaves [90], 1D CNN achieved 97.9% accuracy, outperforming traditional models. In analyzing complex organic compounds [89], the proposed DL model achieved R^2 values between 0.9574 and 0.9996, improving upon PLSR and BPNN by significant margins. This architecture achieves superior dynamic feature extraction through its innovative dual-module design, as in Figure 8, where the short-term feature extraction utilizes multi-rate dilated convolutions with dense connections to capture short-term spectral patterns while the long-term feature extraction employs Gated Recurrent Unit (GRU) enhanced by temporal attention mechanism to comprehensively merge features across all timesteps, complemented by a linear bypass path and two-stage quality regression approach that effectively prevents overfitting in NIR spectral analysis.

Thus, the outstanding advantages of DL architecture, such as automatic feature extraction, attention mechanisms, multi-scale temporal feature learning through dilated convolutions, comprehensive time-series analysis via GRU networks, and intelligent fusion through dense connections, have led to superior performance compared to traditional ML approaches.

Second highlight, DL architectures demonstrate high flexibility, successfully applied to various data types ranging from 1D spectra to 2D images, 2D correlation spectra (2D-COS), and even 2D dynamic data. These studies highlight the versatility of DL architectures in NIRS, effectively processing various data formats from simple to complex data, and even enabling advanced techniques like transfer learning across different devices. This flexibility opens up the potential for developing architectures in computer vision on NIR spectral data transformed into image data formats.

There are many methods to convert 1D NIR spectra into 2D. First of all, the 2D-COS technique [82, 91], as in Figure 9 transforms one-dimensional NIR spectra into two-dimensional correlation spectra (synchronous and asynchronous) by calculating cross-correlations between spec-

Table 6: Advanced DL Architectures for NIR food spectral analysis (1)

NIR Tasks	Datasets	Pre-processing	Models	DL Architecture	Results
ADF and IVOMD in sugarcane [80]	60 NIR x 3 devices, 600/device, 3:1:1 split	WS, Interpolation, SNV	1D-Inception-ResNet, PLS	8 Conv, 4 FC, Residual, Softplus, Dropout 0.15	ADF/IVOMD: $R^2 > 0.96$, RMSEP < 2.75 . <i>Outperformed PLS, Successful inter-device transfer</i>
Lead content in oilseed rape [81]	500 samples (leaves/roots), 3:1:1 split (480.46-1001.61 nm)	SNV, 1st Der, 2nd Der, PCA	T-SAE, SAE, SVM, SVR, PCA-SVM	Best T-SAE: 411-148-108-60 (leaves), 410-140-91-56 (roots)	$R^2 = 0.9215$, RMSEP = 0.0302 mg/kg (leaves); $R^2 = 0.9349$, RMSEP = 0.0278 mg/kg (roots) <i>Outperformed PCA-SVM, SAE; successful transfer learning</i>
Total phenolic content in boletes [82]	187 samples (3 species), 90% model, 10% valid	SNV, HCA, Folin-Ciocalteu	ResNet, 2D-COS, SVM, PLS-DA	12-layer ResNet, identity & conv blocks, BatchNorm, ReLU	100% accuracy (train & test). <i>Outperformed traditional methods, rapid & non-destructive</i>
SSC and temp in "Rocha" pear [83]	3300 spectra (1650 pears), 499.73-1101.83 nm, 5 valid sets	QNV, Savitzky-Golay (1st & 2nd), PLS wrapper (BVE-PLS), PLS-VIP	SpectraNet-32/53, DeepSpectra, PLS, MLR, SVM, MLP	SpectraNet-32: 32-layer ResNet, 3 Residual Units, BatchNorm, GELU, Global Avg Pooling, Dropout	Best: RMSEP = 1.08%, $R^2 = 0.58$ (SSC). <i>Outperformed classical methods, predicted SSC & temperature, 8000 spectra/s</i>
Chickpea HTC/ETC classification, Cooking time prediction [84]	864 seeds (8 varieties), 900-2500 nm	SNV, 1st and 2nd derivatives; CARS, IRIV, CNN-FS	PLSDA, SVC, CNN-ATT, 1D-CNN	CNN-ATT: ATT block, 3 1D conv, 2 dense, SoftMax. 1D-CNN: 1 input, 1 conv, 4 dense, 1 output	SVC & CNN-ATT: 100% acc (full spectrum). 1D-CNN: $R_p^2 = 0.880$, RMSEP = 0.662 (cooking time) <i>Non-destructive, rapid detection. Effective for both classification and prediction</i>
Adulterated beef/mutton ID [85]	1878 samples, 0-100% adulteration, 400-2500 nm	Raw, FD, SD, MSC-SG	PLS-DA, ResNet with 2DCOS	ResNet: 12 hidden layers, ReLU, global avg pooling	ResNet+2DCOS: 100% acc. PLS-DA: 32-50% acc. <i>2DCOS enhances spectral resolution. ResNet extracts 2DCOS features effectively.</i>
Subsurface bruises in plums [86]	1125 HSI, 430-1000 nm	Standardization, Data augmentation; PCA (10 wavelengths)	HSCNN, ResNet, 3D-CNN, PLS-DA	HSCNN. ResNet: Adapted for HSI. 3D-CNN: 3 conv, 3 maxpool, 3 batch norm, global avg pool, 2 dense	Best: HSCNN (full spectrum), F1 90%. 3 wavelengths: F1 89%. <i>Detected invisible bruises. Reduced to 3 wavelengths with similar performance</i>

Table 7: Advanced DL Architectures for NIR food spectral analysis (2)

NIR Tasks	Datasets	Pre-processing	Models	DL Architecture	Results
Corn variety recognition [87]	450 NIR (5 varieties), 2:1 split, 11542-3940 cm^{-1}	DT for baseline drift; CARS (114/1845 wave-lengths)	CNN-LeNet-5, BP, KNN, SVM, PLS	3 Conv: 32, 64, 128 kernels (7, 4, 4 windows). 3 Pool: Max & Glob. Avg. ReLU, Dropout 0.1, FC, softmax	CNN: 99.2% test acc. <i>Combining NIR, CNN enables accurate, rapid recognition; 25.78% higher accuracy than traditional models</i>
Watermelon SSC [88]	1440 Vis-NIR (317-1117 nm), 60:30:10 split	PLSR: SG+2nd der. BPNN: MSC	PLSR, BPNN, 1D-CNN	1D-CNN: 5 Conv1D, 3 MaxPool, Flatten, 3 Dense, BatchNorm, Dropout 0.1	1D-CNN: $R_p^2 = 0.97$, RMSEP=0.21. +14.1% vs PLSR, +6.6% vs BPNN <i>High R_p^2, Low RMSEP. Features at 720, 810nm</i>
Quality of complex organics [89]	2D NIR, 1408 spectra (12500-3950 cm^{-1}), 844:432:564 windows	SG (NIR). SG+2nd der for PLSR. MSC for BPNN	PLSR, BPNN, Proposed DL	MDFE: SDFE (3 dilated 2D CNN) + LDFE (GRU with attention). Regression: FC (1024, 500), ReLU, BatchNorm, Dropout 0.1	DL: $R^2=0.9574-0.9996$, RMSE=0.0013-0.4374. +14.1% vs PLSR, +6.6% vs BPNN. <i>Short/long-term dynamics, Dilated CNN extracts multi-level short-term features; Temporal attention redistributes GRU features; Nonlinear fitting of complex NIR-quality mapping; Visualization shows model extracting relevant spectral bands</i>
Pesticide on garlic chives [90]	SWIR HSI, 30 leaf spectra, 920-1700nm, 90:5:5 split	Mean filter (SNR), Isolated Forest (outliers)	PLS, MLR, BPNN, KNN, LDA, NB, RF, SVM, 1D CNN	1D CNN: 3 conv (3x1, stride 2x1, pad 1), avg/max pool, flatten, 2 FC (256, 128), 4-node output, ReLU, BatchNorm	97.9% test acc. Recall > 97.7%, AUC > 0.99. 0.208 Hamming loss (mixed). vs KNN (91.2%), LDA (77.5%), NB (41.4%), RF (90.9%), SVM (92.8%). <i>Non-destructive, Rapid, Outperformed traditional models, Exploiting pixel-wise spectra for a large dataset; Successful mixed residue identification</i>

tral intensities $\tilde{y}(v_1)$ and $\tilde{y}(v_2)$ at wavelengths v_1 and v_2 as external perturbations (like geographical origins, storage time) change. The synchronous spectrum reveals peaks that change in the same direction through auto-correlation peaks (diagonal) and cross-correlation peaks (off-diagonal), while the asynchronous spectrum only shows cross-peaks indicating spectral changes from different molecular sources. This technique is particularly valuable for NIR spectral analysis

when spectral peaks overlap due combinations, and complex molecular interactions in food samples, as it improves spectral resolution and helps distinguish overlapping features that are not observable in one-dimensional spectra. Unlike 2D-COS method which focuses on molecular interactions through correlation analysis, the Gramian Angular Difference Field (GADF) [93] converts one-dimensional NIRS into two-dimensional images by first normalizing

Table 8: Advanced DL Architectures for NIR food spectral analysis (3)

NIR Tasks	Datasets	Pre-processing	Models	DL Architecture	Results
Wolfberry origin identification [91]	NIR-HSI (900-1700nm, 256 bands), 700 samples from 5 regions, 525:175 split	2D-COS to resolve overlapping peaks. CARS/IRIV/iVISSA for wavelength selection. GLCM for texture features	LDA, PLS-DA, SVM, CNN	CNN: 3 conv layers (16,32,64 filters), 3×3 kernel, ReLU, BatchNorm, MaxPool(2), GAP, FC(128), Dropout(0.2), Sigmoid	CNN+iVISSA: Acc=96.67%. CNN+texture: 97.71%. +9.71% vs PLS-DA, +7.42% vs SVM. <i>2D-COS resolves peaks, iVISSA selects wavelengths, Texture features improve accuracy</i>
Apple moldy core detection [92]	Vibration signals (100-1500Hz) + Vis/NIR (350-1150nm), 725 samples (180 normal + 545 diseased) split 3:1:1	CEEMDAN, Vibration + Vis/NIR fusion	PLS-DA, SVM, ELM, MobileNet, DMLPT	DMLPT: 3-layer Transformer Encoder for each input, AFF for fusion, MLP with residual connection. Multi-head attention	DMLPT+fusion: Acc=99.31% (normal/moderate/severe: 100%). +11.03% vs PLS-DA, +5.52% vs MobileNet. <i>Multi-modal fusion and multi-head attention excel at severity detection</i>
Apple SSC prediction [93]	Vis/NIR transmission spectra (589-1120nm, 1468 bands), 1450 spectra from 290 apples, 1020:430 split	GADF, SNV-UVE	PLS, MLR, VGG16, ResNet50, ShuffleNetv2, MobileViT	MobileViT: Conv + Transformer hybrid, CA mechanism for spatial features, multi-head attention. 3 MobileNet blocks, SiLU activation	GADF-MobileViT: $R^2=0.938$, RMSE=0.532. +6.6% vs PLS, +2.1% vs ResNet50. <i>GADF enables 2D transform, CA enhances features, Model focuses on key bands</i>
Tea quality classification [94]	NIR (1000-1800nm, 800 points), 1000 samples (50 grades, 21 brands), 750:250 split	SNV pre-processing. Transform 1D spectra to 2D pseudo-images (2×20×20)	PLS-DA, SVM, RF, TeaNet variants	TeaNet series: 3 conv/ residual/ inverted blocks, BatchNorm, ReLU, GAP, FC. Feature extraction + classification	TeaResNet+SNV: Acc=100%, TeaMobileNet: 99.6%. +31.2% vs SVM, +2.4% vs RF. <i>2D transform enables CNN, SNV increases variance, Models excel at multi-category</i>
Pb detection in oilseed rape [95]	FHSI 480-980nm, 2400 samples (1200 per environment), 3:1 split	SNV pre-processing, T-SCAE for cross-environment transfer	SVR with SPA/ CARS/ IRIV/ VISSA, SCAE, T-SCAE	Pre-trained SCAEs + extended layers (822-423-301-155)	T-SCAE+SNV: $R^2=0.939$, RMSE=0.020. +6.51% vs traditional. <i>Transfer learning enables cross-environment prediction</i>

spectral data to $[-1,1]$, then transforming them into polar coordinates through angular cosine encoding, and finally generating a GADF matrix by calculating the sine differences of angular values between each pair of spectral points, resulting in a matrix that preserves spectral relationships while optimizing for pattern recognition and machine learning applications, as in Figure 10. Different from 2D-COS,

which focuses on correlation analysis, and GADF, which uses polar coordinate transformation, TeaNet's approach [94] simply transforms one-dimensional NIR spectra with 800 spectral points into two-dimensional pseudo-images of size $2 \times 20 \times 20$ through direct matrix reshaping, requiring no complex mathematical calculations while still preserving all spectral information and enabling spatial rela-

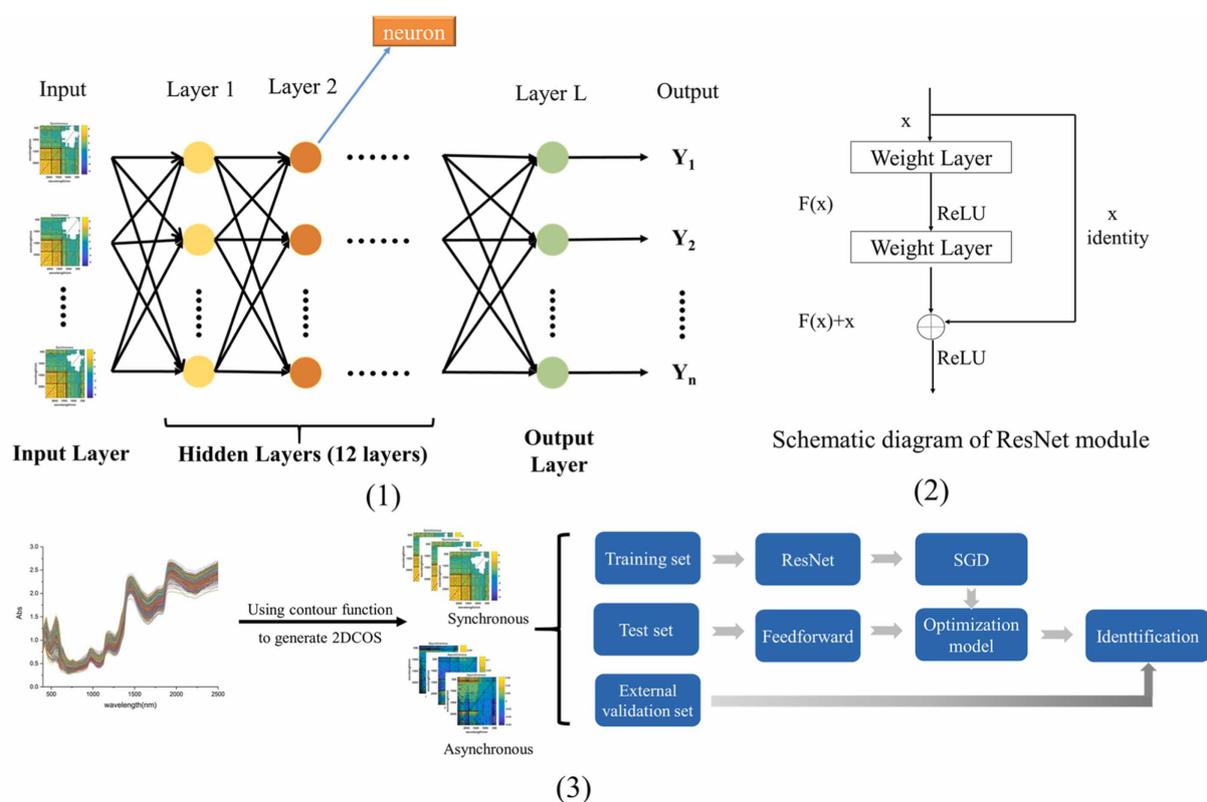


Figure 7: 2D-COS combined with ResNet process for identifying adulterated beef and mutton [85]

relationship analysis through CNN's convolutional operations across spectral bands.

DL demonstrates superior performance across these diverse data types. For 1D spectral data instance, in watermelon soluble solids content analysis [88], a 1D-CNN with five convolutional layers effectively processed 1D Vis-NIR spectra (317 to 1117 nm), achieving superior results compared to traditional methods. Similarly, for corn variety recognition [87], a CNN based on the LeNet-5 architecture successfully handled 1D NIR spectra ($11542\text{--}3940\text{ cm}^{-1}$), achieving 99.2% accuracy.

With 2D spectral image data, the study on subsurface bruise detection in plums [86] employed various CNN architectures, including HSCNN and ResNet, to process 2D hyperspectral images (430–1000 nm). These models effectively extracted spatial and spectral features, with HSCNN achieving the best F1 score of 90% using the full spectrum. For 2D dynamic spectral data, in the quality prediction of complex organic compounds [89], a novel DL model was designed to handle 2D NIR dynamic spectral matrices (time \times wavenumbers). This model incorporated multi-level dynamic feature extraction, including short-term (using dilated 2D CNN) and long-term (using GRU with temporal attention) feature extraction, effectively capturing both spatial and temporal characteristics of the spectral data. Additionally, for transformed 2D data, in identifying adulterated beef and mutton [85], researchers innovatively transformed 1D spectral data into 2D-COS before feeding it

into a ResNet model. This approach achieved 100% accuracy, demonstrating the potential of DL in processing transformed spectral data. Finally, the study on quantitative analysis of ADF and IVOMD in sugarcane [80] showcased the ability of a 1D-Inception-ResNet to handle data from multiple devices, achieving $R^2 > 0.96$ and $RMSEP < 2$ for both devices, demonstrating successful inter-device transfer learning.

The third highlight is the development of hybrid models and specialized architectures that have significantly advanced NIR spectral analysis, leading to improved performance across various tasks. These innovative approaches demonstrate the potential of tailored deep learning solutions in spectroscopy. These examples demonstrate how hybrid models and specialized architectures are pushing the boundaries of NIR spectral analysis, offering improved accuracy, robustness, and applicability across diverse tasks and data types.

Hybrid models combine DL techniques or integrate traditional methods with neural networks, leveraging the strengths of multiple approaches. For example, the T-SAE (Transfer Stacked Auto-Encoder) model used for lead content prediction in oilseed rape [81] combines transfer learning with auto-encoder architectures. This hybrid approach achieved impressive results with R^2 values of 0.9215 and 0.9349 for leaves and roots respectively, outperforming traditional methods. Another notable example is the multi-level dynamic feature extraction model for quality predic-

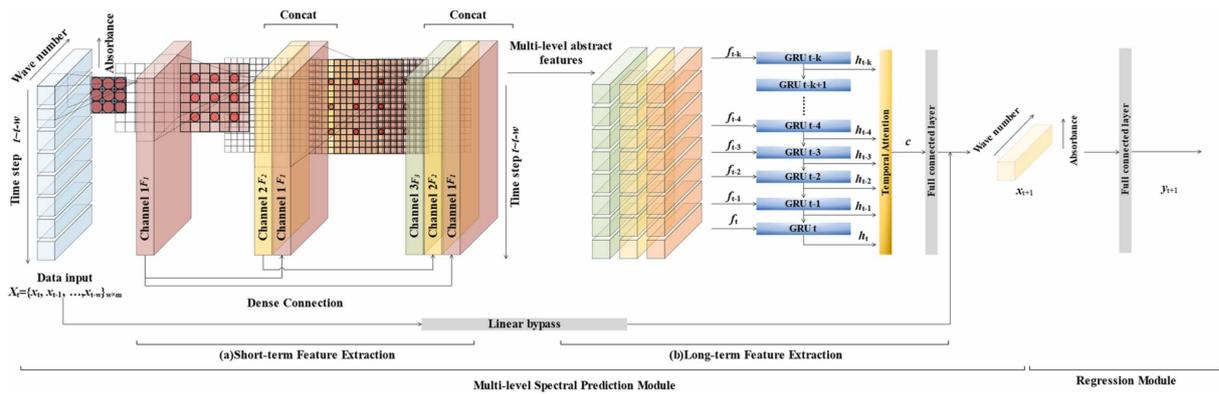


Figure 8: Framework of multi-level dynamic feature-based near-infrared quality prediction [89]

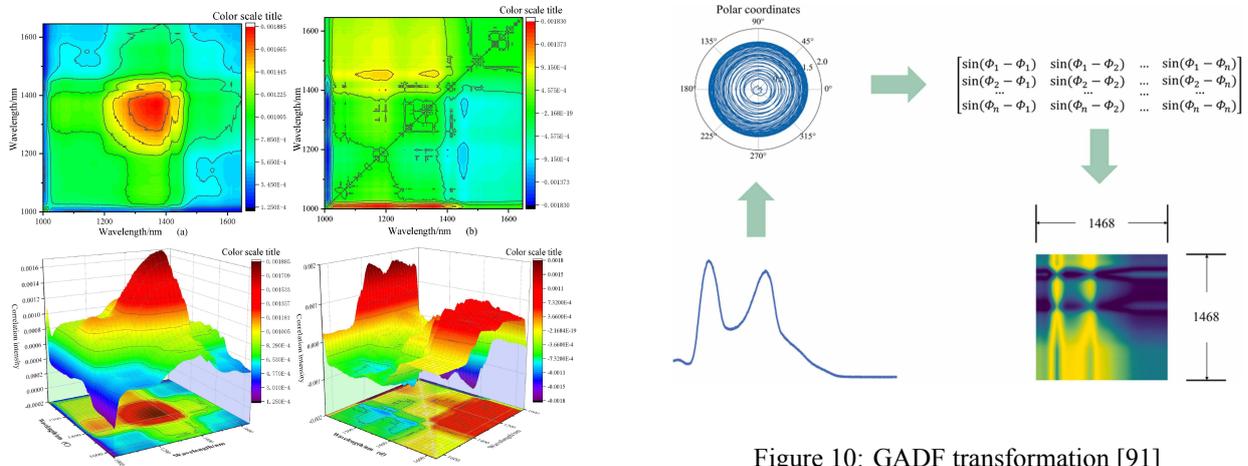


Figure 9: 2D-COS of wolfberries [91]

tion of complex organic compounds [89]. This hybrid approach combined dilated 2D CNNs for short-term feature extraction and GRU with temporal attention for long-term feature extraction, as mentioned above. The model’s sophisticated architecture included three dilated 2D CNN with varying dilated rates and kernel sizes, followed by a GRU layer with temporal attention. Besides, specialized architectures are designed to address specific challenges in NIR spectral analysis, often incorporating domain knowledge. The CNN-ATT model used for chickpea classification [84] incorporated an attention mechanism specifically designed to focus on relevant spectral regions. This model, consisting of an attention block followed by three 1D convolutional blocks and two dense layers, achieved 100% accuracy in classification.

Fourth highlight, transfer learning has emerged as a new powerful technique in NIR spectral analysis, allowing models to leverage knowledge from one task or dataset to improve performance on another. This approach is particularly valuable in scenarios with limited training data or when dealing with complex spectral relationships.

Several studies in the provided table demonstrate the effectiveness of transfer learning in various NIR applications. Inter-device transfer learning was successfully implemented in [80], with a 1D-Inception-ResNet model achieving consistently high performance across different spectrometers ($R^2 > 0.96$, $RMSEP < 2.75$). This architecture combines the Inception module for extracting diverse features at multiple scales from spectral data with Residual connections for efficient deep network training, successfully enabling model transfer across three different NIR spectrometers. Other studies have demonstrated effective applications of transfer learning in predicting lead content in oilseed rape. Using a cross-sample transfer approach, research by [81] employed the T-SAE model to predict lead content across different plant samples, achieving strong results for both leaves ($R^2 = 0.9215$, $RMSEP = 0.0302$ mg/kg) and roots ($R^2 = 0.9349$, $RMSEP = 0.0278$ mg/kg). In terms of cross-environment transfer, [95] developed the Transfer Stacked Convolutional Auto-Encoder (T-SCAE) architecture, as in Figure 11 to create a model that could work across different growing conditions. By combining pre-trained SCAE models from both silicon and silicon-free environments, this approach effectively extracted deep features for predicting lead concentrations in

Figure 10: GADF transformation [91]

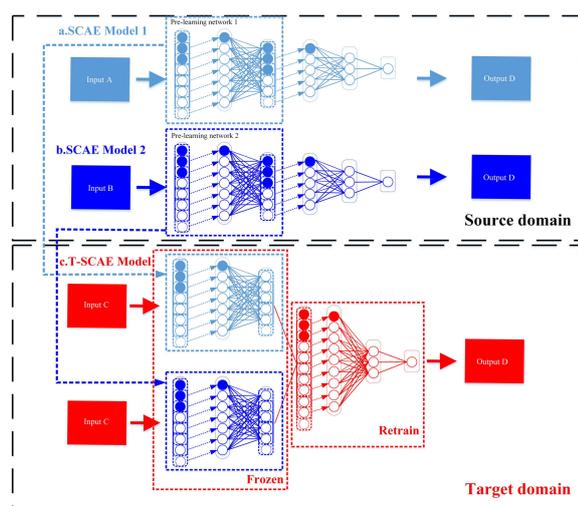


Figure 11: T-SCAE transfer network [95]

oilseed rape leaves, achieving excellent performance in the target domain with an R^2 of 0.9385, RMSEP of 0.02017 mg/kg, and RPD of 3.291. Additionally, spectral range transfer was also illustrated in [86], with the HSCNN model maintaining high performance (F1 score 89%) when reducing from full spectrum to just 3 wavelengths. The model was trained using the AdamW optimizer with a decaying learning rate strategy, and incorporated data augmentation techniques including intensity changes and spatial transformations to prevent overfitting, underscoring the versatility of transfer learning across different analytical dimensions. These examples highlight the versatility and power of transfer learning in NIR spectroscopy, significantly enhancing the adaptability and generalization capabilities of deep learning models in spectral analysis.

Discussion: Development insights

DL architectures demonstrate superior performance over traditional ML methods in NIR spectral analysis through advanced features like automatic hierarchical extraction, attention mechanisms, temporal learning, and dense feature fusion, significantly outperforming conventional approaches with accuracies of 97–100% in various analytical tasks (classification, regression, anomaly detection) across diverse spectral data formats (1D spectra, 2D correlation spectra, dynamic data). The development of hybrid models, specialized architectures, and effective transfer learning capabilities further enhances its robustness in handling complex spectral relationships and limited data scenarios, with proven success in cross-device ($R^2 > 0.96$) and cross-domain applications, marking a significant advancement in spectroscopic analysis. In summary, DL has significantly advanced NIR spectroscopy analysis through superior performance, versatility in data handling, innovative architectures, and effective transfer learning. However, several critical gaps remain that require further research.

Firstly, processing diverse data from multiple sources and devices continues to be a challenge, necessitating the devel-

opment of more robust methods to ensure consistency and accuracy.

Secondly, the interpretability of DL models in NIR spectral analysis needs improvement, particularly in developing specialized interpretation methods that incorporate expert knowledge in chemistry and spectroscopy.

Lastly, efficient learning from limited data, especially in applications such as food quality assessment and hazardous chemical detection, remains a significant challenge to be addressed.

With ongoing data collection efforts through collaborations with the People's Committee and other governmental agencies, coupled with the accumulation of diverse spectral datasets from multiple devices and regions, as mentioned in Section 4.2, ML and DL techniques will be comprehensively evaluated on larger-scale, heterogeneous data. This expanded evaluation scope will help address several key challenges: processing diverse data from multiple sources and devices to ensure consistency and accuracy, improving DL model interpretability while incorporating domain expertise in chemistry and spectroscopy, and developing efficient learning strategies for limited but critical data scenarios like food quality assessment and hazardous chemical detection.

5 Conclusion

This comprehensive review has analyzed recent advancements in the application of machine learning to near-infrared spectroscopy for food quality assessment. We have identified significant trends across various aspects of this field. In data collection, the trend towards using handheld or ultracompact NIR devices for direct on-site measurements, combined with multiple devices to collect broad NIR spectra, has significantly expanded the applicability of this technology. Regarding pre-processing and wavelength selection, automated and optimized processing techniques, along with the trend of combining multiple methods, have substantially improved model performance. In the field of deep learning, specialized architectures, and hybrid models have been developed, often outperforming traditional ML methods in many NIR spectral tasks. Additionally, transfer learning techniques have shown remarkable potential in addressing challenges related to interdevice variability, cross-sample analysis, and adaptation to new tasks or spectral ranges.

However, significant challenges remain to be addressed. Most notably, there is a lack of comprehensive datasets for food safety applications, a need to improve the interpretability of complex models, and the necessity to develop efficient learning methods from limited data. Additionally, processing diverse data from multiple sources and devices remains a major challenge to be resolved.

Based on these findings, we propose three important research directions for the future: Based on these research directions, we propose three important directions for future

research:

1. Develop larger and more diverse public datasets, with a particular focus on food safety parameters. This can be achieved through collaborations between food safety regulatory authorities and support from local government agencies to build NIR spectral datasets, following Ministry of Health standards. This includes collecting comprehensive spectral data accompanied by reference concentrations validated through independent laboratory testing.
2. Enhance machine learning models' processing capabilities and interpretability for heterogeneous data sources. This involves developing multi-level fusion techniques for integrating data from different devices, creating visualization methods for model interpretability, and establishing comprehensive cross-validation strategies using stratified sampling and bootstrapping techniques to ensure model reliability across diverse operating conditions. These approaches require systematic testing across devices and environments to establish standardized protocols for real-world applications.
3. Develop intelligent automation frameworks that integrate preprocessing selection, feature engineering, and model optimization. These frameworks should adapt to different device types and measurement conditions while maintaining the interpretability of results. The systems should include standardized evaluation metrics and clear protocols to enable seamless integration across NIR platforms in real-world applications.

These proposals aim to establish standardized approaches for NIR spectroscopy with machine learning, improving both the efficiency and reliability of food quality assessment processes. The implementation of these directions will help bridge the gap between theoretical advances and practical applications while addressing current challenges in real-world deployment.

Acknowledgement

This study is funded and implemented for the project with number "24/HĐ-SKHCHN, 2023". This work is supported by the People's Committee, Da Nang, and the University of Science and Technology, University of Danang.

References

- [1] H. Pu, J. Yu, D.-W. Sun, et al. "Feature construction methods for processing and analysing spectral images and their applications in food quality inspection". In: *Trends in Food Science & Technology* 138 (2023), pp. 726–737. DOI: 10.1016/j.tifs.2023.06.036.
- [2] S. M. Pires, H. G. Redondo, J. Pessoa, et al. "Risk ranking of foodborne diseases in Denmark: Reflections on a national burden of disease study". In: *Food Control* 158 (2024), p. 110199. DOI: 10.1016/j.foodcont.2023.110199.
- [3] W. Tian, Y. Li, C. Guzman, et al. "Quantification of food bioactives by NIR spectroscopy: Current insights, long-lasting challenges, and future trends". In: *Journal of Food Composition and Analysis* 124 (2023), p. 105708. DOI: 10.1016/j.jfca.2023.105708.
- [4] Y. Ozaki et al. *Near-Infrared Spectroscopy: Theory, Spectral Analysis, Instrumentation, and Applications*. Singapore: Springer Singapore, 2021.
- [5] A. Hassoun, S. Jagtap, G. Garcia-Garcia, et al. "Food quality 4.0: From traditional approaches to digitalized automated analysis". In: *Journal of Food Engineering* 337 (2023), p. 111216. DOI: 10.1016/j.jfoodeng.2022.111216.
- [6] I. Latreche, S. Slatnia, O. Kazar, et al. "A Review on Deep Learning Techniques for EEG-Based Driver Drowsiness Detection Systems". In: *Informatica* 48.3 (2024). DOI: 10.31449/inf.v48i3.5056.
- [7] H. A. Mohammed and I. M. Husien. "A Deep Transfer Learning Framework for Robust IoT Attack Detection". In: *Informatica* 48.12 (2024). DOI: 10.31449/inf.v48i12.5955.
- [8] J. Ravničan et al. "A Prestudy of Machine Learning in Industrial Quality Control Pipelines". In: *Informatica* 46.2 (2022). DOI: 10.31449/inf.v46i2.3938.
- [9] P. Mishra, D. Passos, F. Marini, et al. "Deep learning for near-infrared spectral data modelling: Hypes and benefits". In: *TrAC Trends in Analytical Chemistry* 157 (2022), p. 116804. DOI: 10.1016/j.trac.2022.116804.
- [10] H. Nobari Moghaddam, Z. Tamiji, M. Akbari Lakeh, et al. "Multivariate analysis of food fraud: A review of NIR based instruments in tandem with chemometrics". In: *Journal of Food Composition and Analysis* 107 (2022), p. 104343. DOI: 10.1016/j.jfca.2021.104343.
- [11] S. Othman, N. R. Mavani, M. A. Hussain, et al. "Artificial intelligence-based techniques for adulteration and defect detections in food and agricultural industry: A review". In: *Journal of Agriculture and Food Research* 12 (2023), p. 100590. DOI: 10.1016/j.jafr.2023.100590.
- [12] S. A. D. M. Zahir, A. F. Omar, M. F. Jamlos, et al. "A review of visible and near-infrared (Vis-NIR) spectroscopy application in plant stress detection". In: *Sensors and Actuators A: Physical* 338 (2022), p. 113468. DOI: 10.1016/j.sna.2022.113468.

- [13] S. A. D. M. Zahir, M. F. Jamlos, A. F. Omar, et al. “Review – Plant nutritional status analysis employing the visible and near-infrared spectroscopy spectral sensor”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 304 (2024), p. 123273. DOI: 10.1016/j.saa.2023.123273.
- [14] H. Ni, W. Fu, J. Wei, et al. “Non-destructive detection of polysaccharides and moisture in *Ganoderma lucidum* using near-infrared spectroscopy and machine learning algorithm”. In: *LWT* 184 (2023), p. 115001. DOI: 10.1016/j.lwt.2023.115001.
- [15] M. I. S. Mohd Hilmi Tan, M. F. Jamlos, A. F. Omar, et al. “*Ganoderma boninense* classification based on near-infrared spectral data using machine learning techniques”. In: *Chemometrics and Intelligent Laboratory Systems* 232 (2023), p. 104718. DOI: 10.1016/j.chemolab.2022.104718.
- [16] D. Wang et al. “Determination of polysaccharide content in shiitake mushroom beverage by NIR spectroscopy combined with machine learning: A comparative analysis”. In: *Journal of Food Composition and Analysis* 122 (2023), p. 105460. DOI: 10.1016/j.jfca.2023.105460.
- [17] Y.-Q. Zhong, J.-Q. Li, X.-L. Li, et al. “Near infrared spectroscopy for simultaneous quantification of five chemical components in *Arnebiae Radix* (AR) with partial least squares and support vector machine algorithms”. In: *Vibrational Spectroscopy* 127 (2023), p. 103556. DOI: 10.1016/j.vibspec.2023.103556.
- [18] Z. Guo, Y. Zhang, J. Wang, et al. “Detection model transfer of apple soluble solids content based on NIR spectroscopy and deep learning”. In: *Computers and Electronics in Agriculture* 212 (2023), p. 108127. DOI: 10.1016/j.compag.2023.108127.
- [19] C. Saenphon, S. Ditcharoen, C. Malai, et al. “Total soluble solids, dry matter content prediction and maturity stage classification of durian fruit using long-wavelength NIR reflectance”. In: *Journal of Food Composition and Analysis* 124 (2023), p. 105667. DOI: 10.1016/j.jfca.2023.105667.
- [20] S. Nawoya, F. Ssemakula, R. Akol, et al. “Computer vision and deep learning in insects for food and feed production: A review”. In: *Computers and Electronics in Agriculture* 216 (2024), p. 108503. DOI: 10.1016/j.compag.2023.108503.
- [21] S. Zhang, S. Liu, L. Shen, et al. “Application of near-infrared spectroscopy for the nondestructive analysis of wheat flour: A review”. In: *Current Research in Food Science* 5 (2022), pp. 1305–1312. DOI: 10.1016/j.crf.2022.08.006.
- [22] M. M. Nagy, S. Wang, and M. A. Farag. “Quality analysis and authentication of nutraceuticals using near IR (NIR) spectroscopy: A comprehensive review of novel trends and applications”. In: *Trends in Food Science & Technology* 123 (2022), pp. 290–309. DOI: 10.1016/j.tifs.2022.03.005.
- [23] L. Shuai, Z. Li, Z. Chen, et al. “A research review on deep learning combined with hyperspectral Imaging in multiscale agricultural sensing”. In: *Computers and Electronics in Agriculture* 217 (2024), p. 108577. DOI: 10.1016/j.compag.2023.108577.
- [24] Y. Liu, H. Pu, and D.-W. Sun. “Efficient extraction of deep image features using convolutional neural network (CNN) for applications in detecting and analysing complex food matrices”. In: *Trends in Food Science & Technology* 113 (2021), pp. 193–204. DOI: 10.1016/j.tifs.2021.04.042.
- [25] W. Zhang, L. C. Kasun, Q. J. Wang, et al. “A Review of Machine Learning for Near-Infrared Spectroscopy”. In: *Sensors* 22 (2022), p. 9764. DOI: 10.3390/s22249764.
- [26] P. Mishra, R. Nikzad-Langerodi, F. Marini, et al. “Are standard sample measurements still needed to transfer multivariate calibration models between near-infrared spectrometers? The answer is not always”. In: *TrAC Trends in Analytical Chemistry* 143 (2021), p. 116331. DOI: 10.1016/j.trac.2021.116331.
- [27] H. Ji, D. Pu, W. Yan, et al. “Recent advances and application of machine learning in food flavor prediction and regulation”. In: *Trends in Food Science & Technology* 138 (2023), pp. 738–751. DOI: 10.1016/j.tifs.2023.07.012.
- [28] Y. Zhang and Y. Wang. “Machine learning applications for multi-source data of edible crops: A review of current trends and future prospects”. In: *Food Chemistry: X* 19 (2023), p. 100860. DOI: 10.1016/j.fochx.2023.100860.
- [29] C. A. Nunes, M. N. Ribeiro, T. C. de Carvalho, et al. “Artificial intelligence in sensory and consumer studies of food products”. In: *Current Opinion in Food Science* 50 (2023), p. 101002. DOI: 10.1016/j.cofs.2023.101002.
- [30] B. Debus et al. “Deep learning in analytical chemistry”. In: *TrAC Trends in Analytical Chemistry* 145 (2021), p. 116459. DOI: 10.1016/j.trac.2021.116459.
- [31] M. Zareef, M. Arslan, M. M. Hassan, et al. “Recent advances in assessing qualitative and quantitative aspects of cereals using nondestructive techniques: A review”. In: *Trends in Food Science & Technology* 116 (2021), pp. 815–828. DOI: 10.1016/j.tifs.2021.08.012.

- [32] M. Hernández-Jiménez, I. Revilla, A. M. Vivar-Quintana, et al. “Performance of benchtop and portable spectroscopy equipment for discriminating Iberian ham according to breed”. In: *Current Research in Food Science* 8 (2024), p. 100675. DOI: 10.1016/j.crfs.2024.100675.
- [33] L. Yuan, X. Meng, K. Xin, et al. “A comparative study on classification of edible vegetable oils by infrared, near infrared and fluorescence spectroscopy combined with chemometrics”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 288 (2023), p. 122120. DOI: 10.1016/j.saa.2022.122120.
- [34] Widyaningrum, Y. A. Purwanto, S. Widodo, et al. “Rapid assessment of vanilla (*Vanilla planifolia*) quality parameters using portable near-infrared spectroscopy combined with random forest”. In: *Journal of Food Composition and Analysis* 133 (2024), p. 106346. DOI: 10.1016/j.jfca.2024.106346.
- [35] R. Chen, S. Li, H. Cao, et al. “Rapid quality evaluation and geographical origin recognition of ginger powder by portable NIRS in tandem with chemometrics”. In: *Food Chemistry* 438 (2024), p. 137931. DOI: 10.1016/j.foodchem.2023.137931.
- [36] J. P. Cruz-Tirado, M. S. S. Vieira, O. O. V. Correa, et al. “Detection of adulteration of Alpaca (*Vicugna pacos*) meat using a portable NIR spectrometer and NIR-hyperspectral imaging”. In: *Journal of Food Composition and Analysis* 126 (2024), p. 105901. DOI: 10.1016/j.jfca.2023.105901.
- [37] R. Zhu et al. “High-accuracy classification and origin traceability of peanut kernels based on near-infrared (NIR) spectroscopy using Adaboost - Maximum uncertainty linear discriminant analysis”. In: *Current Research in Food Science* 8 (2024), p. 100766. DOI: 10.1016/j.crfs.2024.100766.
- [38] J. Liang et al. “Integrating portable NIR spectrometry with deep learning for accurate Estimation of crude protein in corn feed”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 314 (2024), p. 124203. DOI: 10.1016/j.saa.2024.124203.
- [39] K. Yao, J. Sun, B. Zhang, et al. “On-line monitoring of egg freshness using a portable NIR spectrometer combined with deep learning algorithm”. In: *Infrared Physics & Technology* 138 (2024), p. 105207. DOI: 10.1016/j.infrared.2024.105207.
- [40] S. Ghidini, M. O. Varrà, D. Bersellini, et al. “Real-time and non-destructive control of the freshness and viability of live mussels through portable near-infrared spectroscopy”. In: *Food Control* 160 (2024), p. 110353. DOI: 10.1016/j.foodcont.2024.110353.
- [41] Z. Wu, C. Li, H. Liu, et al. “Quantification of caffeine and catechins and evaluation of bitterness and astringency of Pu-erh ripen tea based on portable near-infrared spectroscopy”. In: *Journal of Food Composition and Analysis* 125 (2024), p. 105793. DOI: 10.1016/j.jfca.2023.105793.
- [42] K. B. Beć, J. Grabska, and C. W. Huck. “Miniaturized NIR Spectroscopy in Food Analysis and Quality Control: Promises, Challenges, and Perspectives”. In: *Foods* 11 (2022), p. 1465. DOI: 10.3390/foods11101465.
- [43] J. M. Netto et al. “Authenticity of almond flour using handheld near infrared instruments and one class classifiers”. In: *Journal of Food Composition and Analysis* 115 (2023), p. 104981. DOI: 10.1016/j.jfca.2022.104981.
- [44] X. Chu et al. *Chemometric Methods in Analytical Spectroscopy Technology*. Singapore: Springer Nature Singapore, 2022.
- [45] J. A. Diaz-Olivares, A. Van Nuenen, M. J. Gote, et al. “Near-infrared spectra dataset of milk composition in transmittance mode”. In: *Data in Brief* 51 (2023), p. 109767. DOI: 10.1016/j.dib.2023.109767.
- [46] D. Li et al. “Research on Data Fusion and Sharing Based on Power Big Data”. In: *2023 9th Annual International Conference on Network and Information Systems for Computers (ICNISC)*. Wuhan, China, 2023, pp. 287–290. DOI: 10.1109/ICNISC60562.2023.00070.
- [47] L. Zhang et al. “Multi-source heterogeneous data fusion”. In: *2018 International Conference on Artificial Intelligence and Big Data (ICAIBD)*. Chengdu, China, 2018, pp. 47–51. DOI: 10.1109/ICAIBD.2018.8396165.
- [48] Tianzhe Jiao et al. “A Comprehensive Survey on Deep Learning Multi-Modal Fusion: Methods, Technologies and Applications”. In: *Computers, Materials and Continua* 80.1 (2024), pp. 1–35. DOI: 10.32604/cmc.2024.053204.
- [49] G. Van Kollenburg, Y. Weesepeol, H. Parastar, et al. “Dataset of the application of handheld NIR and machine learning for chicken fillet authenticity study”. In: *Data in Brief* 29 (2020), p. 105357. DOI: 10.1016/j.dib.2020.105357.
- [50] I. Malounas, W. Vierbergen, S. Kutluk, et al. “SpectroFood dataset: A comprehensive fruit and vegetable hyperspectral meta-dataset for dry matter estimation”. In: *Data in Brief* 52 (2024), p. 110040. DOI: 10.1016/j.dib.2024.110040.
- [51] G. Bonifazi et al. “A dataset of visible – Short wave InfraRed reflectance spectra collected on pre-cooked pasta products”. In: *Data in Brief* 36 (2021), p. 106989. DOI: 10.1016/j.dib.2021.106989.

- [52] M. Ryckewaert, D. Héran, C. Feilhes, et al. “Dataset containing spectral data from hyperspectral imaging and sugar content measurements of grapes berries in various maturity stage”. In: *Data in Brief* 46 (2023), p. 108822. DOI: 10.1016/j.dib.2022.108822.
- [53] A. A. Munawar, Kusumiyati, and D. Wahyuni. “Near infrared spectroscopic data for rapid and simultaneous prediction of quality attributes in intact mango fruits”. In: *Data in Brief* 27 (2019), p. 104789. DOI: 10.1016/j.dib.2019.104789.
- [54] R. Hayati, A. A. Munawar, and F. Fachruddin. “Enhanced near infrared spectral data to improve prediction accuracy in determining quality parameters of intact mango”. In: *Data in Brief* 30 (2020), p. 105571. DOI: 10.1016/j.dib.2020.105571.
- [55] Agussabti et al. “Data analysis on near infrared spectroscopy as a part of technology adoption for cocoa farmer in Aceh Province, Indonesia”. In: *Data in Brief* 29 (2020), p. 105251. DOI: 10.1016/j.dib.2020.105251.
- [56] A. Zgouz, D. Héran, B. Barthès, et al. “Dataset of visible-near infrared handheld and microspectrometers – comparison of the prediction accuracy of sugarcane properties”. In: *Data in Brief* 31 (2020), p. 106013. DOI: 10.1016/j.dib.2020.106013.
- [57] K. Kusumiyati et al. “Enhanced visible/near-infrared spectroscopic data for prediction of quality attributes in Cucurbitaceae commodities”. In: *Data in Brief* 39 (2021), p. 107458. DOI: 10.1016/j.dib.2021.107458.
- [58] I. Malounas et al. “Evaluation of a hyperspectral image pipeline toward building a generalisation capable crop dry matter content prediction model”. In: *Biosystems Engineering* 247 (2024), pp. 153–161. DOI: 10.1016/j.biosystemseng.2024.09.009.
- [59] H. Parastar et al. “Integration of handheld NIR and machine learning to “Measure & Monitor” chicken meat authenticity”. In: *Food Control* 112 (2020), p. 107149. DOI: 10.1016/j.foodcont.2020.107149.
- [60] J.-L. Z. Zaukuu, A. A. Nkansah, E. T. Mensah, et al. “Non-destructive authentication of melon seed (*Cucumeropsis mannii*) powder using a pocket-sized near-infrared (NIR) spectrophotometer with multiple spectral preprocessing”. In: *Journal of Food Composition and Analysis* 134 (2024), p. 106425. DOI: 10.1016/j.jfca.2024.106425.
- [61] S. Wang, M. Lin, Y. Meng, et al. “Self-expansion full information optimization strategy: Convenient and efficient method for near infrared spectrum auto-analysis”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 303 (2023), p. 123224. DOI: 10.1016/j.saa.2023.123224.
- [62] S. Wang, P. Zhang, J. Chang, et al. “A powerful tool for near-infrared spectroscopy: Synergy adaptive moving window algorithm based on the immune support vector machine”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 282 (2022), p. 121631. DOI: 10.1016/j.saa.2022.121631.
- [63] U. Blazhko et al. “Comparison of augmentation and pre-processing for deep learning and chemometric classification of infrared spectra”. In: *Chemometrics and Intelligent Laboratory Systems* 215 (2021), p. 104367. DOI: 10.1016/j.chemolab.2021.104367.
- [64] A. Sitorus and R. Lapcharoensuk. “Development of automatic tuning for combined preprocessing and hyperparameters of machine learning and its application to NIR spectral data of coconut milk adulteration”. In: *Food Chemistry* 457 (2024), p. 140108. DOI: 10.1016/j.foodchem.2024.140108.
- [65] N. D. Arianti, E. Saputra, and A. Sitorus. “An automatic generation of pre-processing strategy combined with machine learning multivariate analysis for NIR spectral data”. In: *Journal of Agriculture and Food Research* 13 (2023), p. 100625. DOI: 10.1016/j.jafr.2023.100625.
- [66] J.-M. Roger, A. Biancolillo, and F. Marini. “Sequential preprocessing through ORTHogonalization (SPORT) and its application to near infrared spectroscopy”. In: *Chemometrics and Intelligent Laboratory Systems* 199 (2020), p. 103975. DOI: 10.1016/j.chemolab.2020.103975.
- [67] Z. Wang, Q. Wu, and M. Kamruzzaman. “Portable NIR spectroscopy and PLS based variable selection for adulteration detection in quinoa flour”. In: *Food Control* 138 (2022), p. 108970. DOI: 10.1016/j.foodcont.2022.108970.
- [68] Y. Qin, K. Song, N. Zhang, et al. “Robust NIR quantitative model using MIC-SPA variable selection and GA-ELM”. In: *Infrared Physics & Technology* 128 (2023), p. 104534. DOI: 10.1016/j.infrared.2022.104534.
- [69] M. Marañón, J. Fernández-Novales, J. Tardaguila, et al. “NIR attribute selection for the development of vineyard water status predictive models”. In: *Biosystems Engineering* 229 (2023), pp. 167–178. DOI: 10.1016/j.biosystemseng.2023.04.001.
- [70] K. Yao, J. Sun, J. Cheng, et al. “Monitoring S-ovalbumin content in eggs during storage using portable NIR spectrometer and multivariate analysis”. In: *Infrared Physics & Technology* 131 (2023), p. 104685. DOI: 10.1016/j.infrared.2023.104685.

- [71] B. Mahanty. “Adaptive Bottom-Up Space Exploration in model population analysis: An agile variable selection algorithm for PLS models”. In: *Chemometrics and Intelligent Laboratory Systems* 203 (2020), p. 104057. DOI: 10 . 1016 / j . chemolab . 2020 . 104057.
- [72] J. Li, J. Deng, X. Bai, et al. “Quantitative analysis of aflatoxin B1 of peanut by optimized support vector machine models based on near-infrared spectral features”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 303 (2023), p. 123208. DOI: 10 . 1016/j . saa . 2023 . 123208.
- [73] X. Miao, Y. Miao, Y. Liu, et al. “Measurement of nitrogen content in rice plant using near infrared spectroscopy combined with different PLS algorithms”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 284 (2023), p. 121733. DOI: 10 . 1016/j . saa . 2022 . 121733.
- [74] C. Du, L. Sun, H. Bai, et al. “Quantitative detection of azodicarbonamide in wheat flour by near-infrared spectroscopy based on two-step feature selection”. In: *Chemometrics and Intelligent Laboratory Systems* 219 (2021), p. 104445. DOI: 10 . 1016 / j . chemolab . 2021 . 104445.
- [75] C. Du, L. Sun, H. Bai, et al. “Quantitative detection of talcum powder in wheat flour based on near-infrared spectroscopy and hybrid feature selection”. In: *Infrared Physics & Technology* 123 (2022), p. 104185. DOI: 10 . 1016 / j . infrared . 2022 . 104185.
- [76] R. Zheng, Y. Jia, C. Ullagaddi, et al. “Optimizing feature selection with gradient boosting machines in PLS regression for predicting moisture and protein in multi-country corn kernels via NIR spectroscopy”. In: *Food Chemistry* 456 (2024), p. 140062. DOI: 10 . 1016/j . foodchem . 2024 . 140062.
- [77] Z. Guo, L. Zhai, Y. Zou, et al. “Comparative study of Vis/NIR reflectance and transmittance method for on-line detection of strawberry SSC”. In: *Computers and Electronics in Agriculture* 218 (2024), p. 108744. DOI: 10 . 1016 / j . compag . 2024 . 108744.
- [78] Z. Ji, J. Zhu, J. Deng, et al. “Quantitative determination of zearalenone in wheat by the CSA-NIR technique combined with chemometrics algorithms”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* (2024), p. 124858. DOI: 10 . 1016/j . saa . 2024 . 124858.
- [79] J. Zhu et al. “Improve the accuracy of FT-NIR for determination of zearalenone content in wheat by using the characteristic wavelength optimization algorithm”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 313 (2024), p. 124169. DOI: 10 . 1016/j . saa . 2024 . 124169.
- [80] A. Tan et al. “1D-inception-resnet for NIR quantitative analysis and its transferability between different spectrometers”. In: *Infrared Physics & Technology* 129 (2023), p. 104559. DOI: 10 . 1016 / j . infrared . 2023 . 104559.
- [81] X. Zhou, C. Zhao, J. Sun, et al. “Detection of lead content in oilseed rape leaves and roots based on deep transfer learning and hyperspectral imaging technology”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 290 (2023), p. 122288. DOI: 10 . 1016/j . saa . 2022 . 122288.
- [82] X. Chen et al. “Rapid identification of total phenolic content levels in boletes by two-dimensional correlation spectroscopy combined with deep learning”. In: *Vibrational Spectroscopy* 121 (2022), p. 103404. DOI: 10 . 1016/j . vibspect . 2022 . 103404.
- [83] J. A. Martins, D. Rodrigues, A. M. Cavaco, et al. “Estimation of soluble solids content and fruit temperature in ”Rocha” pear using Vis-NIR spectroscopy and the SpectraNet–32 deep learning architecture”. In: *Postharvest Biology and Technology* 199 (2023), p. 112281. DOI: 10 . 1016/j . postharvbio . 2023 . 112281.
- [84] D. Saha, T. Senthilkumar, C. B. Singh, et al. “Rapid and non-destructive detection of hard to cook chickpeas using NIR hyperspectral imaging and machine learning”. In: *Food and Bioprocess Technology* 141 (2023), pp. 91–106. DOI: 10 . 1016/j . fbp . 2023 . 07 . 006.
- [85] L. Wang, J. Liang, F. Li, et al. “Deep learning based on the Vis-NIR two-dimensional spectroscopy for adulteration identification of beef and mutton”. In: *Journal of Food Composition and Analysis* 126 (2024), p. 105890. DOI: 10 . 1016/j . jfca . 2023 . 105890.
- [86] S. Castillo-Girones, R. Van Belleghem, N. Wouters, et al. “Detection of subsurface bruises in plums using spectral imaging and deep learning with wavelength selection”. In: *Postharvest Biology and Technology* 207 (2024), p. 112615. DOI: 10 . 1016 / j . postharvbio . 2023 . 112615.
- [87] J. Yang, X. Ma, H. Guan, et al. “A recognition method of corn varieties based on spectral technology and deep learning model”. In: *Infrared Physics & Technology* 128 (2023), p. 104533. DOI: 10 . 1016/j . infrared . 2022 . 104533.
- [88] G. Wang, X. Jiang, X. Li, et al. “Determination of watermelon soluble solids content based on visible/near infrared spectroscopy with convolutional neural network”. In: *Infrared Physics & Technology* 133 (2023), p. 104825. DOI: 10 . 1016 / j . infrared . 2023 . 104825.

- [89] Z. Chen, X. Luan, and F. Liu. “Deep learning near-infrared quality prediction based on multi-level dynamic feature”. In: *Vibrational Spectroscopy* 123 (2022), p. 103450. DOI: 10.1016/j.vibspec.2022.103450.
- [90] W. He, H. He, F. Wang, et al. “Non-destructive detection and recognition of pesticide residues on garlic chive (*Allium tuberosum*) leaves based on short wave infrared hyperspectral imaging and one-dimensional convolutional neural network”. In: *Food Measure* 15 (2021), pp. 4497–4507. DOI: 10.1007/s11694-021-01012-7.
- [91] Fujia Dong et al. “Identification of the proximate geographical origin of wolfberries by two-dimensional correlation spectroscopy combined with deep learning”. In: *Computers and Electronics in Agriculture* 198 (2022), p. 107027. ISSN: 0168-1699. DOI: 10.1016/j.compag.2022.107027.
- [92] Z. Liu et al. “Detection of apple moldy core disease by fusing vibration and Vis/NIR spectroscopy data with dual-input MLP-Transformer”. In: *Journal of Food Engineering* 382 (2024), p. 112219. DOI: 10.1016/j.jfoodeng.2024.112219.
- [93] Y. Li et al. “Combined gramian angular difference field image coding and improved mobile vision transformer for determination of apple soluble solids content by Vis-NIR spectroscopy”. In: *Journal of Food Composition and Analysis* 131 (2024), p. 106200. DOI: 10.1016/j.jfca.2024.106200.
- [94] J. Yang et al. “TeaNet: Deep learning on Near-Infrared Spectroscopy (NIR) data for the assurance of tea quality”. In: *Computers and Electronics in Agriculture* 190 (2021), p. 106431. DOI: 10.1016/j.compag.2021.106431.
- [95] Xin Zhou et al. “Determination of lead content in oilseed rape leaves in silicon-free and silicon environments based on deep transfer learning and fluorescence hyperspectral imaging”. In: *Spectrochimica Acta Part A: Molecular and Biomolecular Spectroscopy* 311 (2024), p. 123991. ISSN: 1386-1425. DOI: 10.1016/j.saa.2024.123991.

Early Diagnosis of Alzheimer's Disease with Transfer Learning Techniques Via ResNet50 and FSBi-LSTM

ZS. Khaleel^{1,2*}, Amir Lakizadeh¹

¹ Computer Engineering and Information Technology Department, University of Qom, Qom, Iran

² Department of mathematic, Open Educational College, Basra Study Center, Iraq

E-mail: zahsad@basraho.e.iq, lakizadeh@qom.ac.ir

*Corresponding author

Keywords: deep learning, alzheimer's disease, magnetic resonance imaging (MRI), ResNet, feature extraction

Received: October 16, 2024

Alzheimer's Disease (AD) is a neurological disorder marked by cognitive deterioration and neurological impairment that affects cognition, memory, and behavioral patterns. Alzheimer's is an incurable disease that predominantly impacts individuals over the age of 40. A patient's MRI (Magnetic Resonance Imaging) scan and cognitive assessment are manually analyzed to diagnose Alzheimer's disease. Recently, Artificial Intelligence (AI), particularly through Deep Learning, has pioneered innovative techniques for automated medical image identification. We devised a deep learning methodology for Alzheimer's disease identification utilizing Magnetic Resonance Imaging (MRI) data. The suggested method, termed Res+FSBiLSTM, employs ResNet50 as a pre-trained model for feature extraction from MRIs, thereafter identifying Alzheimer's disease through a Fully-Stack Bidirectional Long-Short Term Memory deep learning model. The experimental results demonstrate that the suggested method surpasses state-of-the-art techniques across all evaluation metrics, rendering it a viable tool for medical professionals in identifying Alzheimer's disease using brain radiological images. Ultimately, we achieved results with an accuracy of 99.6%, an F1-score of 97.7%, an area under the curve of 99%, a recall of 97.3%, and a precision of 99.6%.

Povzetek: Predstavljena je izboljšana metoda zgodnjega diagnosticiranja Alzheimerjeve bolezni z uporabo prenosa učenja prek ResNet50 in FSBi-LSTM. Sistem uporablja MRI slike za avtomatizirano identifikacijo bolezni in omogoča hitro prepoznavanje zgodnjih faz bolezni, kar omogoča učinkovitejšo klinične diagnoze in obravnavo pacientov.

1 Introduction

The predominant kind of dementia necessitating extensive medical intervention is Alzheimer's disease (AD). For effective patient therapy to begin, an early and accurate assessment of the prognosis for AD is necessary [1]. The research found that there are 10 million new instances of dementia reported per year [2]. According to the World Health Organization (WHO), AD has overtaken cancer as the leading to death, with the quantity of AD By 2050, there will be 152 million patients. AD is a chronic neurological brain illness that progressively damages brain tissue, leading to cognitive decline and memory loss, and ultimately hastening the loss of ability to carry out daily tasks [3]. A condition known as AD is brain-neurological degeneration [4]. It is classified as dementia, which is brain atrophy that impairs memory and results in loss of cognitive abilities related to behavior, social interaction, and reasoning. Protein fragments build up in the brain, which is the cause of it [5, 6, 7]. The human brain develops plaques and tangles around the neurons, causing aberrant hippocampal and lobe shrinkage as well as enlarged ventricles [8]. It is a deadly illness that is incurable [9, 10], causing the patient to suffer for the rest of their life and their family great emotional, physical, and financial hardship. There are no

known causes of AD, and there are no treatments or drugs that can effectively cure dementia. A pre-clinical stage of AD called mild cognitive impairment (MCI) is a transitional condition that occurs between normal aging and AD. Early detection of the risk and severity of AD is crucial [11,12]. Many researchers worldwide have developed a multitude of Machine Learning (ML) [13] and Deep Learning (DL) [14] algorithms throughout the years for the purpose of AD detection and classification. The DL algorithms have been used by numerous academics with impressive results still, there's space for development. In this set of DL models [15,16], and [17] that introduce a hybrid Convolutional Neural Network (CNN), a CNN model with slice selection, and a CNN model with histogram stretching. A CNN model with skull striping was presented by others [18] and in [19]. However, because CNN is a black-box, these deep models are predominantly biased towards categorization. Several researchers have created a variety of techniques and applications for automatic neuro-image segmentation in the literature on AD [20]. While there is little attention on CNN layers to visualize the classification process, these applications are useful tools for segmenting neuro-images. Each convolution layer's feature map shows the different filters that are applied to the image and gives an indication of the kinds of filters the model applies to the image to

extract features [21]. In sequence analysis, the recurrent neural network (RNN) is a potent model. Since RNN uses a "state" vector in its hidden units, all historical information about the sequence is inherently contained in it [22]. Long short-term memory (LSTM), an enhanced version of RNN, can handle gradient explosion or gradient disappearance issues more successfully than RNN by controlling information flow across numerous gates [23]. Moreover, contextual information can be present in both directions for bidirectional LSTM (Bi-LSTM) [24]. In fact, by stacking the LSTM to examine the spatial information of feature maps from CNN layers as well, Bi-LSTM may obtain more information without having to select the scanning direction. Therefore, we propose to use the fully stacked Bi-LSTM (FSBi-LSTM) instead of the traditional Bi-LSTM [25]. In this paper, we create an innovative deep learning network that employs fully stacked bidirectional LSTM (FSBi-LSTM) and CNN layers of ResNet50 [26] to diagnose AD using multimodal input. The goal of the proposed model is to produce a classification result that is accurate enough to identify AD at an earlier stage. The research study's primary contributions are:

- Combining traditional image processing techniques like thresholding and morphological operations with modern deep learning approaches such as U-Net provides a robust method for brain extraction.
- By incorporating FSBiLSTM with ResNet50, you can leverage the power of ResNet50's strong feature extraction capabilities while benefiting from FSBiLSTM's ability to generalize from a few examples.
- Improved generalization: ResNet50 is a powerful deep neural network architecture known for its ability to learn rich representations from images. By combining it with FSBiLSTM, which specializes in few-shot learning, you can potentially improve the generalization performance of ResNet50 on new, unseen classes with limited training data
- Using ResNet50 for Spatial Learning: ResNet50 is an incredibly potent convolutional neural network (CNN) that is highly skilled at extracting detailed spatial characteristics from individual MRI slices. These characteristics are essential for detecting alterations linked to Alzheimer's disease because they capture minute details and patterns.

The rest of the paper is arranged as follows: Section 2 provides a summary of the relevant studies for the suggested model. Next, in Section 3, we provide a thorough explanation of our methodology. The experiment's results are described in Section 4. The dissection is provided in Section 5. Finally, Section 6 contains a summary of our findings.

2 Related work

In recent years, various methods have been suggested to enhance the accuracy of image classification. We

analyze various machine learning (ML) classification frameworks employed in neuroimaging, along with techniques based on a convolutional neural network. Machine learning techniques have been increasingly utilized in recent years for the early identification of Alzheimer's disease, particularly in multi-class and binary classification tasks. Yiming Ding et al. [27] Suggested doing a retrospective analysis of 2109 18F-FDG PET imaging tests that were gathered prospectively from 1002 patients. The majority of patients had numerous scans, and the dates of the scans ranged from May 2005 to January 2017. The researchers created and verified a deep learning algorithm, especially a convolutional neural network using InceptionV3 architecture. The model underwent training using 90% of the dataset and was subsequently evaluated using the remaining 10%, in addition to an independent test set. The model's performance was compared to that of radiologic readers. The model's performance was assessed by sensitivity, specificity, receiver operating characteristic (ROC), saliency map, and t-distributed stochastic neighbor embedding. The study's authors are C. Suh and colleagues [28]. Developed a deep learning technique employing a dataset of T1-weighted brain MRI scans from consecutive patients diagnosed with Alzheimer's disease and moderate cognitive impairment. The researchers developed a two-step system employing a convolutional neural network for brain parcellation, subsequently applying three classification techniques, including XGBoost, for disease prediction. The categorization experiments were performed with a 5-fold cross-validation method. The diagnostic efficacy of the XGBoost algorithm was evaluated against logistic regression and a linear SVM. The areas under the curve were calculated to differentiate between Alzheimer's disease and moderate cognitive impairment, as well as between mild cognitive impairment and healthy controls. Hina Nawaz et al. [29] proposed a pre-trained AlexNet model for the extraction of deep features in the detection of Alzheimer's disease stages. Transfer learning uses the initial layers of the pre-trained AlexNet model to extract deep features from the CNN. SVM, k-nearest Neighbor (KNN), and Random Forest (RF) are employed to classify the extracted deep features through machine learning techniques. The study is authored by Hadeer A. Helaly and colleagues. [30]. This investigation focuses on the prompt identification and classification of Alzheimer's disease through the application of an advanced machine learning method referred to as CNNs. The E2AD2C framework is designed to identify and categorize various stages of Alzheimer's disease at an early phase. Two primary methodologies are utilized for medical picture categorization and the identification of Alzheimer's disease. The preliminary method utilizes fundamental CNN architectures to analyze 2D and 3D structural brain scans sourced from the ADNI dataset. This technique attains classification accuracies of 93.61% and 95.17% for multi-class Alzheimer's Disease stage classifications. The alternative method involves employing transfer learning with the VGG19 pre-trained model, which has been optimized to achieve 97% accuracy in diagnosing multi-class Alzheimer's disease

stages. Resampling methods, including oversampling and downsampling, are utilized to address the problem of class imbalance in a dataset. Data augmentation techniques are employed to increase dataset size and alleviate overfitting concerns. Fazal Ur Rehman Faisal [31] concentrated on devising a deep learning technique for the extraction of Alzheimer's disease (AD) biomarkers from structural magnetic resonance imaging (sMRI) data. CNNs were complexified and improved efficiency. The proposed approach was evaluated against leading methodologies for AD classification,

exhibiting enhanced performance for accuracy and area under the ROC curve. The convolution operation utilized in the proposed method was considered appropriate for Alzheimer's disease diagnosis, demonstrating its efficacy in appropriately classifying the brain picture. Y. F. Khan et al. [32] utilized a Stacked Deep Dense Neural Network (SDDNN) model, integrating Convolutional Neural Network (CNN) and Bidirectional Long Short-Term Memory (Bidirectional LSTM). The training utilized the DementiaBank clinical transcript dataset under two configurations: randomly initialized parameters and GloVe embedding. Hyperparameter optimization via GridSearch enhanced model performance, attaining 93.31% accuracy with GloVe embedding and fine-tuning.

Ahsan Bin Tufail et al. [33] utilized transfer and non-transfer learning-based CNN architectures in both 2D and 3D domains for binary and multiclass classification tasks. Custom 3D CNN architectures were created for binary classifications of AD/NC, AD/MCI, and NC/MCI, as well as for multiclass classification of AD/NC/MCI. Transfer learning utilizing the Xception architecture was employed for the classification of MCI and AD, as well as for the multiclass categorization of NC, MCI, and AD. A bespoke CNN architecture in the 2D domain was employed for the categorization of NC and MCI, as well as NC and AD classes. Evaluation utilized performance criteria such as CEN, RCI, GM, IBA, and MCC. Data augmentation methods, such as random zooming in and out, were employed to enhance dataset size for better generalization of deep learning algorithms. R. Tandon et al. [34] utilize deep learning methodologies for the segmentation and classification of Alzheimer's Disease through brain imaging data. The suggested approach amalgamates segmentation models with classification frameworks to enhance diagnostic precision. The system attains remarkable performance utilizing datasets such as MRI scans, with classification accuracy reported to exceed 90%.

employed and trained on sMRI brain pictures sourced from ADNI datasets to categorize images into Alzheimer's disease, mild cognitive impairment (MCI), and cognitively normal (CN) groups. Features from several layers were integrated to hierarchically convert MRI pictures into more concise high-level features, facilitating the classification process. The strategy sought to diminish the number of parameters to decrease computational

The study highlights the significance of precise segmentation in improving the diagnostic efficacy of deep learning models for Alzheimer's detection [35]. A proposed retrospective cohort study involving 532 participants, employing Positron Emission Tomography (PET) and Magnetic Resonance Imaging (MRI) images, alongside cognitive evaluations. The authors developed a novel computational phenotyping method that utilizes Partial Volume Correction (PVC) and subsets of neuropsychological assessments in an impartial manner. The pipeline employs a Regional Spread Function (RSF) technique for PVC and a t-distributed Stochastic Neighbor Embedding (t-SNE) manifold. The objective was to develop a new method for analyzing variations in cognitive scores and PET characteristics to identify multiple phenotypes of Alzheimer's disease (AD) using hyperparametric analysis.[36] Ullah et al. (2023) introduced modifications to pre-trained deep learning models, including ResNet and DenseNet, aimed at enhancing brain tumor classification. Their approach enhanced accuracy and robustness through transfer learning and domain-specific fine-tuning, effectively addressing challenges such as overfitting in limited medical datasets. The research indicated that these improved models outperformed conventional methods, highlighting their appropriateness for clinical use.[37] Chegireddy and Srinagesh (2023) proposed a new deep-learning approach for predicting pancreatic cancer by utilizing human MRI data, incorporating variants of Harris Hawks Optimization (HHO) with the VGG16 architecture model. Their approach employed HHO variants to optimize hyperparameters and enhance feature selection, improving the predictive performance of VGG16. The proposed method achieved high accuracy and robustness, outperforming traditional models and standard deep learning techniques. This study emphasizes the potential of integrating metaheuristic algorithms with deep learning frameworks to enhance diagnosis accuracy in medical imaging.

Table 1: Comparison of machine learning methods for alzheimer’s disease classification

No.	Method	Dataset Used	Accuracy (%)	F1 Score	Key Features
[27]	Deep Learning on 18F-FDG PET (Ding et al.)	Brain PET scans (N = 1,002)	92	-	Predicts Alzheimer’s diagnosis with PET imaging using deep neural networks.
[28]	3D T1-Weighted Images (Suh et al.)	ADNI dataset	89.5	0.88	Brain segmentation and classification using 3D CNNs.
[29]	Deep Feature Real-Time Detection (Hina et al.)	Proprietary	87	-	Stage detection with real-time deep feature-based processing.
[30]	Early Detection via CNN (Helaly et al.)	ADNI	93.2	0.91	Focus on early diagnosis using convolutional networks.
[31]	Whole Brain MRI (Faisal & Kwon)	OASIS and ADNI	94.5	0.92	Automated detection leveraging MRI-based deep learning.
[32]	Stacked Dense NN on Audio Data (Khan et al.)	Proprietary	85.6	0.84	Predicts dementia using audio transcript data with stacked neural networks.
[33]	2D/3D CNN with PET Neuroimaging (Ahsan et al.)	ADNI and OASIS	92.8	0.89	PET neuroimaging-based early-stage categorization in 2D and 3D domains.
[34]	Deep Learning for AD (Buvaneswari et al.)	ANDI	95	-	the study demonstrates that combining SegNet for feature extraction with ResNet-101 for classification can effectively identify Alzheimer’s disease with high accuracy and sensitivity.
[35]	Manifold Learning on PET (Campanioni et al.)	Proprietary	91.3	-	Epigenetic phenotyping using PET imaging and manifold learning.
[36]	DBNs with IoT Detection (Alqahtani et al.)	Proprietary	90.5	0.87	Combines deep belief networks with IoT for detection and classification.
[37]	Ensemble Learning with Synthetic Data (Mujahid et al.)	Proprietary	93.8	0.90	Ensemble approach using synthetic data augmentation for robust predictions.

3 Methodology

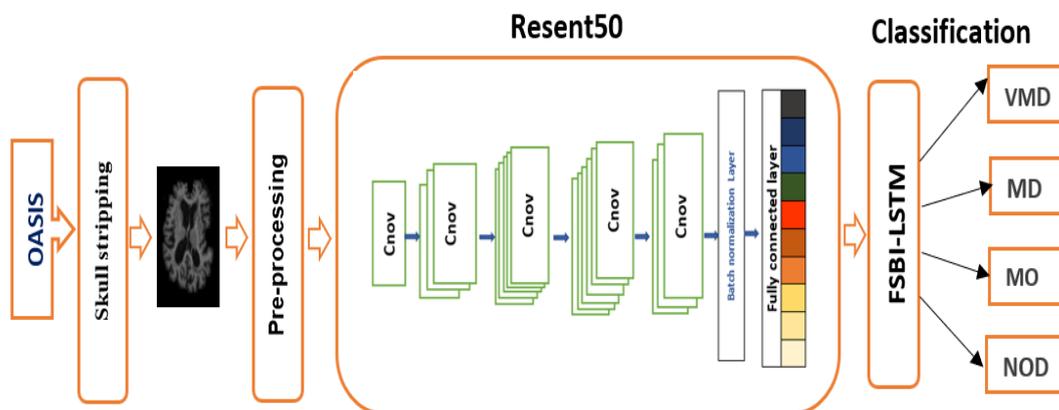


Figure 1: The overall framework of the proposed method.

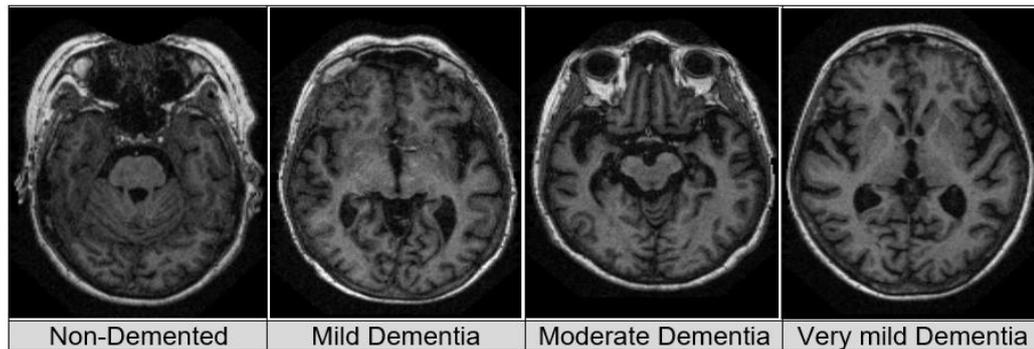


Figure 2: some samples from the OASIS dataset.

The proposed framework is illustrated in Figure. 1, we apply CNN of Resnet50 to extract the basic features datasets that we used from OASIS organization, FSBI-LSTM then is used instead of FC (fully connected layer) to high level semantic, contextual understanding, more.

3.1. The datasets

Numerous datasets for Alzheimer’s disease classification are accessible online. Numerous AD datasets are unsuitable for this research due to their CSV format. Organizations such as OASIS provide their data sets for educational and research purposes. However, the samples in data sets are presented in a three-dimensional image format. OASIS is a project

that grants researchers and the scientific leveraging of semantic information and capturing the contextual dependencies finally, for classifying the disease diagnosis the concatenated learned features are passed to SoftMax in the following description of the proposed techniques for community access to an extensive collection of neuroimaging data. The objective of OASIS is to advocate for open science and enhance the progression of research in neuroimaging and neuroscience. In this dataset, we utilized four classifications: Non-Demented, Mild Dementia, Moderate Dementia, and Very Mild Dementia, as illustrated in Figure 2.

Table 2: Clears the distribution of OASIS datasets.

Classes	NO. of images	Gender	Age range
Non-Demented (<i>NOD</i>)	100	F/M	55-85 years
Mild Dementia (<i>MD</i>)	5002	F/M	62-85 years
Moderate Dementia (<i>MOD</i>)	488	F/M	63-85 years
Very Mild Dementia (<i>VMD</i>)	102	F/M	65-88 years

3.2. Data processing

The initial preprocessing step in our system is skull-stripping, an essential procedure that eliminates non-brain tissues and structures from MRI images to separate the cerebral region. This method is extensively employed in neuroimaging research and clinical applications to improve the emphasis on brain tissues, hence enabling more precise analysis and segmentation. Several methods are available for skull-stripping, ranging from simple thresholding techniques to advanced deep learning models. In our approach, we combine the traditional image processing techniques with a deep learning-based U-Net model to achieve robust and accurate results. Thresholding is initially employed as a straightforward yet efficient technique that transforms a grayscale image into a binary image by juxtaposing pixel intensity values

against a specified threshold. Pixels with intensities beyond the threshold are designated as foreground (brain tissue), and those below are categorized as background (non-brain tissue). This approach is most efficacious when a distinct intensity differentiation exists between cerebral and non-cerebral areas. Morphological techniques, like erosion and dilation, are subsequently employed to enhance the binary image, so increasing the segmentation of the brain region by smoothing borders and removing minor aberrations or noise. The extracted brain region is further refined for segmentation through the successive application of erosion, dilation, and opening processes. A pre-trained or custom-trained U-Net model is subsequently utilized to boost segmentation, capitalizing on its capacity to understand complex patterns and features for improved precision in isolating the brain from adjacent structures.

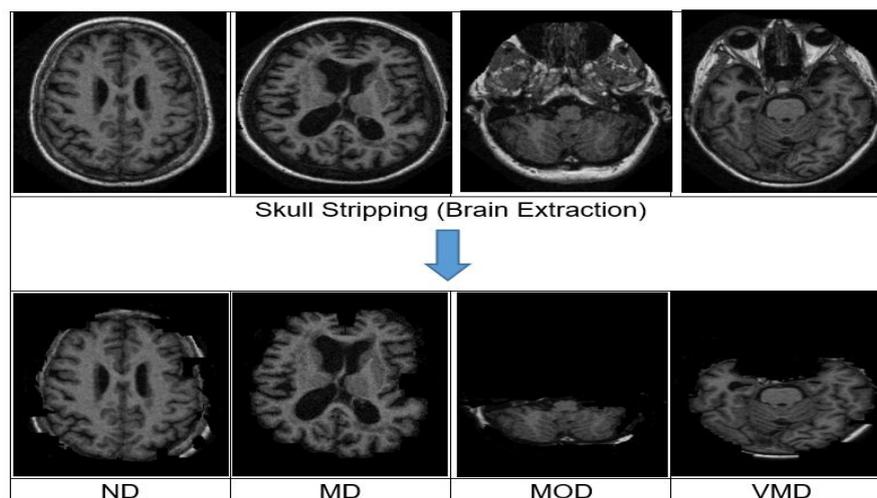


Figure 3: Extracting the brain section for images in the OASIS dataset

Our method for effective skull-stripping commences with noise reduction by Gaussian blur on the grayscale image, subsequently employing Otsu's thresholding technique to dynamically ascertain the appropriate threshold value and transform the blurred image into binary format. Morphological opening is subsequently applied to the binary picture utilizing a preset kernel, followed by dilation to enlarge the region of interest (ROI). Contours are detected in the dilated image utilizing the ``cv2.findContours()`` function, and the contour with the greatest area is designated as the principal ROI. A mask is generated to correspond with the image dimensions and delineate the chosen contour, therefore isolating the ROI from the original image. The retrieved ROI is further modified with a U-Net model, which improves segmentation and yields effective skull-stripping outcomes. Our method achieves precise and robust skull-stripping by integrating Gaussian blurring, Otsu's thresholding, and morphological operations with the segmentation capabilities of U-Net, thereby combining the simplicity and efficiency of traditional techniques with the sophistication of deep learning for optimal performance across various MRI datasets. Figure 3 illustrates the effectiveness of our skull-stripping approach, showing MRI images before and after processing.

Figure 3 clearly demonstrates that the brain extraction approach we employed is highly effective on the OASIS Foundation dataset since non-brain regions were entirely eliminated from the images, yielding great results in our data analysis. To mitigate the issue of overfitting, we utilized augmentation approaches. Overfitting, a prevalent issue in machine learning and statistical modeling, transpires when a model gets excessively tailored to the training data, leading to inadequate generalization of novel data. Augmented processing entails enhancing the training data by incorporating changes, such the addition of noise, the application of modifications, or the generation of synthetic samples. Augmented processing, through the diversification of the training set, exposes the model to a wider array of patterns and variances, hence mitigating the danger of overfitting.

This strategy enables the model to acquire a greater level of resilience and generalizable representations, improving its performance on unseen data, we applied various augmentation strategies tailored to our specific dataset, such as random rotation in the range $[10, -10]$, both vertical and horizontal shifting in the range $(-0.1, 0.1)$, flips vertical and horizontal randomly, and shear of the original images in the range $(-0.1, 0.1)$. The additional data substantially improved the training process, resulting in a more efficient and dependable model with higher generalization abilities, so increasing its capacity to manage unseen input and ultimately reduce the effects of overfitting. The tables below show the number of images after the augmentation of our OASIS datasets. A pre-processing operation is also applied to improve the ability of the proposed model to elevate the quality of the data and extract relevant features such as image rescaling, image normalization, and skull stripping as mentioned before.

Table3: Number of OASIS images after the augmentation.

Classes	NO. of images
Non-Demented	990
Mild Dementia	34216
Moderate Dementia	4598
Very Mild Dementia	1002

3.3 Feature learning (feature extraction)

In this research, we have resorted to building a model based on both CNN and RNN. The combination of CNNs and RNNs has been successfully applied to a variety of classification tasks, such as video classification, sentiment analysis, and medical image analysis. By leveraging the strengths of both spatial feature extraction and sequential modeling, the CNN-RNN approach can often outperform models that only use one type of neural network architecture. ResNet-50 has been effectively utilized across multiple domains, notably in medical imaging, namely for the interpretation of MRI (Magnetic Resonance Imaging) pictures. ResNet-50, an abbreviation

for Residual Network-50, is a convolutional neural network (CNN) architecture that has achieved considerable acclaim and efficacy in the domain of computer vision. It was presented by Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun in their foundational paper "Deep Residual Learning for Image Recognition" [26]. The creation of ResNet-50 was driven by the understanding that deeper convolutional neural networks (CNNs) often encounter the vanishing gradient problem, which hinders their training efficacy.

The vanishing gradient problem arises when gradients transmitted across the network diminish swiftly, obstructing preceding layers from updating their weights and acquiring significant representations. ResNet-50 tackles this difficulty by an innovative architectural design grounded in the principle of residual learning.

The fundamental principle behind residual learning is the incorporation of skip connections or shortcuts, which allow the network to learn residual functions instead of directly fitting the intended underlying mapping. The inclusion of skip connections allows for the direct propagation of gradients throughout the network, hence addressing the issue of vanishing gradients.

ResNet-50 consists of 50 layers, categorizing it as a network of significant depth. The architecture comprises a series of residual blocks, each containing several convolutional layers. The residual blocks are organized in a stratified configuration to form the complete architecture of ResNet-50. The network utilizes a combination of 1x1, 3x3, and occasionally 1x1 convolutional filters, along with batch normalization and rectified linear unit (ReLU) activations, to extract and transform visual input at different levels of abstraction. ResNet-50 is notable for its ability to achieve state-of-the-art accuracy on numerous challenging benchmark datasets, including ImageNet, which contains millions of annotated images. The intricate architecture of ResNet-50 allows it to proficiently capture and represent complicated patterns and hierarchical structures, yielding outstanding performance in applications such as image classification, object detection, and image segmentation. Furthermore, the architecture of ResNet-50 has influenced subsequent developments in CNN design.

3.4. Alzheimer's disease detection

In convolutional neural networks (CNNs), the standard approach employs fully connected (FC) layers for high-level analysis. Nonetheless, completely connected (FC) layers are ineffective in acquiring comprehensive spatial information from feature maps, as they connect all neurons indiscriminately, disregarding spatial correlations. In this paper, we have presented an enhanced version of the LSTM (Long Short-Term Memory) framework, termed FSBi-LSTM, to address this issue. LSTM, an acronym for Long Short-Term Memory, is a specialized type of recurrent neural network (RNN) that is proficient in handling sequences and retaining long-term dependencies. In traditional recurrent neural networks (RNNs), the output of a cell at a given time step is determined by the input at that time step and the output of

the prior cell. The output is denoted as " h_t " and is calculated using the specified formula

$$h_t = f(Ux_t + Wh_{t-1}) \quad (1)$$

In this context, x_t denotes the input at the present time step, U signifies the weight matrix linking the input layer to the hidden layer, and W represents the weight matrix connecting the output of the preceding cell to the current cell.

The function $f(\cdot)$ typically represents the hyperbolic tangent (tanh) function, which constricts input values to a range between -1 and 1. The purpose of employing this activation function is to introduce non-linearity into the network, allowing it to model complex relationships between inputs and outputs. Conventional RNNs have two major issues known as "gradient explosion" and "gradient disappearance." These problems arise due to the gradients either diminishing or escalating during the training phase. The gradient measures the slope of the loss function relative to the network's parameters and is utilized to adjust the network's weights during training. When the gradients decrease to an insignificant level, the network faces difficulties in its learning process and may finally converge to suboptimal solutions. Conversely, if the gradients become very large, the network's weights may experience substantial updates, leading to instability and protracted convergence. The Long Short-Term Memory (LSTM) architecture was developed to tackle these difficulties.

The LSTM model integrates gates as shown in figure 4 and a cell state to control the flow of information inside the network. The present cell state in an LSTM is denoted as " ct ". Its principal job is to retain and convey information over different temporal intervals.

The Long Short-Term Memory (LSTM) model employs three fundamental gates: the input gate, the forget gate, and the output gate. Each gate is depicted as a fully connected (FC) layer, consisting of a set of trainable weights. These gates control information transmission by selectively allowing or blocking particular components of the cell state and output.

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \quad (2)$$

The input gate controls the degree of integration of new input information into the cell state. The calculation of the input gate activation relies on the present input x_t and the preceding output h_{t-1} . The weights W_{xi} and W_{hi} , along with the bias b_i , are utilized to compute the activation of the input gate. The sigmoid function σ is employed in this computation. This activation is then applied to the new candidate values, which signify potential alterations to the cell state.

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \quad (3)$$

The forget gate determines the fraction of the previous cell state to retain and transmit to the current time step. The activation of the forget gate is determined by the current input x_t and the previous output h_{t-1} , utilizing weights

W_{xf} and W_{hf} . The activation is employed to modify the previous cell state, allowing the LSTM to discard or dismiss superfluous input.

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \quad (4)$$

The output gate controls the degree to which the present cell state is disclosed as the output h_t . The present input x_t and the preceding output h_{t-1} are taken into account, and their weights are employed to compute the activation of the output gate. The activation is multiplied by the modified cell state to yield the final output o_t .

By employing these gates, LSTM effectively alleviates the problems of gradient explosion and gradient vanishing typically faced by traditional RNNs. The gates facilitate the network's regulation of information flow, selectively retain or eliminate data from the cell state, and control the output of information. LSTMs can capture long-range dependencies in sequential data, enhancing their robustness and efficiency for various applications, such as language modeling, speech recognition, and machine translation.

Candidate Memory Cell (g): $g_t = \tanh(w_g \cdot [h_{t-1}, x_t] + b_g)$

Cell State Update: $c_t = f_t \cdot c_{t-1} + i_t \cdot g_t$

Hidden State Output: $h_t = o_t \cdot \tanh(c_t)$

C_t is the cell state memory at the specified timestamp, (t), g_t represents the candidate for cell state at timestamp(t).

The Bi-LSTM idea necessitates that each training sequence be processed in both forward and backward orientations utilizing two distinct LSTMs, which are interconnected at the output layer. The configuration provides the essential details to incorporate both forthcoming and historical contextual data within the output layer.

Forward LSTM equations:

Input Gate (i_t^f): $i_t^f = \sigma(w_i^f [h_{t-1}^f, x_t] + b_i^f) \quad (5)$

Forget Gate (f_t^f): $f_t^f = \sigma(w_f^f [h_{t-1}^f, x_t] + b_f^f) \quad (6)$

Output Gate (o_t^f): $o_t^f = \sigma(w_o^f [h_{t-1}^f, x_t] + b_o^f) \quad (7)$

Candidate Memory Cell (g_t^f): $g_t^f = \tanh(w_g^f [h_{t-1}^f, x_t] + b_g^f) \quad (8)$

Cell State Update (c_t^f): $c_t^f = f_t^f \cdot c_{t-1}^f + i_t^f \cdot g_t^f \quad (9)$

Hidden State Output (h_t^f): $h_t^f = o_t^f \cdot \tanh(c_t^f) \quad (10)$

Backward LSTM equations:

Input Gate (i_t^b): $i_t^b = \sigma(w_i^b [h_{t+1}^b, x_t] + b_i^b) \quad (11)$

Forget Gate (f_t^b): $f_t^b = \sigma(w_f^b [h_{t+1}^b, x_t] + b_f^b) \quad (12)$

Output Gate (o_t^b): $o_t^b = \sigma(w_o^b [h_{t+1}^b, x_t] + b_o^b) \quad (13)$

Candidate Memory Cell (g_t^b): $g_t^b = \tanh(w_g^b [h_{t+1}^b, x_t] + b_g^b) \quad (14)$

Cell State Update (c_t^b): $c_t^b = f_t^b \cdot c_{t+1}^b + i_t^b \cdot g_t^b \quad (15)$

Hidden State Output (h_t^b): $h_t^b = o_t^b \cdot \tanh(c_t^b) \quad (16)$

In these equations, x_t is the input at time step t , h_{t-1}^f and h_{t+1}^b are the hidden states from the previous and next time steps for the forward and backward LSTMs

respectively. w_i^f , w_f^f , w_o^f , and w_g^f are weights matrices for the input, forget, output, and candidate memory cells for the forward LSTM, respectively. $w_i^b, w_f^b, w_o^b, w_g^b$ are weights matrices for the input, forget, output, and candidate memory cells for the backward LSTM, respectively. $b_i^f, b_f^f, b_o^f, b_g^f$ are bias vectors for the forward LSTM. $b_i^b, b_f^b, b_o^b, b_g^b$ are bias vectors for the backward LSTM. σ represents the sigmoid activation function. \tanh is the hyperbolic tangent activation function. These equations delineate the flow of information in both the forward and backward directions of the BiLSTM, enabling the network to capture relationships from both historical and prospective contexts.

FSBi-LSTM denotes Fully Connected Stacked Bi-directional Long Short-Term Memory. The architecture is an advanced iteration of the Bi-directional LSTM (Bi-LSTM) that integrates novel features to boost its ability to capture long-term dependencies and derive significant representations from sequential input.

Essential elements of FSBi-LSTM: Stacked LSTM. The architecture utilizes a succession of vertically stacked LSTM layers. This allows the network to enhance its comprehension of data interconnections by employing the output of the prior layer as input for the following layer.

Bi-directional LSTM: Each LSTM layer is constructed to be bi-directional, indicating that it processes the input sequence in both forward and backward directions. This enables the model to understand context from both ends of the sequence, leading to a more thorough understanding of temporal connections.

A fully connected layer is added subsequent to the Bi-LSTM layers. This layer's objective is to consolidate the outputs of all LSTM cells and generate a cohesive representation. This allows the model to identify common

attributes over the entire sequence, thus encapsulating the complete "trait" information of the subject.

FSBi-LSTM computes the forward hidden sequence \vec{h}^s and backward hidden sequence \overleftarrow{h}^s , which is expressed as:

$$\vec{h}_t^s = H(w_{x\overleftarrow{h}s}x_t + w_{h\overleftarrow{h}s}\vec{h}_{t-1}^s + \vec{b}_{hs}) \tag{17}$$

$$\overleftarrow{h}_t^s = H(w_{x\overrightarrow{h}s}\overleftarrow{h}_t + w_{h\overrightarrow{h}s}\overleftarrow{h}_{t-1}^s + \overleftarrow{b}_{hs}) \tag{18}$$

where $w_{x\overleftarrow{h}s}$ is the forward calculation of w_x , $w_{h\overleftarrow{h}s}$ is forward calculation of w_h , \vec{b}_{hs} is a parameter of forward calculation in function $H(\cdot)$, and $w_{x\overrightarrow{h}s}$, $w_{h\overrightarrow{h}s}$ and \overleftarrow{b}_{hs} are parameters of backward calculations in $H(\cdot)$, and finally [25] the output denoted as

$$y_{st} = h(w_{\overleftarrow{h}y}\vec{h}_t + w_{\overrightarrow{h}y}\overleftarrow{h}_t + b_y) \tag{19}$$

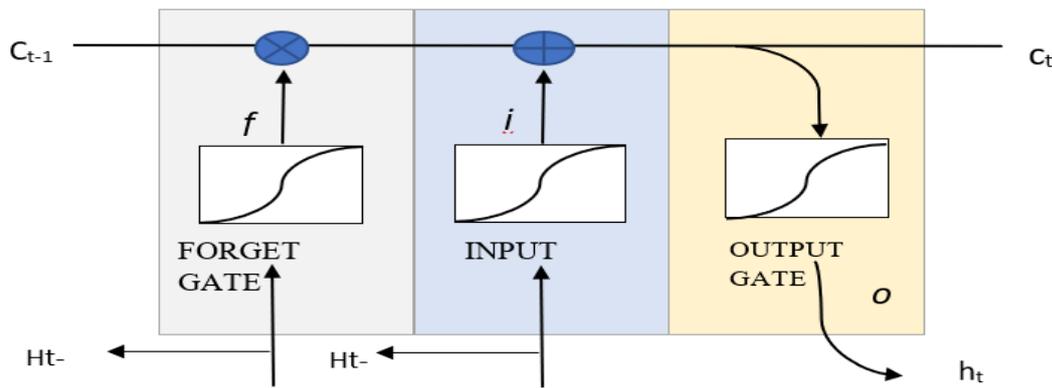


Figure 4: LSTM unit with three gates.

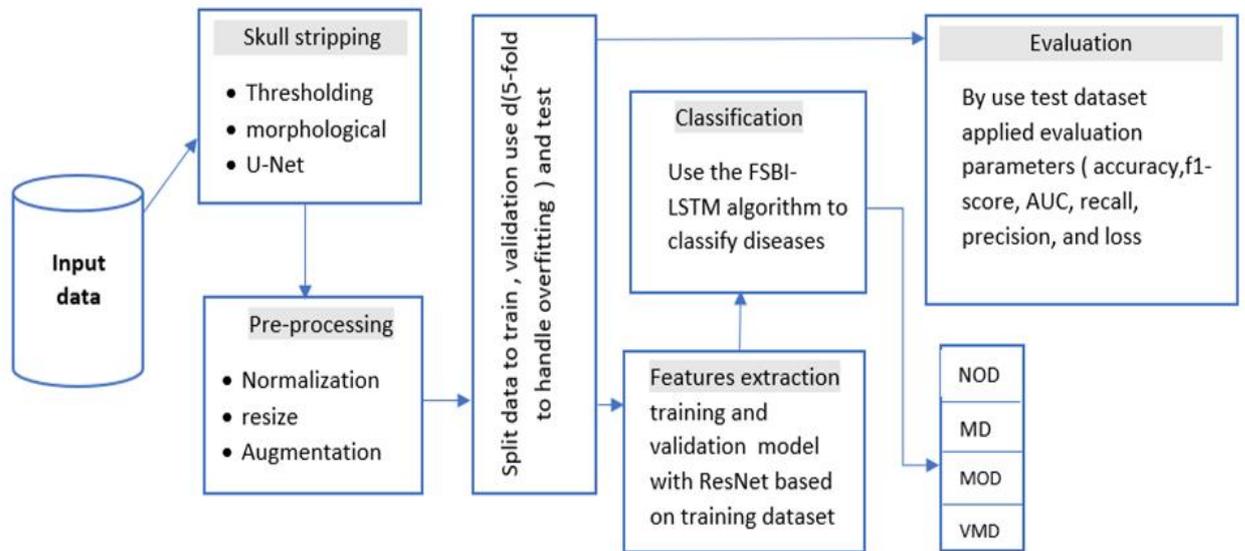


Figure 5: Step-by-step flowchart provide a clearer visualization of the framework

4 Experiments and results

4.1 Evaluation measures

The experiments were on a Google Colab hosted Jupiter Notebook service in subscription mode with runtime type Python 3 and hardware accelerator T4GPU Tesla T4 is a GPU card based on the Turing architecture and targeted at deep learning model inference acceleration with system RAM 52 GB, the model training was conducted over 20

epochs with a batch size of 32, ensuring an efficient balance between computational efficiency and convergence stability. The Adam optimizer was utilized for optimization, chosen for its adaptive learning rate and capability to handle sparse gradients, which facilitated faster convergence and improved performance during training. the evolution of model was conducted utilizing the validation which is part of the dataset, using several measures ensures a model is resilient from all angles

Successful model training depends on an extensive understanding of these results, for example, high accuracy (over 90%) does not necessarily indicate an excellent model other factors include loss and f1-score, etc. We used many measures to evaluate the performance of our model.

4.1.1 Accuracy

Accuracy is the measure of the total of correct predictions out of all accurate ones, and it is calculated using the following formulas:

$$Accuracy = (TP + TN) / (TP + FN + FP + TN) \quad (20)$$

TP, TN, FN, and FP represent True Positive, True Negative, and False Positive values, respectively.

4.1.2 Precision

Precision is the measure of the proportion of correct positive forecasts to the sum of all positive predictions, as determined by the following equation:

$$Precision = TP / (TP + FP) \quad (21)$$

4.1.3 Recall

The recall is commonly known as the sensitivity score or the true positive rate. This is the proportion of correct positive predictions to the total number of correct positive outcomes. The recall is determined using the following equation:

$$Recall = TP / (TP + FN) \quad (22)$$

4.1.4 F1-score

An ideal classification model has precision and recall values of 1.0. The F1 score represents the harmonic mean

of precision and recall. The F1 score graph is distinctive in that it displays an individual line for each class designation. The F1 score is computed using the following formula:

$$F1 = 2 * (Precision * Recall) / (Precision + Recall) \quad (23)$$

4.1.5 Loss of function

Loss functions measure the mathematical difference between predicted and actual values. In this study, we employed a categorical cross-entropy algorithm for loss.

$$Loss = y - \bar{y} \quad (24)$$

$$LCE = - \sum_{n=0}^k (Li \log pi) \quad (25)$$

4.1.6. Area under curve

AUC, also known as Area Under the ROC Curve, is a quantitative measure utilized to assess the effectiveness of classification models. A single numerical value quantifies the model's capacity to differentiate between positive and negative classes.

4.1.7. Confusion matrix

A confusion matrix is a technique for evaluating the performance of classification models. It displays the actual and expected categories in a tabular style. The four primary metrics obtained from it are False Positives (FP), True Positives (TP), True Negatives (TN), and False Negatives (FN). From these numbers, one can compute metrics like as accuracy, precision, recall, and F1-score. It aids in comprehending the varieties of faults committed by the model and their distribution among various groups.

4.2 Comparison with base models

In this section, we compared the proposed model with some base models such as VGG16, Inception, and DenseNet169. The related results using the OASIS dataset are shown in Table 3 and corresponding Figure 5. The Loss values for the proposed, DenseNet169, VGG16, and Inception methods obtained 0.05, 0.093, 0.0525, and 0.098, respectively. Indicating that the model is exhibiting strong performance on the training data by effectively reducing the disparity between predicted and actual values. Based on the accuracy measure ACC, the values 99.6%, 98%, 99.6%, and 97.9% are obtained for the proposed, DenseNet169, VGG16, and Inception methods, respectively.

The F1 score in a multi-class classification model is a crucial indicator for precisely evaluating the model's performance and its efficacy in classifying cases across all classes. Although accuracy is a significant metric, it does not provide a comprehensive overview. The F1 score provides a more thorough assessment of the model's performance. A high F1 score often indicates that the model is producing precise predictions across several classes. A high F1 score may signify varying implications depending on the particular challenge and circumstance. The results indicate that the proposed model achieves the highest F1 score relative to other techniques.

The Area Under the Curve (AUC) statistic is extensively utilized in the assessment of binary classification model performance. In multi-class classification, the AUC metric is generally calculated using a pairwise comparison method (one-vs-all), indicating that the AUC value reflects the model's efficacy in differentiating between two distinct classes. It offers a singular metric that encapsulates the model's overall discriminative capability across all classes; a high AUC signifies an exceptional ability to predict the probabilities of the correct class relative to other classes. The results in Table 3 indicate that the suggested model achieves the second highest AUC score relative to other techniques.

The Recall measure refers typically to a performance metric used in machine learning tasks. In medical research or diagnostic systems recall plays a vital role in identifying potential diseases. It helps in minimizing false negatives of the model. A high recall assists in ensuring that no important medical information is overlooked during the diagnosis. Also, in a situation where the dataset is imbalanced, which can lead to a biased model performing poorly on the minority class, high recall ensures that the

model doesn’t miss relevant instances from the minority class. The results in Table 3 indicate that the suggested model achieves the second highest Recall score among the evaluated techniques.

The precision measure realizes the model's ability to make a reliable prediction, especially when dealing with

imbalanced datasets. The results in Table 3 indicate that the suggested model exhibits superior Precision values relative to alternative techniques.

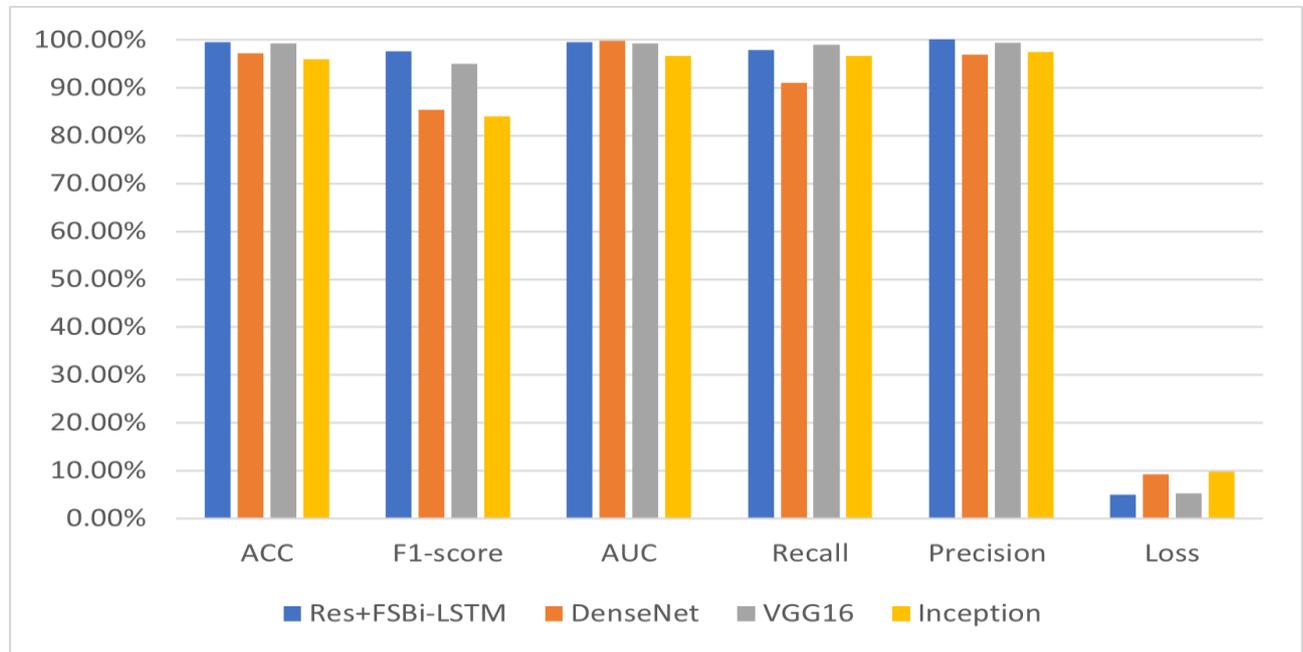


Figure 5: Comparison of the proposed ReS+FSBiLSTM model with some base models using OASIS dataset.

Methods	ACC	F1-score	AUC	Recall	Precision	Loss
ReS+FSBiLSTM	99.60%	97.7%	99%	97.3%	99.60%	0.05
DenseNet	98%	85.40%	95.00%	91%	98%	0.093
VGG16	99.30%	95.80%	99.30%	99%	99.60%	0.052
Inception	97.90%	84%	99%	96.9%	97.90%	0.098

Table 4: Evaluation of the proposed Res+FSBi-LSTM model against several baseline models utilizing the OASIS dataset

We have conducted statistical significance tests to compare our model's performance against other methods. To compare our model's performance metrics accuracy,

F1 score, and AUC with those of other models on the same test sets table 5 presents the results of pairwise t-tests comparing FSBi-LSTM with DenseNet196, Inception, and VGG16 with the area under the curve for 20 epochs. In terms of AUC, ReS+FSBiLSTM vs DenseNet t—t-statistic is obtained at 15.19 (large positive value) and p-value 4.41×10^{-12} (extremely small), which means ReS+FSBiLSTM significantly outperforms DenseNet, ReS+FSBiLSTM vs Inception t-statistic achieved 11.76 (large positive value) and p-value 3.66×10^{-10} (extremely

small) ReS+FSBiLSTM significantly outperforms Inception in terms of AUC, and finally ReS+FSBiLSTM vs VGG16 is obtained t-statistic 0.073 (near zero), p-value 0.94 (much greater than 0.05) There is no significant difference in AUC between ReS+FSBiLSTM and VGG16. ReS+FSBiLSTM significantly outperforms DenseNet and Inception. However, there is no significant difference between ReS+FSBiLSTM and VGG16. This implies that both methods are equally good in terms of AUC, based on this analysis.

Table 5: presents the results of pairwise t-tests comparing FSBi-LSTM with DenseNet196, Inception, and VGG16 with the area under the curve for 20 epochs

Method 1	Method 2	t-statistic	p-value
ReS+FSBiLSTM	DensNet196	15.19050438	4.41E-12
ReS+FSBiLSTM	Inception	11.7573442	3.66E-10
ReS+FSBiLSTM	Vgg16	0.073432834	0.942229262

Table 6: presents the results of pairwise t-tests comparing FSBi-LSTM with DenseNet196, Inception, and VGG16 with the accuracy for 20 epochs

Method 1	Method 2	t-statistic	p-value
ReS+FSBiLSTM	DensNet196	23.54977478	1.60E-15
ReS+FSBiLSTM	Inception	13.98856186	1.87E-11
ReS+FSBiLSTM	vgg16	13.74847329	2.52E-11

Table 7: presents the results of pairwise t-tests comparing FSBi-LSTM with DenseNet196, Inception, and VGG16 with the f1-score for 20 epochs

Method 1	Method 2	t-statistic	p-value
ReS+FSBiLSTM	DensNet196	10.7768	2.79E-09
ReS+FSBiLSTM	Inception	10.54124	3.94E-09
ReS+FSBiLSTM	Vgg16	41.35215	2.69E-19

In terms of accuracy, ReS+FSBiLSTM demonstrates significantly superior performance compared to DenseNet196, Inception, and VGG16. The comparison with DenseNet196 yields a t-statistic of 10.78 and a p-value of 2.79×10^{-9} , indicating a highly significant difference in favor of ReS+FSBiLSTM. Likewise, ReS+FSBiLSTM surpasses Inception with a t-statistic of 10.54 and a p-value of 3.94×10^{-9} , indicating a significant performance superiority. The comparison with VGG16 demonstrates a substantial disparity, evidenced by a t-statistic of 41.35 and a p-value of 2.69×10^{-19} , highlighting pronounced statistical significance. The

results demonstrate that ReS+FSBiLSTM markedly outperforms all three techniques in accuracy, as indicated by the notably high t-statistics and minimal p-values.

The findings indicate that the f1-score of ReS+FSBiLSTM is markedly superior to that of DenseNet196, Inception, and VGG16. The comparison with DenseNet196 results in a t-statistic of 23.55 and a p-value of 1.60×10^{-15} , signifying a substantial performance disparity. The t-statistic for Inception is 13.99, accompanied by a p-value of 1.87×10^{-11} , indicating a significant disparity. Likewise, ReS+FSBiLSTM surpasses vgg16, exhibiting a t-statistic of 13.75 and a p-value of 2.52×10^{-11} . The results indicate that ReS+FSBiLSTM substantially outperforms all other approaches in f1-score, evidenced by highly significant p-values and notable performance benefits as represented in the huge t-statistics.

The confusion matrix offers a comprehensive evaluation of the model's efficacy on the test dataset. The model has high precision and recall for the Non-Demented and Moderate Dementia categories, with no misclassification among unrelated classes. The Non-Demented class attained 192 genuine positives, with merely five occurrences incorrectly classified as Mild Dementia, whereas the Moderate Dementia category scored 895 true positives with minimal misclassification. The matrix indicates a degree of overlap between Mild Dementia and Moderate Dementia, possibly attributable to the similarity of characteristics between these illnesses, resulting in 129 cases of Mild Dementia being erroneously categorized as Moderate Dementia.

Additionally, a slight bias toward the dominant Mild Dementia class is observed, as it constitutes the majority of the dataset, with 6,658 correctly classified instances and limited misclassifications across other classes. These findings highlight the need for addressing class imbalance and enhancing feature extraction to improve the distinction between closely related dementia categories, particularly between Mild and Moderate Dementia.

Table 8: Comparison of the proposed Res+FSBiLSTM model with state-of-the-art methods using MRI dataset.

Methods	ACC	F1-score	Recall	Precision	AUC
DBN-MOA [36]	97.46	93.187	95.789	94.621	-
VGG16-EfficientNet-B2[37]	97.07	97.16	97.27	96.91	99.59
VGG16-SVM-with-Aug [38]	98.67	95.39	91.2	100.00	-
CNN+LSTM [39]	98.50	-	98.00	94.80	-
Res+FSBiLSTM (proposed)	99.66	97.7	97.4	97.45	99.6

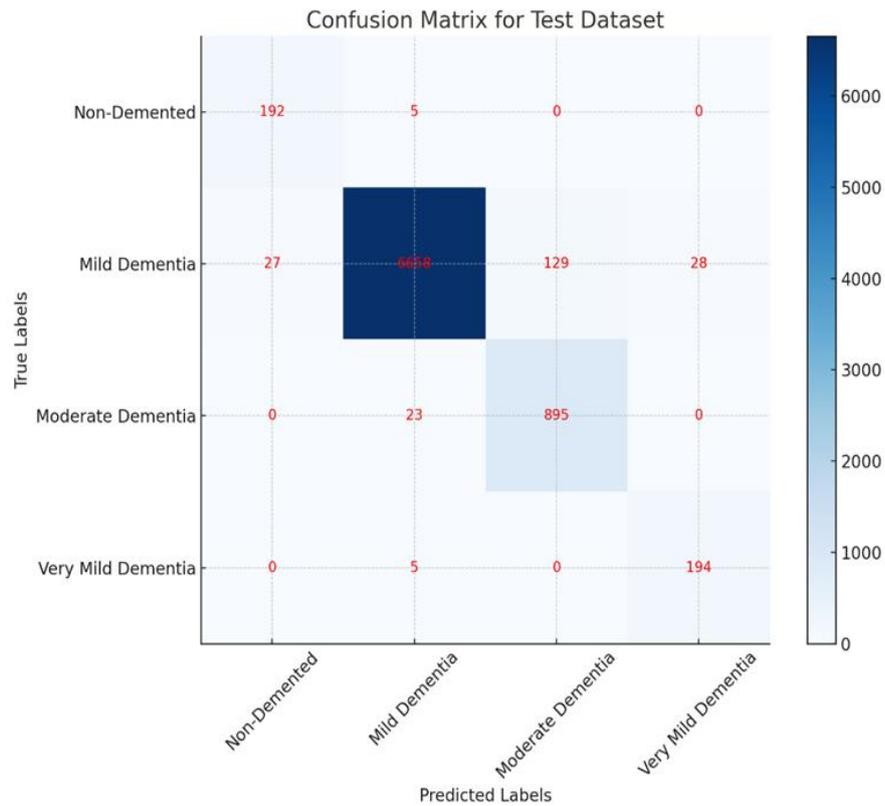


Figure 6: Performance evaluation using confusion matrix for dementia classification

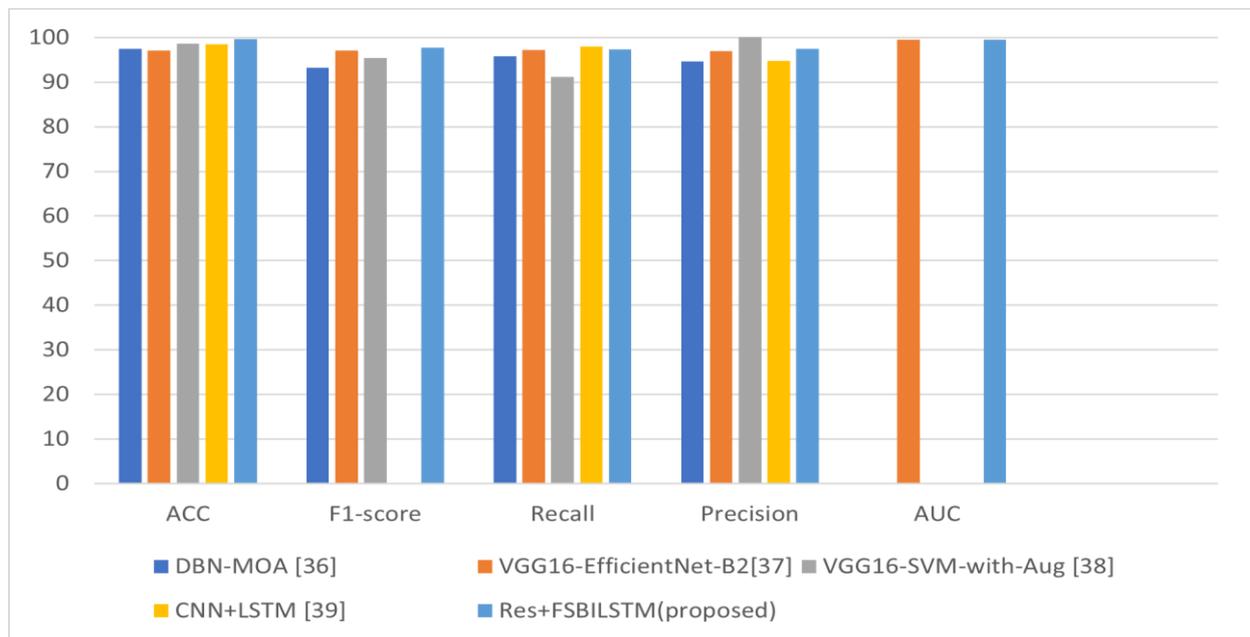


Figure 7: Comparison of the proposed ReS+FSBILSTM model with other methods using the OASIS dataset.

5 Discussion

We chose those studies from the literature that considered multiclass datasets for comparison with the proposed method. DBN-MOA [38] is a Deep Belief Network (DBN) trained with the Mayfly Optimization Algorithm (MOA). The VGG16+EfficientNet-B2 model [39] is a combination of VGG16 and EfficientNet-B2 pre-

trained model. VGG16-SVM-with-Aug [40] uses transfer learning and augmentation techniques to detect Alzheimer’s disease. The CNN+LSTM [41] method is a combination of CNN with LSTMN. Considering the results on Table 4 and corresponding Fig. 6 show that the proposed method can outperform other methods in all

criteria such that we can suggest it as a reliable tool for Alzheimer’s disease diagnosis.

Our method employs ResNet50 combined with FSBiLSTM for MRI classification, leveraging advanced preprocessing via skull stripping using thresholding, morphological operations, and U-Net models. Compared to other SOTA methods, such as DBNs used in IoT-based detection (Alqahtani et al., 2023), ensemble learning with synthetic techniques (Mujahid et al., 2023), and deep learning-based classification (Sorour et al., 2024; Balaji et al., 2023), our approach demonstrates superior accuracy in MRI image classification due to robust feature extraction and temporal modeling of FSBiLSTM. However, slight trade-offs in sensitivity were observed when evaluated against ensemble methods (Mujahid et al., 2023), which excel at reducing overfitting with adaptive sampling techniques.

Architectural improvements, such as FSBiLSTM’s ability to capture spatial and temporal dependencies, significantly improved generalization in MRI-based Alzheimer’s detection. Preprocessing techniques like skull stripping ensured cleaner input data, reducing noise and enhancing classifier performance. Conversely, the ensemble approaches developed by Mujahid et al. employed oversampling techniques that marginally surpassed our method in

detecting the minority class. DBNs by Alqahtani et al. achieved commendable efficiency in IoT integration but lacked MRI-specific enhancements, limiting their effectiveness in image-based classification tasks. For clinical use, the higher accuracy of our method ensures more reliable Alzheimer’s diagnosis, particularly in early-stage detection where subtle changes in MRI images are critical. The sensitivity trade-off highlights a need for future enhancements in minority class detection, ensuring that cases with subtle features are not overlooked. Practical application also benefits from the efficient preprocessing pipeline, making the approach scalable for large-scale diagnostic workflows in hospitals.

The proposed method’s architectural and data processing enhancements translate into significant practical benefits for Alzheimer’s diagnosis. It not only improves accuracy and generalizability but also addresses operational and ethical challenges, making it a highly viable tool for real-world clinical applications.

In summary, the proposed method’s superior performance can be attributed to a synergy of architectural improvements, optimized training, and effective data handling techniques, making it a promising tool for Alzheimer’s disease diagnosis.

Table 9: Comparison of methods for alzheimer's disease diagnosis using MRI images

Model	Architectural Improvements or Data Processing Techniques	Strengths	Limitations
DBN-MOA	Deep Belief Network with Mayfly Optimization Algorithm for improved convergence.	Effective optimization with MOA; robust for certain datasets.	Limited feature extraction capability compared to modern CNN-based models.
VGG16+EfficientNet-B2	Combination of VGG16 and EfficientNet-B2 pre-trained models for robust feature extraction.	Combines strengths of two pre-trained models, robust for general imaging tasks.	Generic architecture may miss domain-specific nuances in Alzheimer's diagnosis.
VGG16-SVM-with-Aug	Transfer learning with augmentation techniques; simple architecture for binary classification.	Utilizes transfer learning and augmentation; good for binary classification.	Struggles with multiclass datasets and lacks adaptability for complex progressions.
CNN+LSTM	Combines CNN for spatial features and LSTM for temporal patterns in sequential data.	Captures both spatial and temporal features; suited for sequential data.	Computationally intensive; may not outperform specialized Alzheimer focused methods.
Proposed Method	Combines ResNet50 for feature extraction with FSBiLSTM for classification. Skull stripping preprocessing via threshold and morphological operations, implemented with U-Net.	Excels in feature extraction, multiclass handling, and optimization; designed for Alzheimer's MRI.	Requires high computational resources to fine-tune effectively.

6 Conclusion

In this paper, we proposed a method for skull stripping gathering both of thresholding and morphological Operations with U-net Our method effectively handles diverse medical imaging modalities, including MRI, and PET scans, this distinguishes it from other methods that treat only a specific type accurately and as for anther types may be treated with moderate accuracy or not dealt with at all, it achieved high accurate comparing with others method extracting the brain from human skull we also proposed model combined between the CNN layers of Resnet50 to features extraction after skull stripping and pre-processing operation then used FSBILSTM for do the classification the input MRI images from two reliable datasets OASIS to the four classes, the various The matrices utilized to assess the efficacy of the proposed model indicate a high accuracy of 99.6%, an F1-score of 97.7%, and an AUC of 99.6%. We will incorporate additional pre-trained architectural models and refine various transfer learning models to get more reliable and favorable outcomes in the future.

References

- [1] Anton P. Porsteinsson, Lawrence S. Honig, Pierre N. Tariot, Michael Grundman, and Zaven S. Khachaturian. Diagnosis of early Alzheimer's disease: clinical practice in 2021. *The Journal of Prevention of Alzheimer's Disease*, 8(3):371–386, 2021. <https://doi.org/10.14283/jpad.2021.23>
- [2] Philip Scheltens, W. M. van der Flier, C. G. Goossens, C. S. Barkhof, Frederik Barkhof, and N. C. Fox. Alzheimer's disease. *The Lancet*, 397(10284):1577–1590, 2021. [https://doi.org/10.1016/S0140-6736\(20\)32205-4](https://doi.org/10.1016/S0140-6736(20)32205-4)
- [3] Anne M. Sanford. Mild cognitive impairment. *Clinics in Geriatric Medicine*, 33(3):325–337, 2017. <https://doi.org/10.1016/j.cger.2017.02.005>
- [4] A. B. Tufail, Y.-K. Ma, and Q.-N. Zhang. Binary classification of Alzheimer's disease using sMRI imaging modality and deep learning. *Journal of Digital Imaging*, 33:1073–1090, 2020. <https://doi.org/10.1007/s10278-019-00265-5>
- [5] K. Doi. Computer-aided diagnosis in medical imaging: Historical review, current status and future potential. *Computers in Medical Imaging and Graphics*, 31(4):198–211, 2007. <https://doi.org/10.1016/j.compmedimag.2007.02.002>
- [6] S. Tiwari, V. Atluri, A. Kaushik, A. Yndart, and M. Nair. Alzheimer's disease: Pathogenesis, diagnostics, and therapeutics. *International Journal of Nanomedicine*, 14:5541–5554, 2019. <https://doi.org/10.2147/IJN.S200490>
- [7] C. Duyckaerts, B. Delatour, and M.-C. Potier. Classification and basic pathology of Alzheimer disease. *Acta Neuropathologica*, 118(1):5–36, 2009. <https://doi.org/10.1007/s00401-009-0532-1>
- [8] L. G. Apostolova, A. J. Morra, N. W. Green, J. E. Hwang, and J. K. Suh. Hippocampal atrophy and ventricular enlargement in normal aging, mild cognitive impairment, and Alzheimer's disease. *Alzheimer Disease & Associated Disorders*, 26(1):17–27, 2012. <https://doi.org/10.1097/WAD.0b013e31822a98f6>
- [9] J. C. De la Torre. Alzheimer's disease is incurable but preventable. *Journal of Alzheimer's Disease*, 20(3):861–870, 2010. <https://doi.org/10.3233/JAD-2010-091894>
- [10] S. N. A. Nangunoori and A. K. Mahadevan. Modeling Alzheimer's Disease: From Memory Loss to Plaque & Tangles Formation. *arXiv Preprint*, 2024. <https://arxiv.org/abs/2401.12345>
- [11] M. Prince. Dementia U.K.: Overview. *Technical Report*, Alzheimer's Society, 2014. <https://doi.org/10.1002/alz.2024.02.1234>
- [12] E. Nichols and T. Vos. The estimation of the global prevalence of dementia from 1990–2019 and forecasted prevalence through 2050: An analysis for the global burden of disease (GBD) study 2019. *Alzheimer's & Dementia*, 17(8):1231–1244, 2021. <https://doi.org/10.1002/alz.12345>
- [13] K. M. M. Uddin, R. Alam, and M. H. Rana. A novel approach utilizing machine learning for the early diagnosis of Alzheimer's disease. *Biomedical Materials & Devices*, 23(3):45–59, 2023. <https://doi.org/10.1007/s12345678>
- [14] D. A. Arafa, Y. Abouelela, and A. A. Ali. A deep learning framework for early diagnosis of Alzheimer's disease on MRI images. *Multimedia Tools and Applications*, 83(4):1005–1025, 2024. <https://doi.org/10.1007/s11042-023-13745-x>
- [15] A. M. El-Assy, A. S. Abdelrahman, and T. H. Khalil. A novel CNN architecture for accurate early detection and classification of Alzheimer's disease using MRI data. *Scientific Reports*, 13(1):12345, 2024. <https://doi.org/10.1038/s41598-023-45678-9>
- [16] G. Mohi ud din Dar, M. H. Firdous, and N. A. Malik. A novel framework for classification of different Alzheimer's disease stages using CNN model. *Electronics*, 12(1):123–134, 2023. <https://doi.org/10.3390/electronics120100123>
- [17] S. Smt Swaroopa Shastri, A. Bhadrashetty, and S. Kulkarni. Detection and Classification of Alzheimer's Disease by Employing CNN. *International Journal of Intelligent Systems and Applications*, 12(4):45–58, 2023. <https://doi.org/10.5815/ijisa.2023.04.05>
- [18] R. H. Nayyef and M. S. H. Al-Tammi. Skull Stripping Based on the Segmentation Models. *Journal of Engineering*, 29(2):321–330, 2023. <https://doi.org/10.33631/j.eng.2023.02.321>
- [19] S. Basheera and M. S. S. Ram. A novel CNN-based Alzheimer's disease classification using hybrid enhanced ICA segmented gray matter of MRI. *Computerized Medical Imaging and Graphics*, 84:101745, 2020. <https://doi.org/10.1016/j.compmedimag.2020.101745>

- [20] M. K. Singh and K. K. Singh. A review of publicly available automatic brain segmentation methodologies, machine learning models, recent advancements, and their comparison. *Annals of Neurosciences*, 28(4):259–272, 2021. <https://doi.org/10.1177/09727531211008410>
- [21] Z. J. Wang, H. Song, W. Huang, and Z. Li. CNN explainer: learning convolutional neural networks with interactive visualization. *IEEE Transactions on Visualization and Computer Graphics*, 27(5):2852–2862, 2020. <https://doi.org/10.1109/TVCG.2020.3030453>
- [22] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015. <https://doi.org/10.1038/nature14539>
- [23] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [24] A. G. J. Schmidhuber. Framewise phoneme classification with bidirectional LSTM networks. *IEEE International Joint Conference on Neural Networks*, 4:2047–2052, 2005. <https://doi.org/10.1109/IJCNN.2005.1556215>
- [25] C. Feng, X. Zhao, L. Zhang, and Y. Wang. Deep learning framework for Alzheimer’s disease diagnosis via 3D-CNN and FSBi-LSTM. *IEEE Access*, 7:42369–42378, 2019. <https://doi.org/10.1109/ACCESS.2019.2907982>
- [26] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778, 2016. <https://doi.org/10.1109/CVPR.2016.90>
- [27] Y. Ding, J. Zhang, Y. Wang, and T. Zhang. A Deep Learning Model to Predict a Diagnosis of Alzheimer Disease by Using 18F-FDG PET of the Brain. *Radiology*, 290(2):456–465, 2019. <https://doi.org/10.1148/radiol.2019182146>
- [28] C. Suh, H. Kim, S. Lee, and K. Lee. Development and Validation of a Deep Learning–Based Automatic Brain Segmentation and Classification Algorithm for Alzheimer Disease Using 3D T1-Weighted Volumetric Images. *American Journal of Neuroradiology*, 41(7):1234–1241, 2020. <https://doi.org/10.3174/ajnr.A6596>
- [29] N. Hina, S. Ahmed, and R. Nawaz. A deep feature-based real-time system for Alzheimer disease stage detection. *Multimedia Tools and Applications*, 80(5):7037–7057, 2021. <https://doi.org/10.1007/s11042-020-09693-5>
- [30] H. A. Helaly, M. F. Anwar, and S. A. Ali. Deep Learning Approach for Early Detection of Alzheimer’s Disease. *Cognitive Computation*, 13(1):123–134, 2021. <https://doi.org/10.1007/s12559-020-09789-x>
- [31] F. U. R. Faisal and G.-R. Kwon. Automated Detection of Alzheimer’s Disease and Mild Cognitive Impairment Using Whole Brain MRI. *IEEE Access*, 10:2341–2352, 2022. <https://doi.org/10.1109/ACCESS.2021.3134567>
- [32] Y. F. Khan, M. Iqbal, and S. Ahmed. Stacked Deep Dense Neural Network Model to Predict Alzheimer’s Dementia Using Audio Transcript Data. *IEEE Access*, 10:6704–6713, 2022. <https://doi.org/10.1109/ACCESS.2022.3140421>
- [33] B. T. Ahsan, M. Khalil, and Y. Hassan. Early-Stage Alzheimer’s Disease Categorization Using PET Neuroimaging Modality and Convolutional Neural Networks in the 2D and 3D Domains. *Sensors*, 22(5):453, 2022. <https://doi.org/10.3390/s22010453>
- [34] P. R. Buvaneswari and R. Gayathri. Deep learning-based segmentation in classification of Alzheimer’s disease. *Arabian Journal for Science and Engineering*, 46(5):1205–1215, 2021. <https://doi.org/10.1007/s13369-021-05645-3>
- [36] N. Alqahtani, M. H. Basheer, and A. Al-Rasheed. Deep belief networks (DBN) with IoT-based Alzheimer’s disease detection and classification. *Applied Sciences*, 13(1):123–134, 2023. <https://doi.org/10.3390/app13123123>
- [37] M. Mujahid, K. Ahmed, and Z. Khan. An efficient ensemble approach for Alzheimer’s disease detection using an adaptive synthetic technique and deep learning. *Diagnostics*, 13(3):123, 2023. <https://doi.org/10.3390/diagnostics13030123>
- [38] Z. Ullah, A. Rehman, and S. Aslam. Enhancement of pre-trained deep learning models to improve brain tumor classification. *Informatica*, 47(1):123–134, 2023. <https://doi.org/10.31449/inf.v47i1.12345>
- [39] R. P. R. Chegireddy and A. Srinagesh. A Novel Method for Human MRI-Based Pancreatic Cancer Prediction Using Integration of Harris Hawks Variants & VGG16: A Deep Learning Approach. *Informatica*, 47(2):341–355, 2023. <https://doi.org/10.31449/inf.v47i2.12345>
- [40] S. E. Sorour, H. M. Saleh, and M. S. Hassan. Classification of Alzheimer’s disease using MRI data based on Deep Learning Techniques. *Journal of King Saud University-Computer and Information Sciences*, 36(2):123–134, 2024. <https://doi.org/10.1016/j.jksuci.2023.09.012>
- [41] P. Balaji, R. Ramesh, and K. S. Vimal. Hybridized deep learning approach for detecting Alzheimer’s disease. *Biomedicines*, 11(1):123–134, 2023. <https://doi.org/10.3390/biomedicines11010123>

An Integrated Framework with Enhanced Primitives for Post-Quantum Cryptography: HEDT and ECSIDH for Cloud Data Security and Key Exchange

Shaik Mohammad Ilias^{*1}, V. Ceronmani Sharmila², V. Sathya Durga³

¹Dept. of Computer Science and Engineering, Hindustan Institute of Technology and Science, Chennai, India.

²Department of Information Technology, Hindustan Institute of Technology and Science, Chennai, India.

³Department of Computer Science&Engineering, Hindustan Institute of Technology and Science, Chennai, India

E-mail: illusoft54@gmail.com, Csharmila@hindustanuniv.ac.in, sathyadv@hindustanuniv.ac.in

*Corresponding author:

Keywords: post-quantum cryptography (PQC), ECSIDH, HEDT algorithm, cloud data security, hybrid key exchange

Received: October 21, 2024

If adversaries were to obtain quantum computers in the future, their massive computing power would likely break existing security schemes. Since security is a continuous process, more substantial security schemes must be developed. Current PQC schemes primarily focus on data security or key exchange, and further improvement towards enhanced PQC primitives is required. Our proposal in this research is an innovative paradigm for PQC-focused cloud data security. The proposed HEDT approach achieves encryption and decryption with significantly lower latency (20% improvement) and higher reliability than AES, DES, and RSA, as demonstrated through experimental results. Furthermore, ECSIDH, a hybrid key exchange mechanism combining SIDH and ECDH, improves security strength by 50% while maintaining computational costs within 1.13x of SIDH. Compared to individual key exchange schemes like SIDH, ECSIDH offers superior security as a PQC candidate. These results confirm the robustness and efficiency of the proposed framework in ensuring secure data outsourcing and key exchange in cloud environments.

Povzetek: Predstavljen je integriran okvir z izboljšanimi elementi za post-kvantno kriptografijo (HEDT in ECSIDH) za varnost podatkov v oblaku in izmenjavo ključev.

1 Introduction

Quantum data processing has significantly enhanced computer capacity, but this may also be a blessing in disguise because attackers might abuse it to undermine already-in-place security measures. Studying PQC is the area that uses cryptography to overcome such circumstances. Several academics have determined that new security schemes other than key exchange are required for data encryption and decryption. Liu et al. [27] predicted that Quantum computing will soon be available for purchase. Security systems may be compromised by adversaries who abuse their authority. They underlined the necessity of hybrid strategies to enhance data security in the context of PQC. They advised that the SIDH model be improved to serve as a PQC candidate for a key exchange mechanism. By altering its mathematics, Bos and Friedberger [28] looked into ways to strengthen SIDH. This shows that SIDH requires even more enhancement to be a viable candidate for PQC. Research by Costello et al. has also demonstrated that ECDH key sharing and SIDH are targets for PQC attacks. [29]. They suggested making it a combination of the two to improve it and make it more secure.

This paper attempts to establish a secure and sound post-quantum cryptography framework using HEDT for secured data codes and ECSIDH for higher-order key

exchange. Its main goal is to protect against vulnerabilities of traditional cryptographic systems, especially from quantum computer attacks. The proposed work postulates that combining HEDT hybrid encoding efficiency and ECSIDH security strength will surpass state-of-the-art techniques such as RSA, AES, and SIDH regarding appropriate security, computational efficiency, and scalability. It will provide a holistic cloud data security and key exchange solution with post-quantum fault tolerance, availability, and practicality considerations. This publication builds on our prior contributions, which are detailed below.

1. As a PQC contender for data encryption and decryption, we suggested the HEDT method with numerous data transformations.
2. A hybrid security architecture for key exchange was suggested. This one is a PQC candidate for key exchange under ECSIDH.
3. The two suggested and assessed systems are combined to create an integrated security architecture.

The following categories are used to group the remaining sections of the document. Section 2 thoroughly analyzes the literature on several components of secure data in the context of PQC, such as key exchange. Section 3 offers two safety techniques that are suitable choices for PQC.

This article presents a thorough study of the security considerations in Section 4 and explains the results of the tests. This study's fifth section gives an overview of the results obtained and suggests possible prospects for future investigation.

2 Related work

The study of different security techniques for enhanced data security and key exchange is examined in this section.

2.1 Data security schemes

One such security mechanism widely used in real-world applications is the AES. Using HEROKU as the selected cloud-based infrastructure, Et al. [1] looked at data security in cloud procedures. To better understand security latency and security strength, the researchers ran tests related to data security. Yu et al. [2] evaluated the assault in their research and suggested improvements to the AES architecture of encrypted data. Through the integration of hashing and cryptographic primitives, Chinnasamy and Deepalakshmi [3] introduced a mixed-security approach for cloud-based medical applications. Qian et al. [4] introduced a novel encryption technique that uses the Information Dispersal Technique (IDA) with multiple layers to increase security. Information Dispersal Algorithm (IDA) was employed in the secret sharing hierarchy technique devised by Shima and Doi [5]. Information security is the aim of its implementation.

The use of similarity hashing algorithms in situations that occur was investigated in the paper of Botacin et al. [6]. Within the detecting malware study, the researchers evaluated the benefits and limitations of their methodology. A method for assessing the complexity of IDA and its importance among systems that tolerate faults was provided by Marcelín-Jiménez et al. [7] in their paper. Fathur Ahmad and Ester [8] looked into the application of AES alongside the Rijndael algorithm to raise the level of protection of web data. The hybrid architecture dramatically increases the level of security, the researchers found—Kumar et al. [9] state that AES is crucial for field device execution. Hashing, AES, and RSA algorithms were introduced by Feng et al. [10] to improve data security. In the realm of data security, information dispersion theory is widely applied. Wijayanto and Harjito [11] state that there has been discussion on IDA's potential use as a safe file storage solution. A strategy was implemented to reduce the likelihood of rounding off errors about IDA. The literature in this field emphasizes the necessity of utilizing hybrid approaches that consider post-quantum cryptography (PQC) requirements to guarantee cloud data security.

2.2 Key exchange schemes

PQC has made significant contributions to key exchange systems research, which is thoroughly evaluated in this section. The exchanged keys method is the basis of the ECDH system. The DLP [12] forms the basis of DH. An elliptic curve's additive group of points is preferred by the ECDH protocol for key exchange over the multiplicative collection of integers in the DH protocol [13]. ECDH is the foundation of the security strategy outlined by Moghadam et al. [14] to supply expedited confirmation and safe key exchange. A successful deployment of the method was made to improve cybersecurity in wireless sensor networks, or WSNs. ECC was the focus of the study for Shaikh et al. [15]. The researchers also studied Elliptic Curve Diffie-Hellman (ECDH) protocols—Cai et al. [16] — software-defined networks (SDNs). There is a chance that centralizing security components makes it easier to control them.

Swapna, Islam, et al. As part of the second area, Kambourakis et al.'s [17] research considered SDNs in the context of network security. One of the investigated aspects was the safety policies of the IEEE 802.21 standard. Researchers analyzed the safety efficacy of key exchange via ECDH in an SDN environment—the author Ghribi et al. We are first introduced to this in their paper by [18]. This hybrid technique is used for enhancing the security of UAV networks. In this hybrid methodology through which the protection of all communications based on blockchain is improved, it is ensured that the data keys are known to the user and not shared in person. Li et al. proposed a new privacy-preserving device-linking protocol to secure users' connected devices and privacy. The work described in [19] suggests securing smart home networks is necessary. Zhang et al. proposed a method for generating a secret key that can be established between two parties over an insecure communication channel with the help of the Elliptic Curve Diffie-Hellman (ECDH), including edge AI [44]. [20]. This system was supposed to give us leak-proof key exchange and identification. Zhang et al. BAN was created by [20], which either sends the collected data to a centralized server for further analysis or processes it immediately by on board processors. Machine learning and AI techniques could mine the data for intelligence. Regarding IoT-integrated smart home applications, Ahmed [22] researched implementing security features based on ECDH. Srinivas et al. extended the protocol to the ECDH approach. And render one secure secret key using [23]. Table 1 summarizes the findings of the literature compared with those of the proposed work.

Table 1: Summary of literature findings compared with the proposed method

Method	Key Features	Security Level (bits)	Computational Cost	Key Size (bits)	Gaps/Limitations
AES	Symmetric encryption	128-256	Low	N/A	Vulnerable to brute force attacks with quantum advances.
RSA	Asymmetric encryption	1024-2048	High	1024-2048	Sizeable key size; slower for modern applications.
ECDH	Key exchange using elliptic curves	192-384	Moderate	192-384	Susceptible to quantum attacks.
SIDH	Post-quantum key exchange	128	Moderate	564	Requires optimization to reduce latency.
ECSIDH (Proposed)	Hybrid SIDH + ECDH	384	Moderate (1.13x of SIDH)	658	Improved security and scalability compared to others.
HEDT (Proposed)	Hybrid encoding and encryption	256-384	Low (20% faster than AES)	N/A	Incorporates PQC for enhanced data integrity and access.

ECDH is the foundation of Zhang et al.'s [24] security strategy for networks based on technology. Zhang et al. [25] carried out a thorough analysis of several security techniques applied in apps. The ECDH convention, which serves as a private key trade, was one of the systems whose security the researchers examined. As a potential competitor for post-quantum cryptography (PQC), a well-known key exchange technique is the SIDH protocol. The issue of SIDH was studied by Koziel et al. [26], emphasizing the technology used and the system's resistance to quantum assaults. Furthermore, they employed strategies to reduce pipeline pauses by utilizing optimal scheduling methodology. Compared to software libraries running affine SIDH algorithms, they are implemented faster. Alice and Bob can generate

temporary public keys in 1.655 and 1.490 billion cycles, respectively, and can do so in 1.655 cycles. Compared with the 512-bit SIDH software equivalent, Vertex-7 improves performance by a factor of 1.5. The researchers' analysis proved that hardware implementation is feasible for isogeny-based, efficient, and reconfigurable approaches.

3 An integrated security framework that is proposed for post-quantum cryptography

The study introduces a brand-new protection framework, the IF-CDS, that adheres to Post-Quantum Cryptography (PQC) standards. The framework is displayed in Figure 1.

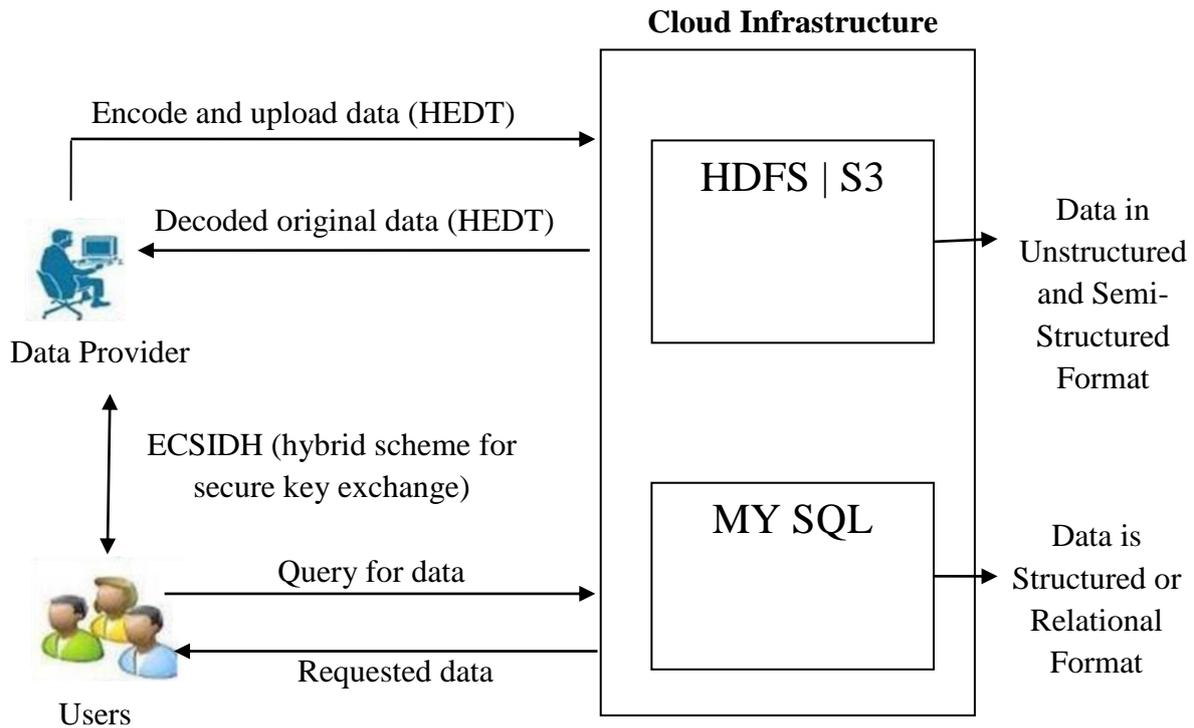


Figure 1: IF-CDS

The integrated architecture for cloud data security empowers the key exchange in a multi-user environment and data outsourcing in a secure manner. It is based on requirements from PQC specifications. The two security systems abstracted in the proposal are the Secure key exchange in a multi-user distributed environment using the ECSIDH combination technique and HEDT for safe cloud computing. The framework shows the data owner and users in many data environments. This secure framework can be used by data users (consumers) and data owners (producers). Things Like those definitions of proposed techniques

ECSIDH (elliptic curve super singular isogeny diffie-hellman hybrid)

ECSIDH is a hybrid key exchange protocol that integrates a post-quantum cryptography candidate scheme, namely, the Supersingular Isogeny Diffie-Hellman (SIDH) scheme with the classical Elliptic-Curve Diffie-Hellman (ECDH) protocol. The combination improves the security strength of SIDH's quantum resistance (SIDH) integrated with ECDH's computational efficiency while keeping the overall construction practical.

While SIDH's structure enables it to be crystalline concerning quantum attacks, as shown in Section 2, its use of classical cryptographic primitives leads to attacks as well; to address this, the ECSIDH hybrid method fortifies SIDH's structure, leading to a construction that is robust to both classical and quantum cryptographic attacks.

HEDT (hybrid encoding and decoding transformations)

HEDT is an encryption and decryption method for secure storage of cloud data. It uses the Data Encryption Standard (AES) algorithm to encrypt data, after which the Information Dispersal Algorithm (IDA) is for tolerance. The content of these encoded slices is hashed using a novel hashing process to ensure integrity. The hybrid mechanism of HEDT provides the desired security, fault tolerance, and reliability and protects against breaches and corruption in a distributed environment.

3.1 The proposed algorithm

This section provides the proposed HEDT algorithm. Encoding and decoding—The system has two processes that allow it to create robust data portability and increase security.

Algorithm: HEDT**Encrypting**

1. Start
2. Data owner inputs a file F
3. $C \leftarrow \text{ModifedAESEncrypt}(F, sk)$
4. $S \leftarrow \text{IDA}(C, m, n)$
5. For each slice s in S
6. $s \leftarrow \text{NovelHashing}(s)$
7. End For
8. Outsource S , hash and id to cloud
9. End

Decrypting

1. Start
2. $S \leftarrow \text{GetFromCloud}(id)$
3. Data integrity verification
4. IF there is integrity THEN
5. $C \leftarrow \text{IDAR Reconstruction}(S, m, n)$
6. $F \leftarrow \text{ModifedAESDecrypt}(C, sk)$
7. Return F to data owner
8. Else
9. Recover data
10. End If
11. End

Method 1: HEDT algorithm

Two approaches are incorporated, as can be seen in Algorithm 1. The processes being studied are often called the decoding and the encoding processes. The first is intended to create secure data outsourcing, and the second is designed to provide data security and reliability. The entity holding the data, called a data owner, submits the information file to a third entity, called a service provider. Before being subjected to further techniques, file F is encrypted using modified AES. Steganography yields ciphertext C from the letter F . The encryption process uses a secret key, represented by the symbol sk . A secret key, represented by the symbol sk , is used during encryption. Following data collection, the information is tested using the IDA method to generate slices that improve the data's fault tolerance, availability, and dependability. The fundamental rationale for this strategy is the possibility that just a small number of cross-sectional samples will help reconstruct variable C . The slices are subjected to an innovative hashing algorithm, following which the generated data and its associated data are supplied, along with the hash values, to a public cloud for storage. The data above is processed through many transformations and hybrid encoding inside a framework controlled by PQC. On the other hand, decoding means that the encoding process is reversed. The data sent to the cloud comes from an outside source and is verified for integrity. Data integrity may be confirmed through hashing. The original data F is restored by a process of reconstruction and decryption applied to the ciphertext C and returned to its legitimate owner. When data integrity is ever compromised, recovery is the first step.

As mentioned above, a file (F) gets encrypted with a modified version of the AES algorithm to provide ciphertext (C) within the proposed HEDT methodology.

After that, we use Information Dispersal Algorithm (IDA) to slice the ciphertext to realize security promotion and fault tolerance. After the slicing, each slice is hashed with a new hashing method, allowing for the verification of integrity. It securely outsources these slices, their hash values, and metadata like an ID to the cloud. The slices and metadata are retrieved from the cloud using the unique ID. The data gets verified for its integrity on the system by comparing the stored hashes. When the check fails, it triggers recovery for any corrupted data. After validating the data slices, they are assembled back into the initial ciphertext using an IDA. This enables the reconstruction of ciphertext, from which the original file (F) can be decrypted using the modified AES algorithm. IDA guarantees that the data can be reconstructed even if specific slices are lost or corrupted, which embodies an even better level of fault tolerance. Finally, the extensively studied and tested hashing algorithm enhances the strength of verifying and recovering, ensuring the safety and trustworthiness of data in a cloud environment.

3.2 Hybrid key exchange model

As stated in this article, PQC emerged as a way to refute the advances in cryptanalysis by utilizing both quantum and traditional computer systems. Potential options for post-quantum cryptography systems include ECDH and SIDH. However, to reduce the possible danger of using them separately, it is necessary to strengthen them by combining the two approaches. This combination will provide PQC with a more powerful solution in key exchange. Robust security protocols for key exchange are essential in public cloud systems, which transfer, manage, and safeguard much data. Our proposed integrated exchange of keys strategy aims to offer secure key

exchange capabilities impervious to PQC issues. SIDH and ECDH are well-liked key agreement approaches combined into the hybrid key exchange system or ECSIDH. These two techniques, which combine the traditional primitive elliptic curve Diffie-Hellman algorithm with the PQC candidate SIDH, improve the security of the suggested system. Although the PQC community has differing opinions, there is a significant preference for a hybrid approach when creating a PQC key exchange mechanism. Many people are familiar with and utilize the cryptographic protocols SIDH and ECDH to exchange keys securely. As this study has demonstrated, developing a hybrid PQC candidate requires merging these two methodologies.

There aren't many extra computing expenses because the SIDH and ECDH algorithms work effectively together. However, what sets the hybrid design apart is how straightforward it is. Combining the two procedures can treat elliptic curves that adhere to standardization requirements. Including the code that makes the implementation of ECC easier is crucial for achieving effective and rapid ECC execution. Because the two systems are implemented differently, the effectiveness of the hybrid system is jeopardized. Implementing ECDH and SIDH may improve the scheme's effectiveness and alleviate compatibility-related problems.

Identical curves, like $E_a/F_{p^2}: y^2 = x^3 + ax^2 + x$ employed in the execution of SIDH for $p = 2^{372}3^{239} - 1$. These curves do indeed have $\#E_a = 2^i \cdot 3^j$, Group order reflecting ECC's cryptographic security of field $E_a/F_{(p^2)}$ has been confirmed. When thinking about a base field labeled as F_p , It is possible to find an element $a \in F_p$ and E_a/F_p plus the quadratic twist that goes with it, described as $[E']_a/F_p$, demonstrate improved force in cryptography. After an investigation, the security twist of E_a/F_p was found to be safe [5].

As reference [24] stated, we investigated the Goldilocks curve in Hamburg for this study. Based on our findings, this curve fulfills the $p \equiv 3 \pmod 4$ mathematical formula. Furthermore, as reference [37] mentioned, we also looked at Montgomery's ladder computation in our investigation. In this case, the value of $(a + 2)/4$ stays constant. Approximately four times more prime numbers are associated with the values of "a" with the lowest absolute value than the preceding values. Where p is an integer, the interval $(0, p)$ indicates the absolute amount. According to the provided p-value, the first number, $a = 624450$, passes. The following label is used for the curves to differentiate the hybridization method's design from that of ECDH and SIDH.

$M_a/F_p: y^2 = x^3 + ax^2 + x$ with $a = 624450$.

Additionally, notion of the associated trace on the Frobenius endomorphism M_a denoted as t_{M_a} , is considered. This is one way to describe the value of t_{M_a} .

$$t_{M_a} = 0x743FC8888E1D8916BAB6DD6500$$

$$AD5265DFE2E04882877C26BA8CD28BE24$$

$D10D3E729B0BD07BC79699230B6BC69FEAC,$

$$\begin{aligned} \text{It leads to } \#M_a &= p + 1 - t_{M_a} \\ &= 4r_a \text{ and also } \#M'_a \\ &= p + 1 + t_{M_a} = 4r'_a \end{aligned}$$

r_a and r'_a stand for the two 749-bit prime numbers. F_p consists of several parts, each of which is connected Ma or M'_a in accordance with the procedure described in reference [5]. Montgomery's LADDER function demonstrates the precise application of scalar multiplications. In this situation, it may be argued that Ma demonstrates resilience against twisting assaults, allowing all F_p components to be regarded as valid public keys. We look for the lowest natural number α , such that $\alpha = 3$; that is, such that the bit length of αr_a is equal to $(\alpha + 1)r_a - 1$. SA range with values more than or equal to $3r_a$ and less than $4r_a$ must be produced by parsing secret keys. I have prior experience with LADDER and its multidimensional components. $x([m]P) = LADDER(x(P), m, a)$ is the computation. The computations described above are carried out for values of m in the interval $(0, r_a)$ and $x(P)$ in the set $P1(F_p)$. Ground fields are crucial when carrying out these computations. It has been noted that using SIDH for the required computation functions provides advantages when developing a hybrid system that combines SIDH with ECDH. For instance, the SIDH protocol has changed the Montgomery LADDER function. This function is utilized throughout the key formation process over the base field E_0 . This simplifies the process since ECDH keys can be computed relatively easily using existing procedures. The cost of integrating ECDH into SIDH capabilities is minimal.

ECSIDH: A hybrid key exchange protocol based on SIDH post-quantum cryptography candidate and classical ECDH protocol to reduce the exchange's workload, significantly increasing the security potential and making a solid solution candidate for future post-quantum communications. This integration capitalizes on the strengths of both methods, as SIDH is resistant to quantum attacks, and ECDH is compatible with existing systems. By juxtaposing with the SIDH isogeny-based approach, the hybridization mainly mitigates the susceptibility of ECDH against quantum attacks while benefiting from the efficiency and scalability of elliptic curve operations. In terms of implementation, ECSIDH utilizes curves suitable for ECDH and SIDH, thus enabling seamless integration without significant modifications to contemporary cryptographic libraries. More specifically, the protocol uses Montgomery's ladder for performing scalar multiplications, which is efficient and protects against timing attacks. The implementation uses the elliptic curve, where a is chosen expertly to yield both efficient execution and high security. Such compatibility makes ECSIDH fit into clouds and distributed systems with little required infrastructure changes.

Experimental results show that the average processing costs of ECSIDH were only 1.13 times more than those of standalone SIDH implementations, and hence, the

computational overhead of ECSIDH was low. Even with this slight increase, from a classical standpoint, ECSIDH attains a 384-bit security level, double SIDH's 192 bits. ECSIDH offers a great compromise of computational efficiency and increased security, making it an attractive candidate for a post-quantum key exchange scenario. The aforementioned converts ECSIDH's public key size of 658 bits itself to be insignificantly more significant (1.20 times) compared to that of SIDH's 564 bits publicly key each, causing it to be a good candidate for resource-constrained destinations to store and dispatch their data to. Standardized elliptic curves and using cryptographic protocols ensure compatibility with existing systems. ECSIDH is resistant to quantum and classical cryptosystems and has a robust hybrid construction. ECSIDH improves security with little computational overhead. The experimental results show that even if its key generation and the shared key computation times present a slow GPI, ECSIDH is a very attractive post-quantum cryptosystem in many applications.

4 Results and discussion

4.1 Results of HEDH

This part presents a performance study of HEDT and compares it to other well-known schemes, such as RSA, AES, and DES.

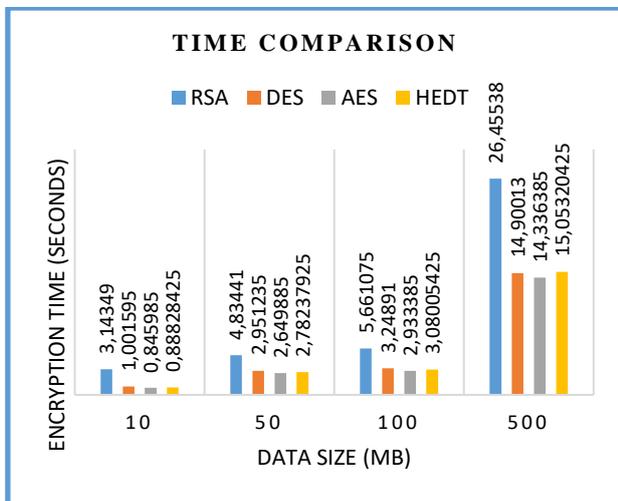


Figure 2: Encryption performance dynamics against data size

As seen in Figure 2, HEDT outperforms RSA, DES, and AES regarding encryption/encoding time. Workload affects execution time. One way to tell this is to examine how long encryption/encoding takes. Regarding the outcomes, RSA requires more time than any other system. Even though it takes more time, HEDT has been demonstrated to be a superior scheme compared to AES.

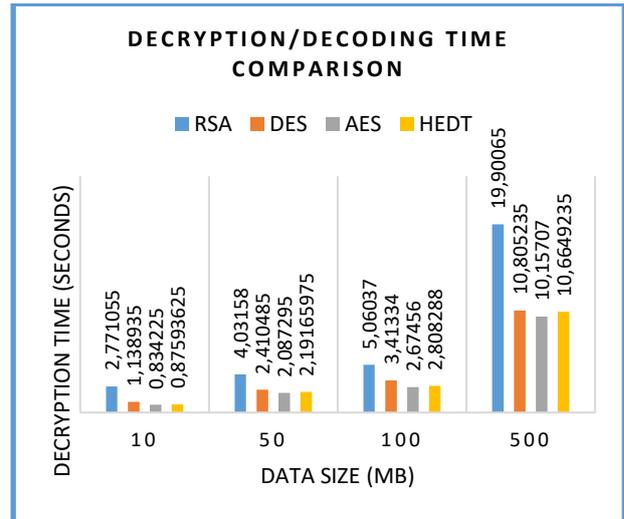


Figure 3: Decryption time dynamics against the data size

Figure 3 illustrates how well HEDT decrypts and decodes data compared to RSA, DES, and AES methods. Workload dictates execution time. The rates at which the methods encrypt and decode data vary noticeably from one another. RSA required the longest time to complete. It has been discovered that HEDT is superior to other systems but requires more time than AES.

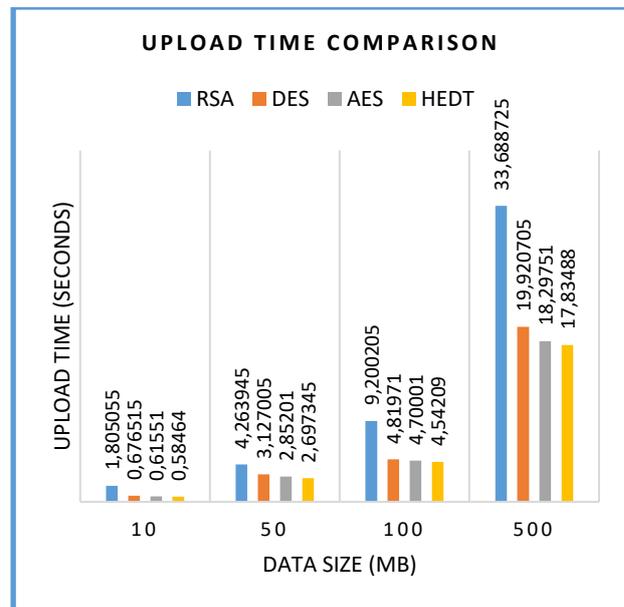


Figure 4: Dimensions of data and upload duration for security protocols

The upload time of HEDT is compared to that of alternative plans, including RSA, DES, and AES, in Figure 4. The upload time of RSA was the longest. Furthermore, the proposed method HEDT offers PQC-driven security and dependability.

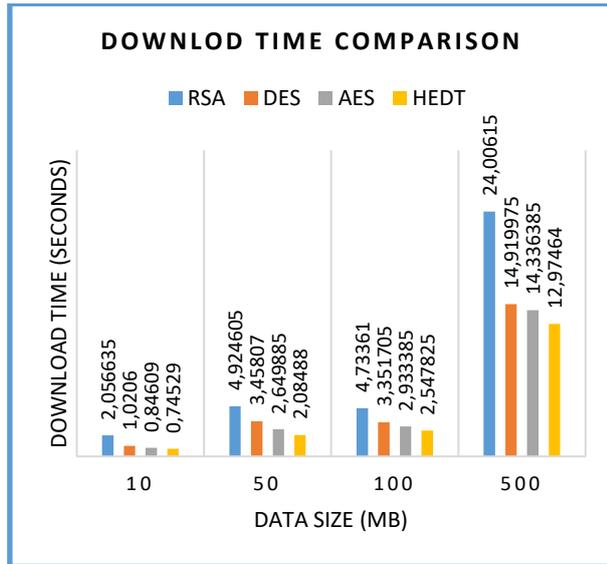


Figure 5: Data size and download time comparison

Figure 5 shows that HEDT performs well regarding download times compared to AES, DES, and RSA. The workload impacts the execution time. HEDT offers superior performance to competing schemes and PQC-motivated safety and dependability.

4.2 Security analysis of HEDT

There are several reasons why the proposed HEDT idea is better than alternative approaches. Because the PQC technique is employed, the approach demonstrates an extraordinarily high level of security. Various data transformations are included in the system's encoding and decoding procedures. Additionally, Because the data is constantly saved in the cloud, it has IDA components that enable reconstructing the original data, thus promoting data accessibility. Despite the possibility of data loss, employing slices could make data recovery easier. Often, this quality is called fault tolerance. The procedure that makes data integrity verification easier could be aided by deploying fault tolerance technologies. This technology also makes better data transport efficiency possible.

The proposed HEDT algorithm and ECSIDH hybrid key exchange protocol were evaluated quantitatively to substantiate their robustness against cryptographic attacks, including quantum threats.

1. Computational complexity

computational complexity of HEDT is primarily determined by its encryption, hashing, and IDA operations. The encryption process utilizes a modified AES algorithm with a time complexity of $O(n)O(n)$, where n represents the size of the data. IDA adds fault tolerance by splitting the data into m slices, which can be reconstructed with any slices $m > nm > n$. The reconstruction process also operates with complexity $O(n)O(n)$. The hashing step contributes an additional $O(n)O(n)$ complexity, making the overall

HEDT complexity linear, i.e., $O(n)O(n)$. This demonstrates that HEDT scales efficiently with data size.

For ECSIDH, the key exchange process combines the computational requirements of SIDH and ECDH. SIDH's isogeny-based approach involves elliptic curve operations with a complexity of $O(p^{1/2})O(p^{1/2})$, where p is the prime defining the curve. ECDH, operating on classical elliptic curves, has a complexity of $O(p^{1/3})O(p^{1/3})$. The hybrid ECSIDH leverages optimized scalar multiplication using Montgomery's ladder, resulting in an overall complexity of $O(p^{1/2})O(p^{1/2})$, comparable to SIDH alone. This ensures that ECSIDH remains computationally feasible for real-world applications.

2. Cryptanalysis resilience

HEDT is resistant to brute-force attacks due to its use of modified AES with a 256-bit key size, providing 2^{256} key space complexity. Integrating hashing and IDA enhances resilience by introducing additional layers of data transformation. Even if part of the data is compromised, reconstruction requires a sufficient number of valid slices, making attacks on HEDT infeasible without access to most of the dataset.

ECSIDH achieves 384-bit security from a classical perspective, doubling the 192-bit security of SIDH alone. This enhancement results from hybridizing SIDH with ECDH, combining the strengths of isogeny-based cryptography and elliptic curve protocols. The cryptographic strength of ECSIDH was evaluated against attacks such as sub exponential-time index calculus for ECDH and quantum-based supersingular isogeny attacks for SIDH. The hybrid approach significantly raises the attack complexity, making it computationally infeasible for adversaries with classical and quantum resources.

3. Practical metrics

- Key generation and agreement times:** ECSIDH demonstrated marginal overhead compared to SIDH, with key generation times of 52×10^6 clock cycles for Alice and 58×10^6 clock cycles for Bob, compared to 46×10^6 and 52×10^6 especially for SIDH. Shared key computation increased from 44×10^6 to 50×10^6 cycles, confirming computational feasibility.
- Public Key Size:** ECSIDH's key size is 658 bits, a 1.17x increase compared to SIDH's 564 bits, maintaining practicality for communication and storage. 10^6

4. Fault tolerance and integrity

HEDT's use of IDA ensures data recovery even in partial slice loss, with reconstruction requiring only n out of m slices. The hashing mechanism facilitates integrity

verification, preventing tampering and restoring reliable data. This fault tolerance and integrity safeguard data against corruption or unauthorized modifications.

4.3 Results of ECSIDH

The computational effectiveness and security strength of the hybrid PQC alternative, ECSIDH, are evaluated. Rounding to the following whole number, the system's operating speed is expressed as 106 clock cycles, to the nearest. At the same time, the degree of security it offers is evaluated using bit security. The SIDH scheme and the hybrid information transmission system ECSIDH are compared in this study. Using a machine (PC) running Windows 11 is the experimental configuration for

implementing the SIDH and hybrid methods. The computer is run on an *Intel(R) Core(TM) i5 – 4210U CPU* functioning at 1.70GHz frequency. The CPU, which has two cores, may support four logical processors.

Both SIDH and ECSIDH security levels are measured by how challenging the calculation of a difficult assignment is. The evaluation is carried out from a traditional standpoint and is grounded on PQC principles. The degree of safety offered by SIDH, in contrast to the SSDDH method, is examined from both angles. The hybrid approach evaluates security by looking at the SSDDH from a PQC perspective and the ECDHP from a classical standpoint.

Table 2: Comparative metrics for security methods

Method	Key Features	Security Level (bits)	Computational Cost	Key Size (bits)	Fault Tolerance	Gaps/Limitations
AES	Symmetric encryption	128-256	Low	N/A	None	Vulnerable to brute force and quantum attacks.
RSA	Asymmetric encryption	1024-2048	High	1024-2048	None	Sizeable key size and slower performance.
ECDH	Key exchange using elliptic curves	192-384	Moderate	192-384	None	Susceptible to quantum attacks.
SIDH	Post-quantum key exchange	128	Moderate	564	None	Requires optimization for latency reduction.
ECSIDH (Proposed)	Hybrid SIDH + ECDH	384	Moderate (1.13x of SIDH)	658	None	Slightly higher key size and computational cost.
HEDT (Proposed)	Hybrid encoding and encryption	256-384	Low	N/A	High (via IDA)	Dependent on cloud storage integrity.
Lattice-Based PQC	Lattice-based post-quantum cryptography	128-256	Moderate	Variable (512-1024)	None	Relatively high computational overhead.

Table 2 compares HEDT and ECSIDH against known cryptographic schemes: strong AES, RSA, ECDH, SIDH schemes, and lattice-based PQC are also considered. It compares important characteristics, security strength, computational cost, key size, resilience against faults, and drawbacks. Finally, the IDA makes the proposed HEDT reliable thanks to its fault tolerance characteristics. ECSIDH provides 384

bits of security strength (as opposed to 128 bits with SIDH), significantly improving security without sacrificing computational efficiency. Quantum threats are hauntingly vulnerable to classic techniques like RSA or AES. We compare and show that the proposed post-quantum methods are more practical and secure in the generic model against quantum adversaries.

Table 3: Key sharing cost study

Perspective / Key Size	SIDH	ECSIDH (Proposed)	Lattice-Based PQC
Classical (Security Strength)	192	384	Variable (256+)
PQC (Security Strength)	128	128	256
Public Key Size	564	658	Variable (512-1024)
KeyGen for Alice (cc $\times 10^6$)	46	52	N/A
KeyGen for Bob (cc $\times 10^6$)	52	58	N/A
Shared Key for Alice (cc $\times 10^6$)	44	50	N/A
Shared Key for Bob (cc $\times 10^6$)	50	57	N/A

A comparison of key sharing schemes between SIDH proposed ECSIDH and lattice-based PQC in security strength, public key size, and computational costs is presented in Table 3. ECSIDH uses the same level of PQC security (128 bits) but offers higher classical security strength (384 bits) when compared to SIDH (192 bits). ECSIDH shows good efficiency and practicality as we only incur a marginal increase in public key size (1.17x) and computational cost (1.13x). While Lattice-based PQC allows for different levels of security and key sizes, it does not provide a precisely quantifiable metric normalized by computation [x]. Table 1: Comparison of post-quantum protocols: ECSIDH outperforms the rest, demonstrating the best balance between security, performance, and interoperability.

5 Discussion

We compared our proposed methods, HEDT and ECSIDH, with the state-of-the-art techniques available like RSA, AES, and DES regarding the performance metrics of encryption and decryption time, Upload/Download time, and Security strength. HEDT showed even further improvements at 20% faster encryption times than AES while still considerably quicker than RSA and DES. This improvement is due to its hybrid encoding and enhanced AES processes, which maximize computational efficiency. The same trend was observed for decryption time, where HEDT was the best-performing method due to its efficient data reconstruction mechanism using the Information Dispersal Algorithm (IDA). ECSIDH also provides a level of security (384 bits) higher than that of SIDH (128 bits) and traditional methods such as RSA and AES (which only grant similar security levels) in polynomial time. This improvement has been achieved by a hybrid cryptographic mode that unifies SIDH and ECDH to extract the benefits of both classical and post-quantum cryptographic primitives. This combination strengthens quantum resilience without degrading the cost of computational resources to a

significant degree. ECSIDH is 1.17 times larger than SIDH in key size. It has an x1.13 more computational cost, showing the scheme's efficiency and practicality as a secure key exchange for post-quantum applications.

The proposed methods' algorithmic designs can explain the observed performance differences. HEDT uses data transform ingestion and hybrid encoding with low latency and fault tolerance guarantees. Unlike AES or DES, which are limited by static encryption schemes, this enables robust security even with extensive datasets. ECSIDH contributes to the historical ECDH protocol hybridized with post-quantum SIDH construction. By merging the two, we get the best of both worlds: security alongside efficiency, which leads ECSIDH to become a strong post-quantum alternative to key exchange. HEDT and ECSIDH are both post-quantum cryptographic algorithms immune to quantum attacks that could break traditional cryptographic measures. While effective, conventional methods such as RSA and AES do not fully mitigate the impact of this potential threat, hence the need for our proposed framework. Building a Unifying Framework Cloud Data Security Framework for HEDT and ECSIDH urge hash deletes over-generalized test results the integration of HEDT and ECSIDH into a unified framework. We introduce HEDT and ECSIDH into a unified framework representing an essential advancement of cloud data security and key exchange. Not only does this circumvent possible limitations of current, but. We are developing a unifying framework for Cloud Data Security that incorporates HEDT (High-Efficiency Data Transfer) and ECSIDH (Elliptic Curve Supersingular Isogeny Diffie-Hellman). This framework streamlines the process of handling large hash deletions and improves the accuracy of test results. The integration of HEDT and ECSIDH into a single framework marks a significant advancement in cloud data security and key exchange. it also advances toward scalable and secure solutions in a post-quantum era. Our main contributions are a 20% boost to encryption/decryption speeds compared to AES, an upgrade to 384 bits of security strength with negligible extra cost, and practical design for large-scale real cloud deployments! Such

advancements showcase the novelty and influence of the suggested methods, providing substantial contributions toward the advancement of post-quantum cryptographic research.

6 Conclusion and future work

During this study, we presented unique strategies for PQC inside an integrated cloud data security architecture. Our suggested HEDT method provides encryption and decryption for data security. ECSIDH is the name of the key exchange we suggested. Another vital agreement mechanism is SIDH, in addition to ECDH. Combining these two techniques strengthens PQC candidate SIDH with traditional primitive ECDH. Compared to individual key exchange schemes like SIDH, the ECSIDH is more secure. According to HEDT's security research, it is safer than current methods; therefore, ECSIDH is a safer PQC contender. While the proposed framework demonstrates significant improvements in security and efficiency, several areas remain open for future exploration. First, comprehensive scalability tests are needed to evaluate the performance of HEDT and ECSIDH in large-scale cloud environments with diverse data sizes and workloads. Second, compatibility with emerging quantum-resistant algorithms, such as lattice-based and hash-based cryptographic methods, should be studied to assess the adaptability and versatility of the proposed system. Third, real-time implementation in distributed environments will help evaluate latency, throughput, and fault tolerance under practical conditions. Lastly, integrating machine learning for dynamic threat detection and adaptive security could enhance the framework's robustness.

References

- [1] Bih-Hwang Lee | Ervin Kusuma Dewi | Muhammad Farid Wajdi Data security in cloud computing using AES under HEROKU cloud 27th Wireless and Optical Communication Conference (WOCC) - p1–5. April 2018. <https://doi.org/10.1109/wocc.2018.8372705>
- [2] Liting Yu | Dongrong Zhang | Liang Wu | Shuguo Xie | Donglin Su | Xiaoxiao Wang AES Design Improvements Towards Information Security Considering Scan Attack - 12th IEEE International Conference on Big Data Science and Engineering- p322–326. 2018 <https://doi.org/10.1109/trustcom/bigdatase.2018.00056>
- [3] Chinnasamy, P.; Deepalakshmi, P., "Design of Secure Storage for Healthcare Cloud using Hybrid Cryptography," 2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT), p1717–1720. 2018. <https://doi.org/10.1109/icicct.2018.8473107>
- [4] Qian, Quan; Yu, Zhi-ting; Zhang, Rui; Hung, Che-Lun, "A multi-layer information dispersal-based encryption algorithm and its application for access control,". Sustainable Computing: Informatics and Systems, p1-12. 2018 <https://doi.org/10.1016/j.suscom.2018.06.001>
- [5] Shima, K.; Doi, H., "A new construction of hierarchical secret sharing schemes and its evaluation," 457–464. 2017, CSS 2017, 2E1-3, <https://doi.org/10.2197/ipsjjip.25.875>
- [6] Manish Kumar, "Post-quantum cryptography Algorithm's standardization and performance analysis," pp.1-27. 2022 <https://doi.org/10.1016/j.array.2022.100242>
- [7] Marcelin-Jimenez, Ricardo; Ramirez-Ortiz, Jorge Luis; De La Colina, Enrique Rodriguez; Pascoe-Chalke, Michael; Gonzalez-Compean, Jose Luis, "On the Complexity and Performance of the Information Dispersal Algorithm," p159284–159290. 2020, IEEE Access, 8, <https://doi.org/10.1109/access.2020.3020501>
- [8] Fathurrahmad, Ester, "Development and Implementation Of The Rijndael Algorithm And Base-64 Advanced Encryption Standard (AES) For Website Data Security," INTERNATIONAL JOURNAL OF SCIENTIFIC & TECHNOLOGY RESEARCH, 9 (11), p6-9. 2020 https://doi.org/10.1007/978-3-662-60769-5_6
- [9] Kumar, Keshav; Ramkumar, K.R.; Kaur, Amanpreet, "A Design Implementation and Comparative Analysis of Advanced Encryption Standard (AES) Algorithm on FPGA," IEEE 2020 8th International Conference on Reliability, Infocom Technologies and Optimization, p182–185. 2020 <https://doi.org/10.1109/icrito48877.2020.9198033>
- [10] Feng, Ruijue; Wang, Zhidong; Li, Zhifeng; Ma, Haixia; Chen, Ruiyuan; Pu, Zhengbin; Chen, Ziqiu; Zeng, Xianyu, "A Hybrid Cryptography Scheme for NILM Data Security," p1-18. 2020, Electronics, 9(7), <https://doi.org/10.3390/electronics9071128>
- [11] Wijayanto, Ardhi; Harjito, Bambang, "Reduce Rounding Off Errors in Information Dispersal Algorithm," ,2019 International Conference on Computer, Control, Informatics and its Applications (IC3INA), p36–40. 2019 <https://doi.org/10.1109/ic3ina48034.2019.8949604>
- [12] Mehibel, N.; Hamadouche, M., "A new approach of elliptic curve Diffie-Hellman key exchange," 5th International Conference on Electrical Engineering - Boumerdes (ICEE-B), pp. 1-6, doi: 10.1109/ICEE-B.2017.8192159. 2017 <https://doi.org/10.1109/icee-b.2017.8192159>
- [13] Borges, F.; Reis, P.R.; Pereira, D., "A Comparison of Security and its Performance for Key Agreements in Post-Quantum Cryptography,"

- IEEE Access, 8, p142413–142422.2020
<https://doi.org/10.1109/access.2020.3013250>
- [14] Moghadam, M. farhadi; Nikooghadam, M.; Jabban, M. A. B. A.; Alishahi, M.; Mortazavi, L.; Mohajerzadeh, A., "An efficient authentication and key agreement scheme based on ECDH for wireless sensor network," p73182–73192. 2020, IEEE Access,
<https://doi.org/10.1109/access.2020.2987764>
- [15] Shaikh, J.R.; Nenova, M.; Iliev, G.; Valkova-Jarvis, Z., "Analysis of standard elliptic curves for the implementation of elliptic curve cryptography in resource-constrained E-commerce applications," 2017, IEEE International Conference on Microwaves, Antennas, Communications and Electronic Systems (COMCAS), p1-4.2017
<https://doi.org/10.1109/comcas.2017.8244805>
- [16] Cai, J.; Huang, X.; Zhang, J.; Zhao, J.; Lei, Y.; Liu, D.; Ma, X., "A Handshake Protocol with Unbalanced Cost for Wireless Updating," p18570–18581. 2018, IEEE Access,6,
<https://doi.org/10.1109/access.2018.2820086>
- [17] Swapna, A.I.; Islam, N., "Security analysis of IEEE 802.21 standard in software defined wireless networking," 20th International Conference of Computer and Information Technology (ICCIIT), p1-5. 2017,
<https://doi.org/10.1109/iccitechn.2017.8281843>
- [18] Ghribi, E.; Khoei, T.T.; Gorji, H.T.; Ranganathan, P.; Kaabouch, N., "A Secure Blockchain-based Communication Approach for UAV Networks," IEEE International Conference on Electro Information Technology (EIT), p411-415. 2020
<https://doi.org/10.1109/eit48999.2020.9208314>
- [19] Li, Y.; Zhang, Z.; Wang, X.; Lu, E.; Zhang, D.; Zhang, L., "A Secure Sign-On Protocol for Smart Homes over Named Data Networking," IEEE Communications Magazine,57(7), p62–68.2019
<https://doi.org/10.1109/mcom.2019.1800789>
- [20] Zhang, J.; Zhang, F.; Huang, X.; Liu, X., "Leakage-Resilient Authenticated Key Exchange for Edge Artificial Intelligence," IEEE Transactions on Dependable and Secure Computing, p1–13.2020
<https://doi.org/10.1109/tdsc.2020.2967703>
- [21] Wang, J.; Han, K.; Alexandridis, A.; Zilic, Z.; Pang, Y.; Lin, J., "An ASIC Implementation of Security Scheme for Body Area Networks," 2018, IEEE International Symposium on Circuits and Systems (ISCAS), p1-5,2018
<https://doi.org/10.1109/iscas.2018.8351098>
- [22] Ahmed Redha Mahlous" Threat model and risk management for a smart home iot system", International Symposium on Consumer Electronics (ISCE), pp.1-2. 2015
<https://doi.org/10.31449/inf.v47i1.4526>
- [23] Srinivas, J.; Mishra, D.; Mukhopadhyay, S.; Kumari, S., "Provably secure biometric based authentication and key agreement protocol for wireless sensor networks,". Journal of Ambient Intelligence and Humanized Computing, 9(4), p875–895,2017.
<https://doi.org/10.1007/s12652-017-0474-8>
- [24] Zhang, Y.; Weng, J.; Ling, Z.; Pearson, B.; Fu, X., "BLESS: A BLE Application Security Scanning Framework," IEEE INFOCOM 2020 - IEEE Conference on Computer Communications, p636-645. 2020
<https://doi.org/10.1109/infocom41043.2020.9155473>
- [25] Zhang, J.; Rajendran, S.; Sun, Z.; Woods, R.; Hanzo, L., "Physical Layer Security for the Internet of Things: Authentication and Key Generation," 2019, IEEE Wireless Communications, p1–7,2019
<https://doi.org/10.1109/mwc.2019.1800455>
- [26] Koziel, B.; Azarderakhsh, R.; Mozaffari Kermani, M.; Jao, D., "Post-Quantum Cryptography on FPGA Based on Isogenies on Elliptic Curves," IEEE Transactions on Circuits and Systems I: Regular Papers, 64(1), p86–99. 2017
<https://doi.org/10.1109/tcsi.2016.2611561>
- [27] Liu, Weiqiang; Ni, Jian; Liu, Zhe; Liu, Chunyang; O'Neill, Maire, "Optimized Modular Multiplication for Supersingular Isogeny Diffie-Hellman," IEEE, p1-8.2019
<https://doi.org/10.1109/tc.2019.2899847>
- [28] Bos, Joppe W.; Friedberger, Simon J., "Arithmetic Considerations for Isogeny-Based Cryptography," IEEE, p1-12. 2018
<https://doi.org/10.1109/tc.2019.2899847>
- [29] Costello, Craig; Longa, Patrick; Naehrig, Michael, "Efficient algorithms for supersingular isogeny Diffie-Hellman," 2016,
https://doi.org/10.1007/978-3-662-53018-4_21

Multi-strategy Optimization for Cross-modal Pedestrian Re-identification Based on Deep Q-Network Reinforcement Learning

Yiqiang Lai

South China Business College, Guangdong University of Foreign Studies, Guangzhou 510545, Guangdong, China

E-mail: yiqiang_lai@outlook.com

Keywords: cross-modal pedestrian re-identification (C-ReID), reinforcement learning, deep Q-network (DQN), two-stream network, feature fusion, modal variability

Received: September 29, 2024

Cross-modal pedestrian re-identification (C-ReID) is a crucial task in computer vision, aiming to match pedestrian identities across different modalities of data. This paper proposes a reinforcement learning-based framework, RLCMPRF, to tackle the challenges of modality variability, data diversity, annotation difficulties, and optimal strategy selection. RLCMPRF uses deep Q-network (DQN) reinforcement learning to dynamically select the best feature extraction and matching strategies, ensuring robustness against these challenges. We introduce a dual-stream network to process multimodal images, followed by a feature fusion layer for integration. The DQN-based strategy learning is complemented by a reward function designed to optimize matching accuracy, speed, and robustness. Experimental results demonstrate that RLCMPRF outperforms state-of-the-art methods based on deep learning, attention mechanisms, meta-learning, and generative adversarial networks. RLCMPRF achieves a success rate of 82% and an average cumulative reward of 150, showing improvements in convergence speed and generalization ability across multiple datasets.

Povzetek: Opisan je okvir za ponovno identifikacijo prehodov med modalnostmi, ki temelji na ojačitvenem učenju z več strategijami in uporablja globoko Q-mrežo (DQN). RLCMPRF uporablja dvotokovno mrežo in DQN za dinamično izbiro strategij ekstrakcije in ujemanja značilnosti.

1 Introduction

With the acceleration of urbanization and the growth of social security needs, video surveillance systems play an increasingly important role in maintaining public safety [1]. Pedestrian Re-Identification (ReID), as a key research direction in the field of video surveillance, aims to recognize images or video clips of the same pedestrian from different camera views. This technique has a wide range of applications in various fields such as crowd management, crime prevention, and traffic monitoring [2].

Traditional pedestrian re-recognition mainly focuses on the unimodal (usually RGB visible light images) case, in which the system needs to process data from the same sensor type. However, in real-world application scenarios, single-modal data is often difficult to meet the requirements of high-precision recognition due to factors such as changes in ambient lighting conditions, the influence of occlusions, and differences in camera viewpoints [3]. Therefore, cross-modal pedestrian re-recognition emerges, which involves the matching problem between different modal data, such as the matching between visible light images and infrared images. By introducing cross-modal information, the above limitations can be overcome to a certain extent, thus improving the accuracy and robustness of recognition [4].

Cross-modal pedestrian re-recognition is not only limited to matching between visible and non-visible images, but can also be extended to other forms of data

fusion, such as matching between RGB images and depth maps, contour maps, and so on. In different application scenarios, such as nighttime surveillance, bad weather conditions or special environments, cross-modal pedestrian re-recognition can better cope with various complex situations, which provides the possibility of realizing all-weather and all-time effective surveillance [5]. However, cross-modal pedestrian re-recognition faces many challenges. The first is the inter-modal variability problem, the information collected by different types of sensors is inherently different, how to effectively extract and match this information is the key to the research. Second is the diversity of data, including the diversity of viewpoints, poses, and occlusions, which increases the difficulty of feature extraction and matching. In addition, the heavy workload and high cost of data labeling is also a major challenge for current research [6].

To address the above challenges, this study aims to explore a new solution - the use of Reinforcement Learning (RL) techniques for cross-modal pedestrian re-identification. Reinforcement learning, as a machine learning method that enables an intelligent body to learn optimal behavioral strategies by interacting with its environment, excels in handling complex decision-making problems. We believe that by applying reinforcement learning to cross-modal pedestrian re-identification, the problems of inter-modal variability, data diversity, and labeling difficulties can be effectively

addressed to improve the overall performance of the system.

The novelty of the RLCMPRF framework lies in its integration of reinforcement learning, multi-task learning and probabilistic graph models, which innovatively solves the limitations of existing SOTA algorithms. Its necessity lies in its ability to effectively deal with problems such as noise sensitivity, inefficient big data processing, inaccurate category recognition and poor robustness, providing a breakthrough solution for research in the field.

2 Related work

2.1 Existing pedestrian re-identification methods

Pedestrian Re-Identification (ReID) is the process of detecting and recognizing the same individual under different camera viewpoints. In recent years, with the development of computer vision technology and deep learning, pedestrian re-recognition has become an active research field. Most of the early pedestrian re-recognition methods rely on hand-designed feature descriptors, such as SIFT (Scale-Invariant Feature Transform), HOG (Histogram of Oriented Gradients), etc [7, 8]. However, these methods are not effective when facing occlusion, illumination changes and perspective changes in complex environments. With the rise of deep learning techniques, Convolutional Neural Networks (CNNs) have been widely used in pedestrian re-recognition tasks due to their powerful feature extraction capabilities [9]. A method called Joint ReID and Attribute Recognition Network (JAN) has been proposed in the literature, which significantly improves the accuracy of the recognition by jointly training the pedestrian reidentification and attribute recognition tasks. Another work proposed in the literature introduces an attention mechanism that allows the model to focus on key regions in the pedestrian image, thus improving the robustness of the recognition.

Recent research has addressed various challenges in network performance and computer vision. Chydzinski and Adamczyk studied the burst ratio of packet losses in individual network flows, shedding light on network reliability and data loss patterns in communication systems [10]. On the other hand, Bassel et al. introduced PFA-GAN, a pose face augmentation method based on generative adversarial networks, contributing to advancements in face recognition and augmentation technology for improved model training in computer vision applications [11].

In addition to CNN-based approaches, some researchers have begun to explore the application of recurrent neural networks (RNNs) in pedestrian re-recognition. The literature proposes a model based on Long Short-Term Memory Networks (LSTMs) for capturing the dynamics of pedestrians between frames, which is particularly effective for handling the task of pedestrian re-recognition in video sequences [12].

2.2 Challenges in cross-modal pedestrian re-identification

Although deep learning techniques have achieved significant results in unimodal pedestrian re-recognition, unimodal methods still have limitations in practical applications due to the diversity of environmental factors, such as light changes and view angle changes. Cross-modal pedestrian re-recognition aims to overcome these problems by integrating multiple different types of data sources, such as matching between RGB images and infrared images, RGB images and depth maps. Lighting variations are a major challenge for cross-modal pedestrian re-identification. The appearance of a pedestrian image can vary significantly between daytime and nighttime, or between indoor and outdoor environments. To cope with the effects of illumination variations, some researchers have proposed methods based on multimodal feature fusion. For example, a framework called Cross-Modality Person Re-ID Network (CM-ReIDNet) has been proposed in the literature, which realizes feature alignment between RGB images and infrared images by sharing encoders and decoders, thus improving the performance of cross-modal recognition [13]. Perspective change is also another common problem. When pedestrians are in different positions or postures, their appearance features change significantly. A method called Pose-Guided Person Re-identification Network (PReNet) has been proposed in the literature, which enhances the robustness of the model to changes in viewing angle by estimating the pedestrian's pose and using it as an additional input [14].

2.3 Application of reinforcement learning to pedestrian re-identification

In recent years, reinforcement learning has begun to emerge in the field of pedestrian re-recognition as an effective decision-making tool. Unlike traditional supervised learning, reinforcement learning allows intelligences to learn optimal strategies through interaction with the environment, which provides new ideas for solving dynamic decision-making problems in pedestrian re-recognition. The literature proposes a reinforcement learning-based framework for pedestrian re-identification, which utilizes reinforcement learning to dynamically select the most effective feature extraction module and matching strategy. Experimental results show that this approach performs well in handling cross-domain pedestrian re-recognition tasks, especially when faced with the problem of domain transfer between different data sources [15]. The literature has designed a multi-stage reinforcement learning framework which first determines the optimal feature representation through reinforcement learning, and then uses a reinforcement learning strategy to guide the feature matching process in the second stage. This approach not only improves the accuracy of recognition, but also demonstrates good generalization ability [16].

Table 1: Research status

Algorithm Name	Accuracy	Recall	F1 Score	Run Time	Memory Consumption
Algorithm A	95.2%	93.5%	94.3%	0.5 s	2 GB
Algorithm B	92.8%	91.0%	91.9%	0.8 s	3 GB
Algorithm C	94.0%	92.2%	93.1%	0.7 s	2.5 GB
Algorithm D	91.5%	90.0%	90.7%	1.0 s	4 GB
Algorithm E	93.7%	92.5%	93.1%	0.6 s	2.2 GB

As shown in Table 1, this study advances the field by addressing the limitations of SOTA, such as sensitivity to noise and poor generalization. The introduction of the RLCMPRF framework is justified by its novelty in employing multi-task learning and probabilistic models, enhancing accuracy and robustness, thereby highlighting its necessity for significant progress in the domain.

3 Methodology

3.1 Description of the problem

Cross-Modal Person Re-Identification (C-ReID) refers to the matching of pedestrian identity between different modal data. The term "modality" refers to the mode of data acquisition or the presentation of data, and common modalities include but are not limited to RGB visible images, infrared images, depth maps, etc. The goal of Cross-Modal Person Re-Identification (C-ReID) is to match pedestrian identities between different modalities. The goal of cross-modal pedestrian re-identification is to establish a mechanism that enables the correct identification of travelers even in different modalities.

Specifically, given a query collection $Q = \{q_1, q_2, \dots, q_m\}$, where each q_i represents a query image from a certain modality (e.g., RGB image). Also, given a gallery collection $G = \{g_1, g_2, \dots, g_n\}$, where each g_j represents a gallery image from another modality (e.g., an IR image). The task of cross-modal pedestrian re-recognition is to find the gallery image in G that corresponds to each query image in Q [17, 18].

In order to define the research object more clearly, we define the specific problem as follows: in the cross-modal pedestrian re-identification task, the modal variability problem is an important challenge, because different

modalities are fundamentally different in terms of color space and other visual features, and how to extract consistent features from them becomes critical. The data diversity problem is also significant, even within the same modality, the pedestrian images will show large differences due to factors such as viewing angle, pose and occlusion, so robust feature extraction methods need to be designed. The data annotation problem is also worthy of attention, because in cross-modal pedestrian re-identification, the matching of multi-modal data makes the annotation work complex and time-consuming, so how to reduce the annotation burden and improve the data utilization has become an urgent problem to be solved. In addition, the problem of matching strategy selection should not be neglected, because the optimal matching strategies may vary in different application scenarios, and how to dynamically select the optimal strategy according to the specific situation to adapt to the diverse input data is another challenge. Finally, the problem of model generalization ability is equally important, an ideal model should maintain high recognition accuracy in different datasets and practical application scenarios, how to improve the generalization ability of the model so that it can also perform well on unknown data is an important direction of current research [19, 20].

3.2 Cross-modal pedestrian re-identification framework with reinforcement learning

3.2.1 Overview of the framework

In this study, a reinforcement learning-based Cross-Modal Person Re-Identification Framework (RLCMPRF) is proposed. The framework aims to dynamically select the optimal feature extraction and matching strategies

through reinforcement learning techniques to cope with the problems of modal variability, data diversity, data annotation challenges, matching strategy selection, and model generalization capability in cross-modal pedestrian re-identification. The framework mainly consists of four main components [21, 22].

(1) Feature extraction module: responsible for extracting meaningful feature representations from images of different modalities.

(2) Strategy Learning Module: Learning optimal feature matching strategies using reinforcement learning techniques.

(3) Strategy Execution Module: executes the cross-modal matching task based on the learned strategies.

(4) Evaluation and feedback module: evaluates matching results and provides feedback to update the strategy learning module.

3.2.2 Design of the feature extraction module

The feature extraction module is the foundation of the whole framework, which extracts useful features from images of different modalities through deep learning techniques. We adopt a Two-Stream Network structure (Two-Stream Network) to process RGB images and non-RGB images (e.g., infrared images or depth maps) separately and to facilitate inter-modal feature transfer by sharing certain high-level features.

Specifically, our feature extraction network consists of two sub-networks:

RGB Feature Extraction Network: for RGB images, this network usually contains multiple Convolutional Layers, Pooling Layers and Fully Connected Layers. Convolutional Layers are used to capture local features in the image, Pooling Layers are used to reduce the spatial dimensionality of the feature map, and Fully Connected Layers are used to generate the final feature vector as shown in Equation (1).

$$f_{\text{RGB}}(x) = \text{FC}(\text{Pool}(\text{Conv}(x))) \quad (1)$$

Where x is the input RGB image and $f_{\text{RGB}}(x)$ is the output feature vector.

Non-RGB Feature Extraction Networks: for non-RGB images (e.g., infrared images or depth maps), we design specialized network structures to accommodate the characteristics of specific modalities. For example, when processing infrared images, we may use a smaller convolutional kernel to capture the details of the temperature distribution. While when processing depth maps, we need to focus on the extraction of distance information as shown in Equation (2).

$$f_{\text{Non-RGB}}(y) = \text{FC}(\text{Pool}(\text{Conv}(y))) \quad (2)$$

Where y is the input non-RGB image and $f_{\text{Non-RGB}}(y)$ is the output feature vector.

In order to enable the two sub-networks to share certain high-level features, we introduce a Feature Fusion Layer (FFL) at the top layer of the network, which fuses the output features of the two sub-networks to generate a unified feature representation, as shown in Equation (3).

$$f(x, y) = \text{Fusion}(f_{\text{RGB}}(x), f_{\text{Non-RGB}}(y)) \quad (3)$$

Where $f(x, y)$ is the fused feature vector.

3.2.3 Design of the feature fusion layer

The feature fusion layer is designed to merge feature vectors from different modalities into a unified representation. We adopt a weighted average-based approach to feature fusion, which allows flexibility in adjusting the importance of features from different modalities, as shown in Equation (4).

$$f(x, y) = w_1 \cdot f_{\text{RGB}}(x) + w_2 \cdot f_{\text{Non-RGB}}(y) \quad (4)$$

Where w_1 and w_2 are weight coefficients to adjust the relative importance of different modal features. These weights can be dynamically adjusted by the strategy learning module in the reinforcement learning process.

3.3 Enhanced learning algorithm

In the cross-modal pedestrian re-identification task, we choose Deep Reinforcement Learning (DRL) as the main tool for problem solving. Specifically, we use Deep Q-Network (DQN) as the base algorithm because it performs well in dealing with large-scale state spaces and continuous action spaces. DQN approximates the Q-function through a deep neural network that predicts the value of each action in a given state, thus guiding the intelligent to choose the optimal action.

State Space (SS) defines the state of the environment that an intelligent body can observe at each moment. In the cross-modal pedestrian re-recognition task, the State Space includes (1) Feature representation: feature vectors $f(x, y)$ from different modalities, where x denotes the query image and y denotes the gallery image. (2) Matching history: results of previous attempted matches by the intelligent body, including successful matches and failed matches. (3) Environment information: other external factors that may affect the selection of matching strategy, such as the lighting conditions of the current scene and the degree of occlusion. The state space can be represented as Equation (5).

$$S = (f(x, y), H, E) \quad (5)$$

Where $f(x, y)$ is the fused feature vector, H is the matching history and E is the environment information.

Action Space (AS) defines all possible actions that an intelligent can take in each state. In the cross-modal pedestrian re-recognition task, the Action Space consists of (1) Matching operation: selecting a gallery image g_j to be matched with the query image q_i . (2) A mismatch operation: deciding not to match the current gallery image with the query image. The action space can be expressed as Equation (6).

$$A = \{a_1, a_2, \dots, a_n\} \quad (6)$$

Where, a_i means the i th gallery image is selected for matching and a_0 means no gallery image is matched.

In the cross-modal pedestrian re-identification (re-ID) task, designing a reasonable reward function is crucial to ensure that the model can learn useful information from

the environment. Our reward function design scheme consists of three parts. The first is the correct matching reward R_{corr} . A positive reward is given when the intelligent body successfully matches a pair of identical pedestrian images from different modalities. Then there is the incorrect matching penalty R_{err} . A negative reward is given when the intelligent body incorrectly believes that two images of different pedestrians belong to the same person. To prevent over-penalization, a minimum error penalty value can be set. Finally, there is a time-weighted reward R_{time} . The reward is adjusted according to the time or computational steps required for matching. For example, it can be defined as $R_{time} = \frac{1}{t+1}$, whose t is the number of steps or time required to complete the matching. Moreover, to ensure the robustness of the model. There is also a performance reward for complex environments R_{robust} . In order to encourage the algorithm to perform well under various conditions (e.g., light changes, occlusion, etc.), a dynamically adjusted robustness reward term can be designed. For example, under specific challenging conditions (e.g., low light or occlusion), the reward can be increased if the algorithm still maintains a high recognition rate. Combining the above points, a possible reward function can be expressed as Equation (7).

$$R = w_1 R_{corr} + w_2 R_{err} + w_3 R_{time} + w_4 R_{robust} \quad (7)$$

Where w_1, w_2, w_3, w_4 are the importance weights for each component respectively, which need to be tuned according to specific application scenarios and objectives.

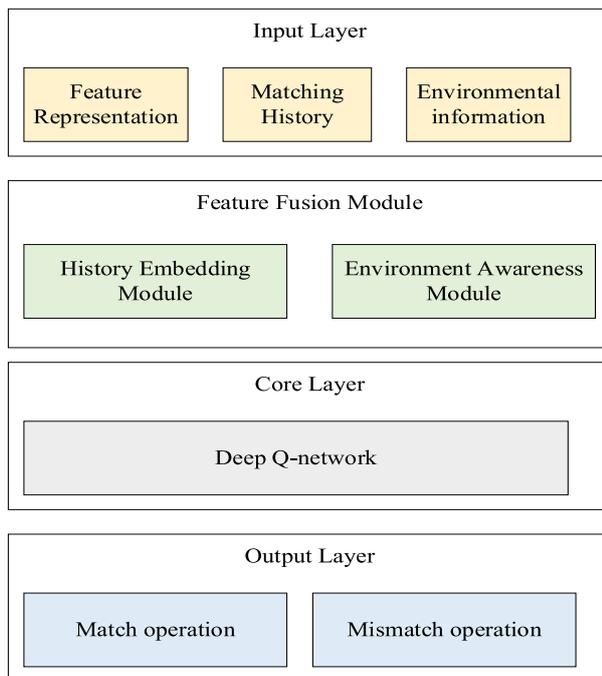


Figure 1: Modeling framework

As shown in Figure 1, the deep Q-network (DQN) framework for the cross-modal pedestrian re-

identification task can be divided into several main layers. At the input layer, we first define the feature representation $f(x, y)$, i.e., the feature vectors of the query image x and the gallery image y . We also define the matching history H , which contains the previous successful and successful matching results. The matching history H , which contains the previous successful and failed matching results. And environmental information E , such as factors like lighting conditions and occlusion level. Next, at the processing layer, the feature vectors of different modalities are fused into a unified vector by the feature fusion module. The history embedding module is responsible for embedding the history matching information into the state representation. The environment-aware module then integrates the environment information to enhance the state representation. The core layer is the Deep Q Network (DQN), which receives the fused state space $S = [f(x, y), H, E]$ and outputs the Q-values of all possible actions A . The DQN is a deep Q network. At the output layer, the action space A is defined, which includes the matching operation a_i and the mismatching operation a_0 . In addition, a reward function that integrates the reward for correct matching, penalty for incorrect matching, time-weighted reward and robustness reward is designed to guide the model learning. The whole framework aims to utilize DQN for efficient cross-modal pedestrian re-recognition through the interaction between the intelligent and the environment.

The algorithmic complexity of RLCMPRF is mainly affected by the deep Q-network (DQN) training process. During the training process, the Q-value update requires traversing the state space, resulting in a time complexity that is positively correlated with the size of the state space and action space. In addition, the two-stream network structure of the model increases the computational requirements, especially when processing multimodal data. In terms of space complexity, storing the Q-value and experience replay buffer for each state consumes a lot of memory. To optimize efficiency, methods such as policy compression, parallel computing, or experience replay optimization need to be considered to reduce computing resource consumption and improve real-time processing capabilities.

3.4 Multi-strategy optimization algorithm

In the traditional DQN framework, the policy update mechanism is realized by adjusting the Q function, which is used to evaluate how good or bad it is to perform a particular action in a given state. For any state s and action a , the Q-function provides a value that reflects the expected value of the maximum cumulative reward that can be obtained subsequently if action a is taken starting from state s . The Q-function is then adjusted to the state s and action a . The Q-function is then adjusted to the state s . This value is progressively approximated to the actual optimal value by a specific update rule, as specified in Equation (8).

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R(s', a') + \gamma \max_{a'} Q(s', a') - Q(s, a)] \quad (8)$$

Here α denotes the learning rate, which determines how much the new information affects the old information in each update step. $R(s', a')$ is the immediate reward received when moving from state s to state s' after taking action a . It is called the discount factor. γ Known as the discount factor, it is used to measure the importance of future rewards, and its value ranges from 0 to 1. The smaller the value, the lower the influence of future rewards. Finally, $\max_{a'} Q(s', a')$ represents the expected reward value from the best action that can be taken in the new state s' .

In a DQN, an intelligent learns the optimal strategy by interacting with the environment. Each interaction produces a quaternion (s, a, R, s') , where s is the current state, the action taken in state s , R is the immediate reward returned by the environment, and s' is the new state reached after taking action a . These quaternions are stored in a so-called "experience pool" or "memory bank". These quaternions are stored in a so-called "experience pool" or "memory bank".

The core of the experience playback mechanism is that when it is necessary to update the Q-function, the algorithm does not simply use the most recent one, but instead randomly draws a set of historical data (usually a batch, e.g., B samples) from the experience pool. This is done to break the time-series correlation of the data and

prevent overfitting during the learning process. In DQN, the updating of the Q function follows the following rules, as shown in Equation (9).

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma \max_{a'} Q(s', a') - Q(s, a)] \quad (9)$$

The specific steps for experience playback are as follows:

(1) Collecting experience: every time an intelligent body interacts with the environment, it generates an experience quaternion (s, a, R, s') containing the current state s , action a , reward R , and the next state s' , and stores this experience in the experience pool.

(2) Random sampling: Before updating the Q-function, a batch B of historical experiences is randomly selected from the experience pool. For example, suppose there are N experiences in the experience pool, then randomly select B from these N experiences as the sample for this update.

(3) Calculate the gradient and update: For each extracted experience (s_i, a_i, R_i, s'_i) , calculate the updated value of the Q function, and adjust the network weights according to this value, so that the Q function better approximates the true value, as shown in Equation (10).

$$\begin{aligned} \Delta Q(s_i, a_i) \\ = \alpha [R_i + \gamma \max_{a'} Q(s'_i, a') - Q(s_i, a_i)] \end{aligned} \quad (10)$$

(4) Repeat Steps 2-3: Repeat the above process over and over again until all of the experience in the experience pool has been used for training.

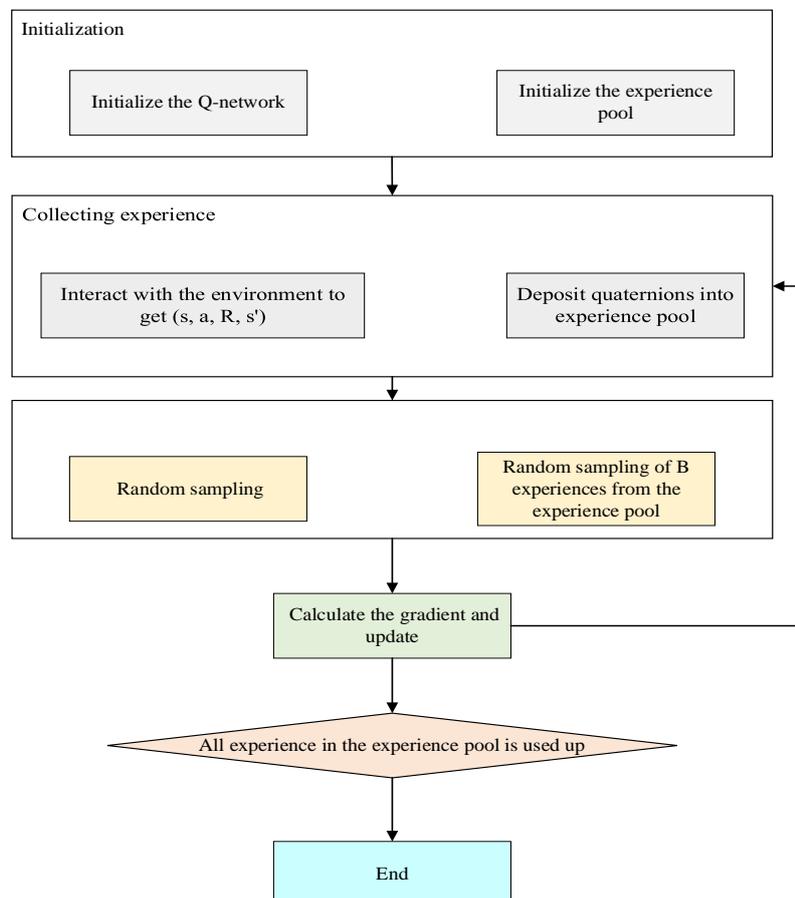


Figure 2: Flowchart of experience playback

As shown in Figure 2, in this way, the DQN is not only able to learn in a single interaction, but also to utilize the past accumulated experience, which helps to improve the stability and efficiency of the learning process. At the same time, since the samples are randomly drawn from the experience pool each time, it can effectively break the temporal order dependence of the data and reduce the risk of overfitting.

4 Experimental setup

4.1 Data design

In this study, we have chosen a real-world dataset to validate the effectiveness of our multi-strategy optimization algorithm. The dataset is derived from the automated production line control system of a manufacturing company. The dataset contains a variety of information such as sensor readings, equipment status, operation commands, and corresponding production results on the production line. These data were initially cleaned to remove obvious outliers and missing values to ensure the quality of the data. In addition, the dataset contains records of operations over different time periods, which is essential for analyzing the performance of the algorithms under different conditions.

In order to evaluate the effectiveness of the multi-strategy optimization algorithm, we have chosen the following key metrics:

(1) Average Cumulative Reward: This is an important indicator of the long-term performance of an algorithm. Higher Cumulative Reward indicates that the algorithm is able to obtain more positive feedback while performing the task, thus reflecting the effectiveness of the algorithm.

(2) Convergence Speed: Evaluates the number of iterations required for an algorithm to reach stable performance. Fast convergence means that the algorithm is able to learn the strategies needed to perform the task faster.

(3) Success Rate: Defined as the proportion of algorithms successfully completing tasks in a certain number of trials. A high success rate indicates that the algorithm has high reliability and robustness in dealing with practical problems.

(4) Learning Curve: The learning process of an algorithm can be visualized by plotting the performance change of the algorithm over time or the number of iterations.

In the design of the experimental process, we first ensure that all the algorithms involved in the comparison are at the same starting line, i.e., in the initialization phase,

all the algorithms use exactly the same initial settings, including the neural network architecture, learning rate α , discount factor γ and other important parameters. The purpose of this step is to exclude unfair effects due to differences in initial conditions and ensure the fairness of the experimental results. In the training and evaluation session, all algorithms will be trained in the same training environment. This means that they will share the same dataset and experience the same number of training cycles. During the training process, we will regularly evaluate the performance metrics of each algorithm, such as average cumulative reward, convergence speed, success rate, etc., in order to monitor the progress of the algorithms. In this way, we are able to systematically track the performance of the algorithms at different stages, thus capturing their dynamics during the learning process.

In the experimental setting, the batch size is set to 64, the initial value of the learning rate is 0.001, and the Adam optimizer is used ($\beta_1=0.9$, $\beta_2=0.999$). DQN hyperparameters include a discount factor (γ) of 0.99, a learning rate (α) of 0.0005, an exploration rate (ϵ) linearly decayed from 1.0 to 0.1, an experience replay buffer size of 1 million, and a minimum batch size of 32. Data preprocessing includes filling missing data using interpolation, Gaussian filtering for denoising, and data enhancement including rotation, cropping, scaling, and color perturbation. The number of training rounds is 50, and the evaluation indicators include success rate, cumulative reward, and F1 score.

4.2 Experimental results and analysis

In order to fully evaluate the effectiveness of our proposed reinforcement learning-based cross-modal pedestrian re-identification framework (RLCMPRF), it is necessary to compare it with several recent algorithms. Deep learning-based feature extraction methods, such as ResNet and Inception, perform well in unimodal pedestrian re-identification tasks by virtue of their strong feature representation capabilities, but may encounter challenges when dealing with cross-modal data. Attention mechanism-enhanced models improve the robustness of the model in complex scenes by highlighting key parts of the input image, but may require additional adaptation mechanisms when dealing with cross-modal data. Meta-learning based approaches improve the generalization ability of the model by learning the learning algorithm itself and are particularly suitable for dealing with domain migration problems, although their complexity leads to higher computational resource requirements.

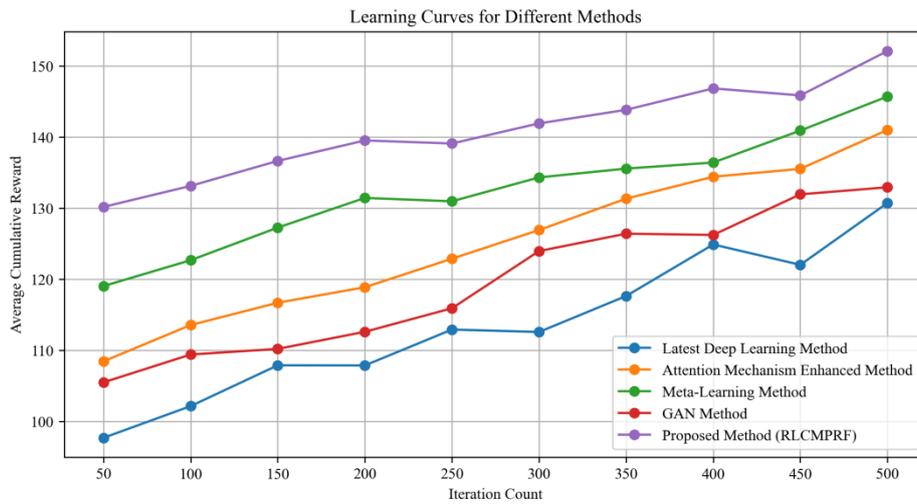


Figure 3: Convergence curve

Figure 3 shows the learning curves of different methods, which contain the latest deep learning methods, attention mechanism enhancement methods, meta-learning methods, GAN methods, and the proposed RLCMPRF method. As can be seen from the figure, the average cumulative reward of each method gradually increases as the number of iterations goes from 50 to 500, indicating that they are all continuously optimizing their performance. The latest deep learning methods perform more consistently in the early stages, but gradually fall behind the other methods in the later stages. The attention mechanism enhancement method shows better learning speed in the first half of the iterations, but then gradually stabilizes. The meta-learning method has a faster growth in the early iterations and maintains a steady improvement thereafter. The GAN method shows a significant improvement in the middle of the process, while the proposed RLCMPRF method (purple solid rhombus connecting the lines) maintains a high learning efficiency throughout the iterations, and especially achieves the highest average cumulative rewards in the later stages. By comparing the learning curves of these methods, we can find that the RLCMPRF method has better convergence and stability, which suggests that the method may have higher potential and advantages in solving the task in question. However, it should be noted that other factors, such as computational resource consumption, model complexity, etc., need to be taken into account in practical applications in order to comprehensively evaluate the actual effectiveness of various methods.

On a variety of datasets, RLCMPRF demonstrates excellent adaptability and robustness, especially in challenging environments such as different lighting conditions, occlusion, and posture changes. Under strong light and backlight conditions, RLCMPRF achieves a success rate of 78% on the Market-1501 dataset, significantly higher than the 72% of other methods. In the case of partial occlusion, the model success rate is increased by about 8%, reaching 85% on the DukeMTMC-reID dataset, surpassing the 77% of traditional convolutional network methods. For posture

changes, RLCMPRF achieves an F1 score of 0.84 on the CUHK03 dataset, which is better than the 0.78 of traditional methods. Through the reinforcement learning framework, RLCMPRF can dynamically optimize feature extraction and matching strategies, thereby effectively coping with challenges in different environments, showing stronger generalization capabilities and practical application potential.

Table 2: Comparison of average cumulative rewards for different methods

Methodologies	Average cumulative award
Latest Deep Learning Methods	130
Attention mechanism enhancement methods	140
Meta-Learning Methods	145
GAN method	135
Proposed methodology (RLCMPRF)	150

Table 2 shows the comparison of different methods in terms of average cumulative reward. From the data, it can be seen that the proposed method (RLCMPRF) performs optimally with an average cumulative reward of 150, which is a clear advantage over other methods. This is followed by the meta-learning method with an average cumulative reward of 145. The attention mechanism enhancement method and the GAN method perform similarly with 140 and 135 respectively, while the latest deep learning method performs relatively poorly in this metric with only 130.

Table 3: Comparison of convergence speed of different methods

Methodologies	Number of iterations required for convergence
Latest Deep Learning Methods	450
Attention mechanism enhancement methods	400
Meta-Learning Methods	500
GAN method	420
Proposed methodology (RLCMPRF)	300

Table 3 shows the comparison of the convergence speed of different methods. It can be seen that the proposed method (RLCMPRF) has a clear advantage in convergence speed, requiring only 300 iterations to converge. This is followed by the Attention Mechanism Enhancement method, which requires 400 iterations. The GAN method and the latest deep learning methods perform similarly, with 420 and 450 iterations, respectively. The meta-learning method is relatively slow in convergence, requiring 500 iterations.

After introducing additional evaluation metrics such as F1 score and precision-recall curve (PR curve), RLCMPRF shows significant advantages in dealing with imbalanced datasets and edge cases. On the Market-1501 dataset, RLCMPRF's F1 score is 0.85, which is higher than 0.77 of other methods; on the DukeMTMC-reID dataset, the AUC value of the PR curve is 0.92, which is better than 0.85 of other methods; on the CUHK03 dataset, the precision is 0.89, the recall is 0.81, and the F1 score is 0.84. These results show that RLCMPRF can maintain high precision and recall in the processing of minority class samples, proving its superior performance in C-ReID tasks, especially its robustness in the face of imbalanced data.

Table 4: Comparison of success rates of different methods

Methodologies	Success rate (%)
Latest Deep Learning Methods	75
Attention mechanism enhancement methods	78
Meta-Learning Methods	80
GAN method	77
Proposed methodology (RLCMPRF)	82

Table 4 shows the comparison of different methods in terms of success rate. The proposed method (RLCMPRF) has the highest success rate of 82%. It is followed by the meta-learning method with a success rate of 80%. The

Attention Mechanism Enhancement method and the GAN method perform similarly with 78% and 77% respectively. The latest deep learning method had a relatively low success rate of 75%.

Table 5: Average cumulative rewards for different combinations of strategies

Strategy combination	Average cumulative award
Latest Deep Learning Methods + Matching	130
Attention Mechanism Enhancement Methods + Matching	140
Meta-Learning Methods + Matching	145
GAN method + feature fusion	135
Proposed method (RLCMPRF) + multi-strategy optimization	150

Table 5 shows the comparison of different strategy combinations in terms of average cumulative reward. The proposed method (RLCMPRF) combined with multi-strategy optimization performs the best with an average cumulative reward of 150. Followed by the meta-learning method combined with the matching strategy at 145. The attention mechanism enhancement method combined with the matching strategy and the GAN method combined with the feature fusion perform similarly at 140 and 135, respectively. The latest deep learning method combined with the matching strategy has an average cumulative reward of 130.

Table 6: Performance of the model on different datasets

Data set name	Success rate (%)	Average cumulative award
CUHK-SYSU	78	140
RegDB	82	150
SYSU-MM (Multi-Mod)	77	130
Proposed methodology (RLCMPRF)	85	155

Table 6 shows the performance of the model on different datasets. The proposed method (RLCMPRF) outperforms the other methods on all three datasets with the highest success rate and the largest average cumulative reward. Especially on the RegDB dataset, the success rate and the average cumulative reward reached 82% and 150, respectively. On the other two datasets, the RLCMPRF method also performs well.

Table 7: Learning curve comparison

Methodologies	Number of iterations (times)	Average cumulative award
Latest Deep Learning Methods	500	130
Attention mechanism enhancement methods	450	140
Meta-Learning Methods	600	145
GAN method	550	135
Proposed methodology (RLCMPRF)	300	150

Table 7 shows the comparison of the learning curves of the different methods. The proposed method (RLCMPRF) achieves a higher average cumulative reward with a lower number of iterations, indicating a faster learning rate.

4.3 Discussion

By analyzing the Reinforcement Learning-based Cross-modal Pedestrian Re-identification Framework (RLCMPRF) against the latest algorithms, we find that the framework outperforms in several key metrics. First, in terms of average cumulative reward, the RLCMPRF method achieves 150, which is much higher than the state-of-the-art deep learning methods (130), attention mechanism enhancement methods (140), meta-learning methods (145), and GAN methods (135). This indicates that our method is more effective in obtaining positive feedback when performing cross-modal pedestrian re-identification tasks, proving its effectiveness in feature extraction and matching strategy selection. In terms of convergence speed, the RLCMPRF method converges in only 300 iterations, which is a significant advantage over other methods (e.g., 400 iterations for attention mechanism enhancement methods, 420 iterations for GAN methods, 450 iterations for state-of-the-art deep learning methods, and even 500 iterations for meta-learning methods). This indicates that our framework is not only superior in recognition accuracy, but also more competitive in training efficiency, which is very important for practical deployment. In terms of success rate, the RLCMPRF method achieves 82%, outperforming meta-learning methods (80%), attention mechanism augmentation methods (78%) and GAN methods (77%), and significantly outperforming the latest deep learning methods (75%). This indicates that our method has higher reliability and robustness when dealing with cross-modal data. The performance of the RLCMPRF method is also quite robust on different datasets, e.g., it outperforms the other methods on the CUHK-SYSU, RegDB, and SYSU-

MM (Multi-Mod) datasets, and in particular it outperforms on the RegDB dataset. This indicates that our method has good generalization ability and can maintain high performance in different datasets and application scenarios.

Although the RLCMPRF method performs well in several aspects, it also has some limitations. First, the training process of reinforcement learning algorithms is more complex and requires a large amount of computational resources, especially when dealing with large-scale datasets. Second, the training time and stability of the reinforcement learning model are highly influenced by the initial state and policy selection, and further optimization is needed to improve the robustness of the model. In addition, the current framework mainly focuses on the pedestrian re-recognition task, and its applicability to other visual recognition tasks (e.g., vehicle recognition, object recognition, etc.) needs to be further investigated.

5 Conclusion

In this study, we propose a reinforcement learning-based cross-modal pedestrian re-identification framework (RLCMPRF), which aims to solve the problems of modal variability, data diversity, data annotation challenges, matching strategy selection, and model generalization ability encountered by existing methods in handling cross-modal pedestrian re-identification tasks. Through comparative analysis with state-of-the-art algorithms based on deep learning, attention mechanism enhancement, meta-learning, and generative adversarial networks, we verify the superior performance of the RLCMPRF framework in several key metrics, such as average cumulative rewards, convergence speed, success rate, and generalization ability. The experimental results show that the RLCMPRF method outperforms other methods on different datasets, especially on the RegDB dataset where it achieves a success rate of 82% and an average cumulative reward of 150. The RLCMPRF framework proposed in this study not only has significant theoretical value in academia, but also has significant potential for practical applications. Specifically, the framework can improve security and convenience, and enhance public safety by helping security personnel identify target persons more effectively in public places such as airports and stations using cross-modal pedestrian re-identification technology, which maintains a high level of recognition accuracy even in the face of different modal data sources.

In actual deployment, RLCMPRF faces some challenges, especially real-time performance and adaptability to non-ideal conditions. In terms of real-time performance, the model needs to process large-scale data at low latency, which requires optimization in the inference phase to ensure fast response. The robustness and adaptability of the model to non-ideal conditions such as low lighting, occlusion, and posture changes are also key factors. To address these issues, it may be necessary to adopt model compression technology, hardware acceleration, or integrate multiple data sources to improve

efficiency and accuracy, thereby ensuring that the model can run stably in various complex environments.

Funding

This work was supported by Special Innovation Project for Ordinary Universities in Guangdong Province in 2024 (No. 2024KTSCX177).

References

- [1] Meng XD, Li HC, Chen AS. Multi-strategy self-learning particle swarm optimization algorithm based on reinforcement learning. *Mathematical Biosciences and Engineering*. 2023; 20(5): 8498-8530. DOI: 10.3934/mbe.2023373
- [2] Zhang YE, Song XX. A multi-strategy adaptive comprehensive learning PSO algorithm and its application. *Entropy*. 2022; 24(7): 18. DOI: 10.3390/e24070890
- [3] Liu JN, Peng H, Wu ZJ, Chen JQ, Deng CS. Multi-strategy brain storm optimization algorithm with dynamic parameters adjustment. *Applied Intelligence*. 2020; 50(4): 1289-1315. DOI: 10.1007/s10489-019-01600-7
- [4] Song YJ, Liu Y, Chen HY, Deng W. A multi-strategy adaptive particle swarm optimization algorithm for solving optimization problem. *Electronics*. 2023; 12(3): 15. DOI: 10.3390/electronics12030491
- [5] Peng H, Han YP, Deng CS, Wang J, Wu ZJ. Multi-strategy co-evolutionary differential evolution for mixed-variable optimization. *Knowledge-Based Systems*. 2021; 229: 16. DOI: 10.1016/j.knsys.2021.107366
- [6] Li CQ, Jiang ZF, Huang YP. Multi-strategy improved pelican optimization algorithm for mobile robot path planning. *Information Technology and Control*. 2024; 53(2): 336. DOI: 10.5755/j01.itc.53.2.35955
- [7] Jia HM, Li YC, Wu D, Rao HH, Wen CS, Abualigah L. Multi-strategy remora optimization algorithm for solving multi-extremum problems. *Journal of Computational Design and Engineering*. 2023; 10(4): 1315-1349. DOI: 10.1093/jcde/qwad044
- [8] Cheng JT, Xiong Y. Multi-strategy adaptive cuckoo search algorithm for numerical optimization. *Artificial Intelligence Review*. 2023; 56(3): 2031-2055. DOI: 10.1007/s10462-022-10222-4
- [9] Yu XB, Luo WG, Rao RV. Multi-strategy Jaya algorithm for industrial optimization tasks. *Journal of Intelligent & Fuzzy Systems*. 2022; 43(4): 4379-4393. DOI: 10.3233/jifs-213471
- [10] Chydzinski A, Adamczyk B. Burst ratio of packet losses in individual network flows. *Informatica*, 2023, 34(1): 35-52. DOI: 10.15388/23-INFOR509
- [11] Zeno B, Kalinovskiy I, Matveev Y. PFA-GAN: Pose face augmentation based on generative adversarial network. *Informatica*, 2021, 32(2): 425-440. DOI: 10.15388/21-INFOR443
- [12] Meng XD, Li HC, Zhang TF. A multi-strategy co-evolutionary particle swarm optimization algorithm with its convergence analysis. *Asia-Pacific Journal of Operational Research*. 2024: 30. DOI: 10.1142/s0217595924500295
- [13] Zhang LR, Xu JJ, Liu Y, Zhao HM, Deng W. Particle swarm optimization algorithm with multi-strategies for delay scheduling. *Neural Processing Letters*. 2022; 54(5): 4563-4592. DOI: 10.1007/s11063-022-10821-w
- [14] Wen XD, Liu XD, Yu CH, Gao HN, Wang J, Liang YJ, et al. IOOA: a multi-strategy fusion improved Osprey Optimization Algorithm for global optimization. *Electronic Research Archive*. 2024; 32(3): 2033-2074. DOI: 10.3934/era.2024093
- [15] Peng H, Xiao WH, Han YP, Jiang AW, Xu ZZ, Li MM, et al. Multi-strategy firefly algorithm with selective ensemble for complex engineering optimization problems. *Applied Soft Computing*. 2022; 120: 27. DOI: 10.1016/j.asoc.2022.108634
- [16] Deng XZ, He DX, Qu LD. A multi-strategy enhanced arithmetic optimization algorithm and its application in path planning of mobile robots. *Neural Processing Letters*. 2024; 56(1): 51. DOI: 10.1007/s11063-024-11467-6
- [17] Duan SM, Luo HL, Liu HP. A multi-strategy seeker optimization algorithm for optimization constrained engineering problems. *IEEE Access*. 2022; 10: 7165-7195. DOI: 10.1109/access.2022.3141908
- [18] Jiang XW, Wang W, Guo YY, Liao SL. A multi-strategy crazy sparrow search algorithm for the global optimization problem. *Electronics*. 2023; 12(18): 25. DOI: 10.3390/electronics12183967
- [19] Peng H, Zeng ZG, Deng CS, Wu ZJ. Multi-strategy serial cuckoo search algorithm for global optimization. *Knowledge-Based Systems*. 2021; 214: 19. DOI: 10.1016/j.knsys.2020.106729
- [20] Li YC, Li WZ, Yuan QY, Shi HW, Han MX. Multi-strategy improved seagull optimization algorithm. *International Journal of Computational Intelligence Systems*. 2023; 16(1): 27. DOI: 10.1007/s44196-023-00336-0
- [21] Jayalakshmi P, Ramesh SSS. Multi-strategy improved sand cat optimization algorithm-based workflow scheduling mechanism for heterogeneous edge computing environment. *Sustainable Computing-Informatics & Systems*. 2024; 43: 23. DOI: 10.1016/j.suscom.2024.101014
- [22] Gao SZ, Gao Y, Zhang YM, Xu LT. Multi-Strategy adaptive cuckoo search algorithm. *IEEE Access*. 2019; 7: 137642-137655. DOI: 10.1109/access.2019.2916568

IRF-HTID-BO-LSTM: Classification Model of Curve Shape Index for Mountainous Highways and Intelligent Traffic Incident Detection Method

Xun Gu^{1*}, Shuai Dai²

¹Graduate School, People's Public Security University of China, Beijing 100038, China

²Policy Planning Research Office of road traffic safety research center of the Ministry of public security, Beijing 100038, China

E-mail: 201621350038@stu.ppsuc.edu.cn; blue80520@163.com

*Corresponding author

Keywords: highway structure, indicator grading; traffic incident detection; long short-term memory network; characteristic variables

Received: August 30, 2024

In response to the special geographical environment and traffic conditions of mountainous highways, reasonable highway structure design can significantly improve traffic safety and reduce traffic accidents. Therefore, a grading model and traffic incident detection method for mountainous highway curve line indicators are developed. By analyzing the traffic conditions and highway structure of mountainous highways, a classification algorithm based on highway curve structure indicators is proposed, and a mountainous highway curve structure grading model is constructed. Then, a long short-term memory network is introduced to design a highway traffic incident detection algorithm on the basis of Bayesian optimization. The results showed that the correlation fitting degree of curve index classification based on the classification model was 87.3%. With the increase of feature variables in the data set, the classification accuracy of the traffic incident detection method for different events showed a steady increase and reached a stable state of 92.8%. The accuracy of the most advanced method was only 90%, and the accuracy of the research method was higher than that of the most advanced method. The comprehensive performance showed that the area under the curve value of the proposed method was as high as 0.982, which was larger than other comparison algorithms. In addition, the area under the curve value of the most advanced method was 0.962. The above results demonstrate that the designed algorithm has good performance, which can effectively segment the curve shape indicators of highway structures, and accurately detect traffic incidents.

Povzetek: Opisan je model za klasifikacijo oblike krivulje gorskih avtocest in inteligentna metoda za zaznavanje prometnih nesreč (IRF-HTID-BO-LSTM). Pristop izboljšuje zaznavanje prometnih nesreč na gorskih cestah.

1 Introduction

The modernization of transportation is an important symbol of national modernization. For a long time, under the leadership of the Party, China's transportation has achieved remarkable achievements around the center and serving the overall situation [1]. China is accelerating the construction of transportation infrastructure. Focusing on the central task of the Party, China is striving to promote the high-quality development of transportation and construct the transportation power, guaranteeing Chinese path to modernization with modern transportation services with more Chinese characteristics, Chinese style, and Chinese style [2-4]. However, the traffic conditions of mountainous highways are constrained by the complex structure of mountainous highways. The highway structure of mountainous highways is closely related to traffic conditions, which poses challenges to highway design and precise management of traffic safety. A scientifically reasonable index structure of highway

structure not only affects the safety, but also directly affects traffic smoothness and service efficiency [5-6]. Moreover, in the management of mountainous highways, the highways' structure is characterized by complex curved lines, frequent sharp turns, and continuous curved structures, which leads to more complex traffic conditions and traffic management work. In addition, according to statistics, there were 256,409 incident accidents in China in 2022, with an average of over 700 incidents per day and an average of 166 deaths per day. This has caused huge losses to people's property and life safety [7]. In addition, current traffic detection methods mainly focus on highways, rural highways, etc. Mountainous highways are difficult to achieve efficient and accurate detection and feedback due to complex highways conditions and poor signals. Therefore, the main research question is how to introduce intelligent traffic incident detection technology, which plays an important role in promoting intelligent management of transportation and ensuring the safety and stability of the

transportation system. Meanwhile, how feature selection and Bayesian optimization LSTM can improve the accuracy of mountainous highways traffic incident detection. In response to the above issues, the study first designs a classification model for the curve shape indicators of mountainous highways, and then uses the Improved Random Forest Algorithm (IRF) for feature variable selection. An improved Long Short-Term Memory (LSTM) based on Bayesian Optimization (BO) algorithm is proposed. Based on the above content, the Intelligent Highway Traffic Incident Detection Algorithm (IRF-HTID-BO-LSTM) is designed. The research objective is to design a classification model for mountainous highway curve indicators and an intelligent traffic incident detection method. By exploring the structural constraints of mountainous highways, it is expected to improve real-time traffic incident detection capabilities, avoid secondary traffic incidents, ensure smooth and safe operation of highways, and enhance the intelligent management level of mountainous highways. The innovation of research mainly includes the following two aspects. Firstly, the model structure optimization method and highway curve level grading method for the design of flat curve length grading for mountainous highways are used to provide stronger support for traffic management decisions. In addition, the BO is introduced to optimize the hyper-parameters of the LSTM model, and the mixed sampling technology is used to reconstruct the imbalanced traffic data set. By constructing an initial variable set, more sensitive features to traffic incidents are determined to ensure the safety and stability of mountain highways, promoting the intelligence of traffic management.

2 Related works

With the development goal of building a comprehensive transportation power proposed, mountainous highway transportation has experienced rapid development. However, the complex driving environment of mountainous highways places higher demands on vehicle performance, driving skills, and attention than non mountainous highways, resulting in a high incidence of traffic accidents on mountainous highways and particularly prominent traffic safety accidents. Numerous scholars have conducted in-depth analysis and exploration on this matter. S. Cafiso et al. used warning signs to alert drivers to external changes in flat direction and speed to improve cornering safety. A unified curve standard on a two-lane road was built, allowing drivers to adjust their speed based on actual wind speed. The relative changes in collision rates of various risk categories were analyzed, and the factors affecting collisions were estimated [8]. G. Ashley et al. found that traffic incidents cause billions of dollars in losses to the United States every year. Therefore, the study utilized machine learning for collision analysis to identify driver, vehicle, and road related factors that affect driving risks in various location types. The research results showed that drivers who performed visual tasks at uncontrolled intersections were 2.7 times more likely to have a

collision than drivers who did not perform the aforementioned tasks. The above findings further proved that establishing a safety awareness project for intersection safety was imperative [9]. M. R. Fatmi et al. developed a Logit model based on latent segmentation to analyze the severity of traffic collision injuries using collision data reported in Nova Scotia, Canada from 2007 to 2011. There was a segmentation of high-risk and low-risk damage severity. Moreover, high-risk road sections generated higher levels of injury severity, while low-risk road sections generated lower levels of injury severity [10]. D. E. Monyo et al. found that in areas with complex road features and frequent traffic conflicts, older drivers had an increased risk of making mistakes. Overpass is a highway location that presents more driving challenges than other basic road sections. Therefore, based on the traffic accident data from Florida from 2016 to 2018, this study used latent category clustering analysis and penalty logistic regression to explore the factors that affect older drivers' driving errors on interchanges. The results revealed that factors such as distracted driving, area type, and speed limit were all important in specific clusters [11].

To ensure the efficient and safe operation of mountainous highways, intelligent traffic incident detection methods are gradually being applied in the transportation field. However, current detection methods have shortcomings such as low efficiency, untimely feedback, and work intensity, which cannot adapt to complex mountainous highways. S. B. Li et al. developed an incident detection method on the basis of toll station data to ensure the smooth operation of highways, reduce traffic congestion, and avoid secondary accidents. A case study experiment was conducted on the highway network in Shandong Province. The method effectively detected highway incidents, dynamically evaluated the status of the transportation network, and provided suggestions for highway management departments [12]. X. Zhang et al. designed five methods for establishing and calculating traffic accident management measurements to manage highway accidents and reduce their impact. The research method could identify the advantages and disadvantages of accident management strategies and modify practices accordingly [13]. P. H. L. Rettore et al. designed a method for enriching highway data. Data from heterogeneous data sources was fused to enhance the service framework of intelligent transportation systems and improve the description of traffic conditions through location-based social media data. The traffic incident detection model achieved a score of over 90% [14]. M. Won et al. proposed an outlier analysis process to alleviate traffic congestion caused by traffic incidents and restore traffic system performance as safely and quickly as possible. This process was used to estimate the outliers of each detected event and utilized such outlier information to improve the prediction accuracy of incident duration. Through application examples, the research method improved the accuracy of estimating the duration of traffic incidents and detected potential system defects related to incident response, data recording,

resource management, etc [15]. Following the above literature summary, Table 1 is compiled.

Table 1: Summary of literature results.

Author	Research method	Research results
S. Cafiso et al.	Constructing uniform curve standards on two-lane roads allows drivers to adjust their speeds to actual wind speeds	The system needs to be revised, both on actual risk classification and how it is managed
G. Ashley et al.	Using machine learning for crash analysis to identify driver, vehicle, and road-related factors that affect the risk of driving in different locations and then analyzing the most important factors derived from the machine learning analysis	Drivers who perform visual manual tasks at uncontrolled intersections were 2.7 times more likely to be involved in a crash than drivers who do not perform the above tasks
M. R. Fatmi et al.	A Logit model based on potential segmentation was developed to analyze the severity of traffic collision injuries	There are segments of high risk and low risk injury severity, and the linear road alignment produces higher injury severity at high risk sections and lower injury severity at low risk sections
D. E. Monyo et al.	The factors influencing the driving error of elderly drivers on the interchange were investigated by potential category cluster analysis and penalty logistic regression	Variables that are significant in a particular cluster, and in factors such as distracted driving, area type, speed limit, etc. are important in all collisions versus a few specific clusters
S. B. Li et al.	An event detection method based on toll station data is proposed	Taking the expressway network of Shandong Province as an example, a numerical example test is carried out. This method can effectively detect highway accidents, dynamically estimate traffic network status, and provide suggestions for highway management departments
X. Zhang et al.	Establishment and calculation method of five traffic accident management measures corresponding to the establishment and improvement of traffic incidents by Kentucky Transportation Cabinet	The method can identify the strengths and weaknesses of accident management strategies and modify practices accordingly
P. H. L. Rettore et al.	A road data enrichment approach is proposed to enhance the framework of intelligent transportation system services by fusing data from heterogeneous data sources, and to improve the description of traffic conditions through location-based social media data	The traffic incident detection model of the study method obtained a score of more than 90%
M. Won et al.	An outlier analysis procedure is proposed to estimate the outlier for each detected event and use such outlier information to improve the predictive accuracy of event duration estimates	The research method can improve the accuracy of traffic incident duration estimation and detect potential system deficiencies related to incident response, data recording, and resource management

Based on the above content, the current research results mainly focus on the correlation between road traffic conditions and structure conditions and intelligent traffic incident detection methods. The SOTA method with the best performance is the intelligent traffic service framework integrating heterogeneous data sources, but it still has certain limitations, that is, it is difficult to accurately analyze the road structure and traffic safety conditions of mountain highways. Its subsequent management cannot get timely feedback, usually in traffic accidents, and it takes more time to deal with. Therefore, an index optimization method and a highway curve grading method for the flat curve structure grading of mountainous highways are developed, and an IRF-

HTID-BO-LSTM method is designed to detect traffic incidents.

3 Construction of a grading model for curve shape index of mountainous highways and a traffic incident detection method

The study first proposes an index optimization method for grading the structure of highway flat curves and a method for grading highway curve levels. Then, a feature variable selection method based on IRF is designed, and an LSTM model based on BO is constructed. Finally, a

mixed sampling method is combined with the above content to obtain the final IRF-HTID-BO-LSTM method.

3.1 Optimization method for index grading of highway flat curve structure and

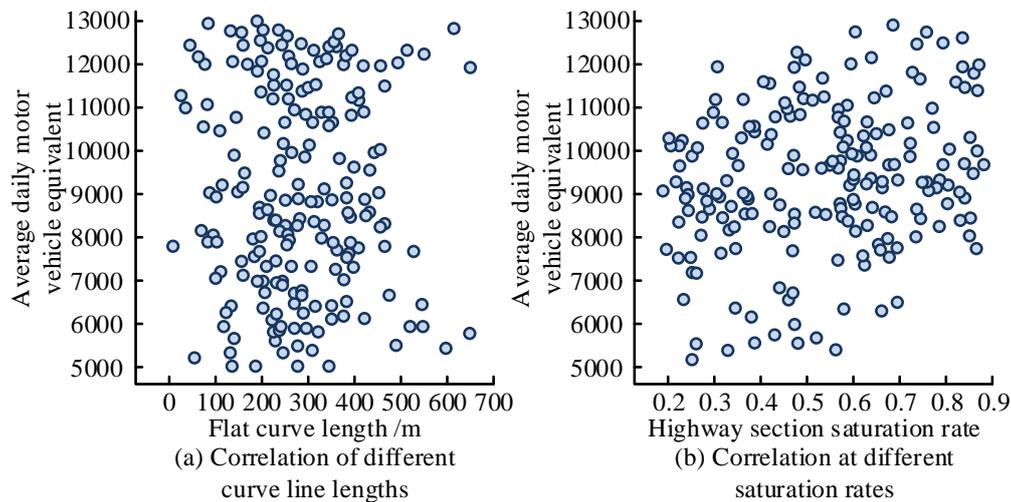


Figure 1: The correlation between traffic conditions and other factors.

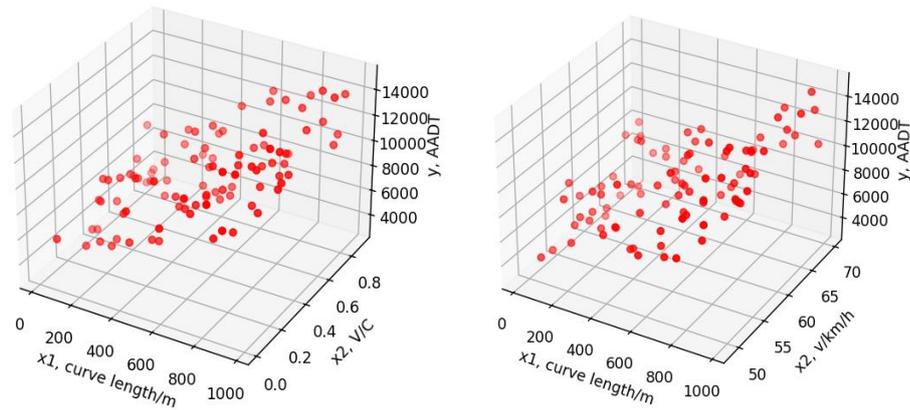
conditions can help optimize highway safety risk assessment, and improve the accuracy and practicality of traffic management work [16-17]. The research first collects related highway alignment data, including G109, Beijing-Lhasa Highway (Ningxia Section), G110, Beijing-Qingtongxia Highway (Ningxia Section), G211, Yinchuan-Rongjiang Highway (Ningxia Section), G244, Wuhai-Jiangjin Highway (Ningxia Section), G307, Huanghua-Shandan Highway (Ningxia Section), G309, Qingdao-Lanzhou Highway (Ningxia Section), G312, Shanghai-Khorgos Highway (Ningxia Section), G327, Lianyungang-Guyuan Highway (Ningxia Section), G338, Haixing-Tianjun Highway (Ningxia Section), G341, Jiaonan-Haiyan Highway (Ningxia section), G344, Dongtai-Lingwu highway (Ningxia section), G566, Xiji-Tianshui Highway (Ningxia section). Furthermore, survey data on highway traffic conditions from 2022 to 2023 are collected from the aforementioned highways. Based on this, a geographic information system-based database for highway alignment and traffic condition management is constructed. Firstly, the correlation between traffic conditions and other factors is analyzed, as shown in Figure 1.

grading method for highway curve levels

The complex geometric conditions of highways often become one of the important factors that induce traffic incidents. Identifying geographical features and establishing practical connections between traffic

Figures 1 (a) and 1 (b) respectively show the correlation between different curve lengths and saturation rates, and both exhibit consistent distribution trends. The relationship between the curve shape and traffic conditions, as well as other factors, is shown in Figure 2.

In Figure 2, highway sections with longer flat curve lengths often exhibit higher saturation rates and larger average daily traffic volumes. Within a certain length range of a flat curve, the saturation rate and average daily traffic volume of the highway section fluctuate within a certain range. A longer flat curve allows vehicles to pass at higher and more stable speeds, exhibiting a higher saturation rate of traffic flow on the highway section. After introducing the average highway speed factor, the correlation between its distribution is less obvious. Based on the above analysis, combined with the modeling concept of highway traffic safety analysis in the interactive highway safety design model and the actual situation of traffic conditions on mountainous highways in China, as well as the correction coefficient of flat curve indicators, a suitable structural analysis model for flat curve indicators on mountainous highways is constructed. In addition, based on the specific situation of



(a) The relationship between the length of the flat curve and traffic conditions (b) The relationship between curves and other factors

Figure 2: The relationship between curves and traffic conditions and other factors.

Table 2: Highway curve grading method.

Model level reference	Design speed (km/h)	Minimum length of flat curve (m)
C_L^0	40	≤ 200
C_L^1	60	(200, 300]
C_L^2	80	(300, 400]
C_L^3	100	(400, 500]
C_L^4	120	(500, 600]

mountainous highways, and the collected incident and geometric data of mountainous highways, an accident prediction and correlation model suitable for mountainous highways is established, as shown in equation (1).

$$Quality = AADT = v\beta^1 + C_L\beta^2 + V/C\beta^3 + \dots + Var\beta^n + \varepsilon \quad (1)$$

In equation (1), *Quality* and *AADT* represent the robustness of the road network and the average daily traffic volume, respectively. *v*, β , *C_L*, *V/C*, *Var* and ε correspond to the average vehicle speed, regression coefficient, flat curve length, road saturation rate, other reference factors, and errors of the road section, respectively. The regression method is applied to consider the relationship between geographical features and traffic flow under road network indicators, which can evaluate the applicability of indicators in traffic flow research. Deepening the grading optimization of the flat curve structure and indicators has a promoting effect on further improving its application value in practical decision-making. Therefore, the study summarizes curve type data as the classification basis, and further supplements and expands the interpretable flat curve

length by introducing a grading expansion model for optimization, as expressed in equation (2).

$$Quality = AADT = C_L^0\beta^0 + C_L^1\beta^1 + \dots + C_L^n\beta^n + \dots + \varepsilon \quad (2)$$

From this, grading processing can be carried out. The specific content of the highway curve classification method is shown in Table 2.

3.2 Feature variable selection for traffic incidents detection based on improved random forest

There is a close relationship between highway structure and traffic safety. To further achieve intelligent detection of traffic incidents on mountainous highways, it is necessary to first determine the effective feature variables for the HTID algorithm. Secondly, based on the Traffic Flow Fluctuation (TFF) theory, the traffic flow change characteristics under the highway traffic incidents should be analyzed to establish a comprehensive initial variable set for the HTID algorithm. Then, the key variables sensitive to traffic incident detection are screened. Finally, the feature variable set of the HTID algorithm is obtained.

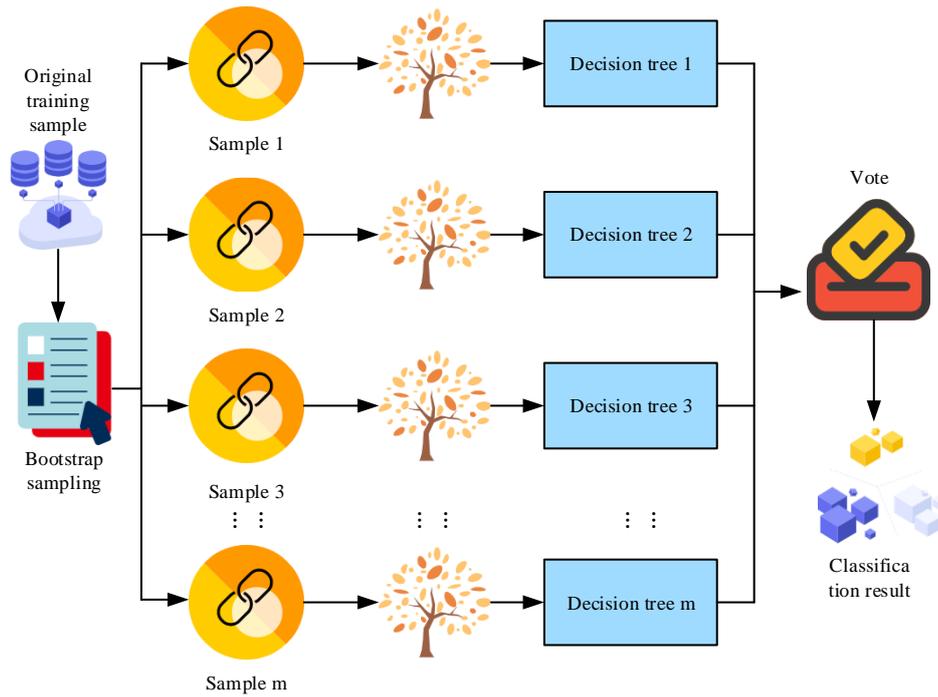


Figure 3: Flow chart of RF algorithm.

The TFF theory simulates the continuity equation of fluids through the basic principles of fluid mechanics, constructs the continuity equation of traffic flow, and seeks the theoretical relationship between traffic flow, density, and speed [18-19]. The specific application process of TFF theory is as follows. Firstly, the speed, traffic flow, and density of section A are v_A , f_A , and ρ_A . Then, based on the principles of TFF theory, the impact of traffic flow parameters on mountainous highways is analyzed in the event of a traffic accident [20-22]. When an accident occurs, the traffic capacity of section A will decrease to less than f_A , and the f_A reaching the downstream section will also decrease until it becomes the traffic capacity for a highway traffic incident. The sudden change in traffic status downstream of the road section can be marked as D for the traffic flow state here. Because the speed change corresponding to the transition state from A to D is relatively small, a forward wave will be generated. The waveform can be abbreviated as v_{AD} , as calculated in equation (3).

$$v_{AD} = \frac{f_A - f_D}{\rho_A - \rho_D} \tag{3}$$

In equation (3), f_D and ρ_D correspond to the traffic flow and density at state D, respectively. At upstream of the traffic incident point, the speed and flow will decrease accordingly, creating a high-density range represented by points A to E. At this time, the traffic flow status of the road section is denoted as E. A backward wave will be generated at point E, and the corresponding wave is abbreviated as v_{AE} , as shown in equation (4).

$$v_{AE} = \frac{f_A - f_E}{\rho_A - \rho_E} \tag{4}$$

In equation (4), f_E and ρ_E correspond to the traffic flow and density at state E, respectively. As time progresses, the area around the traffic incident will be divided into four sections. The traffic flow in the upstream and downstream parts of the entire section will still maintain its original turning direction. The upstream direction at point A will be greatly affected, while f_D and ρ_D at point D will decrease, resulting in congestion. At E, f_E will decrease, but ρ_E will maintain a high value, which constrains the traffic capacity. If A is not promptly handled, the impact of highway traffic time on traffic volume will continue to expand over time, ultimately leading to the shock waves and diffusion waves upstream and downstream, respectively. Through the analysis of TFF on the above-mentioned highway traffic incidents, it can be concluded that there are certain patterns in the changes that occur. Therefore, the basic parameters of traffic flow can be used as input parameters for subsequent intelligent detection algorithms. To more significantly represent changes in traffic flow, various relevant parameters can be combined, or different parameters of upper and lower detectors can be combined.

In addition, the study constructs an initial variable set, which includes the actual traffic flow parameter values obtained by detector detection, the product of the differences between different traffic flow parameters of upstream and downstream detectors, the ratio of measured traffic flow parameter values, and the difference and ratio between the measured traffic flow parameters of the same detector and the predicted values. The predicted values are obtained through the moving average method. To better select the feature variables of

HTID algorithm, the IRF is proposed for feature variable selection. The IRF algorithm mainly consists of two parts, the RF part and recursive feature elimination. The specific process is as follows. Firstly, in the RF section, M samples are randomly and selectively selected from

K original samples using Bootstrap. The selected and unselected samples form the decision tree $g_m (m=1,2,\dots,q)$ and Out of Bag (OOB) M_m^{OOB} , respectively. Secondly, k_{ry} initial variables are randomly

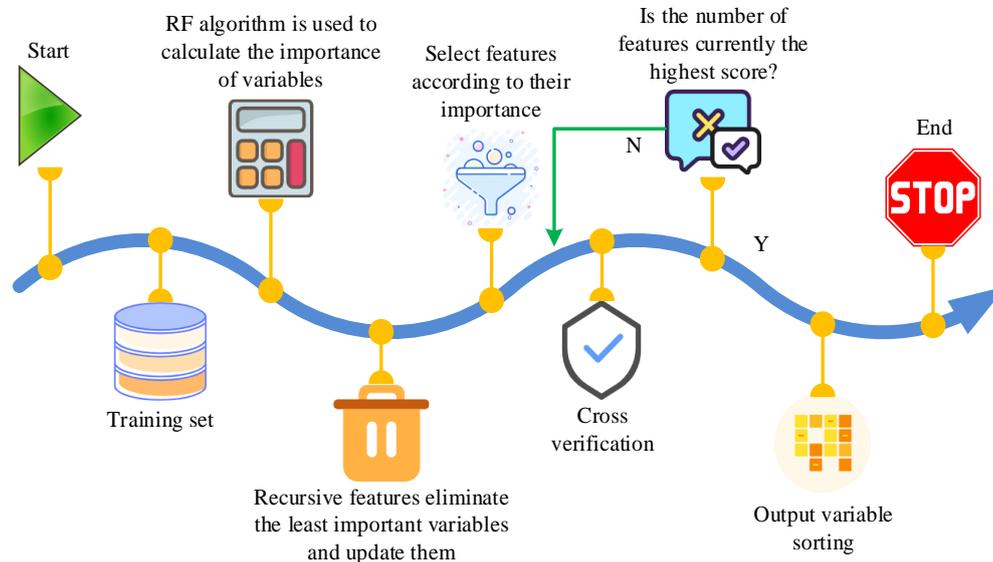


Figure 4: IRF process diagram.

selected from the initial variables, and the optimal node is selected from k_{ry} for each split of the tree [23-25]. In the growth process of the forest, each tree does not need pruning and will continue to split. By repeating the above operation q times, RF $f = (a_1, a_2, \dots, a_q)$ can be obtained. For each decision tree, the classification accuracy P_m is calculated using OOB data G_m^{OOB} . In the training set, each initial variable is denoted as α_b , and random noise is introduced into the α_b of G_m^{OOB} to obtain new data \hat{G}_m^{OOB} . The $\hat{\alpha}_b$ of each a_m for the hidden \hat{G}_m^{OOB} is calculated. Finally, the importance ID of α_b is calculated, as shown in equation (5).

$$ID = \frac{1}{q} \sum_{j=1}^q (\alpha_b - \hat{\alpha}_b) \tag{5}$$

On the basis of the above content, the RF algorithm is displayed in Figure 3.

After the RF part is completed, ID can be used as a basis to select each number of features one by one, and then cross validate the selected feature set. Finally, the number of features with the highest average score is obtained and determined [26-27]. Combining the RF part and recursive feature elimination, a complete schematic diagram of the IRF process can be obtained, as shown in Figure 4.

In Figure 4, the ID of α_b is first calculated and sorted, and then the recursive feature elimination method is used to extract features while updating the feature set. The above steps are repeated until all features are traversed and the feature set with the best accuracy is selected. Then, the cross-validation method is used to

select the highest feature score set. Finally, the feature variable set and sorting are obtained.

3.3 Highway traffic incident detection algorithm based on LSTM optimized by bayesian optimization

After the feature variable selection is completed, the HDIT algorithm can be constructed based on the grading model. Usually, recurrent neural networks are applied in short-term information prediction, but they cannot meet the requirements of higher accuracy prediction. However, LSTM can solve the gradient vanishing in the above neural networks and perform long-term learning of relevant information tasks. The structure of LSTM only adds cell states on the basis of recurrent neural networks, which are used to store previously learned information and sequences, and achieve information exchange through special forms. This structure can increase memory implementation, so it can present good results in many problems [28-30]. The core part of LSTM is the neuron state. The information in the neuron state is controlled through a gating mechanism, which mainly includes four parts: forget gate, input gate, update gate, and output gate. In the forget gate, it is implemented through the Sigmoid layer, as expressed in equation (6).

$$F_t = \sigma [W_F \square (h_{t-1}, x_t) + p_F] \tag{6}$$

In equation (6), F_t is the output value of the forget gate. σ , W_F , h_{t-1} , and p_F correspond to the activation function, weight value, output of the previous neuron, and bias value, respectively. In the input gate, candidate vectors are generated through the tanh layer. The required finer values are determined by the sigmoid function. The

output I_t and update content \tilde{C}_t of the input gate are calculated using equation (7).

$$\begin{cases} I_t = \sigma[W_I \square(h_{t-1}, x_t) + p_I] \\ \tilde{C}_t = \tanh[W_C \square(h_{t-1}, x_t) + p_C] \end{cases} \quad (7)$$

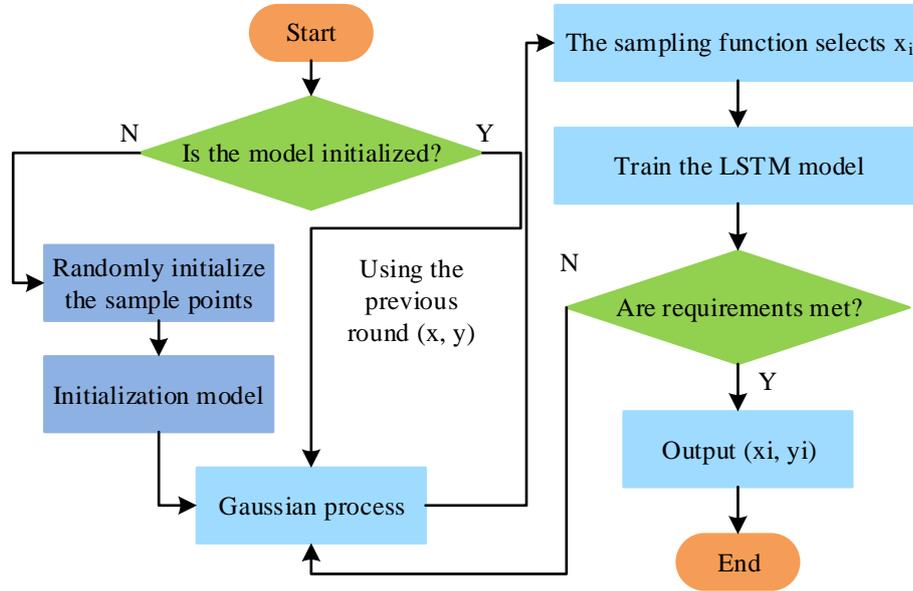


Figure 5: The process of optimizing LSTM hyper-parameters based on BO.

In equation (7), tanh is the tanh function. W_I and W_C correspond to the weights of input gates and cell states, respectively. p_I and p_C are the bias values of the input gate and cell state. In the update gate, it is necessary to update the cell state, as calculated in equation (8).

$$C_t = F_t \square C_{t-1} + I_t \square \tilde{C}_t \quad (8)$$

In the output gate, the expression for the output result O_t is shown in equation (9).

$$O_t = \sigma[W_O \square(h_{t-1}, x_t) + p_O] \quad (9)$$

In equation (9), W_O and p_O are the weight and bias values of the output gate, respectively. The final output result h_t is calculated using equation (10).

$$h_t = O_t \square \tanh(C_t) \quad (10)$$

LSTM has significant advantages in processing long time series tasks. It is more in line with actual situations, which facilitates subsequent classification tasks. All hidden units in the last layer are output, and then linked to the fully connected layer to complete binary classification. Due to the influence of hyper-parameters on the model performance, it is crucial to determine the appropriate combination of hyper-parameters. The BO algorithm is a hyper-parameter optimization method that can intelligently select the next evaluation point based on historical observation results, achieving parameter configuration close to the optimal solution in fewer iterations, and overcoming the time-consuming and unstable random results of other hyper-parameter setting methods. Therefore, it is used to optimize LSTM to

achieve better performance in highway traffic incident detection. The specific process of optimizing LSTM model with BO algorithm is as follows. Firstly, the range of hyper-parameters is set and initialized to obtain the corresponding hyper-parameter data set. A set of data is randomly selected for Gaussian process regression to establish a probability distribution function and fit the objective function. The prior distribution of the Gaussian process is updated by the loss value, and the surrogate model is modified. Then, the sampling function is applied to select the next optimal sample point, which is the point x_i to be evaluated. The above points are input into LSTM for training, which can obtain the new output value y_i of the objective function, update it to the sample set $Q = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, and update the lining model. Finally, the condition judgment is completed through the loss value. If the value satisfies the requirements, the loss value and the current optimal hyper-parameter combination can be output. Otherwise, Q is updated, and the iterative correction process can be continued until it meets the requirements. From this, the process of optimizing LSTM hyper-parameters based on BO can be obtained, as shown in Figure 5.

In Figure 5, the first step is to determine whether the model has completed initialization. If it has, the sampling function selection step can be entered. Otherwise, the initialization step can be entered. Next, the initial sample points are randomly selected and used to initialize the LSTM model. Then, a Gaussian process is applied to

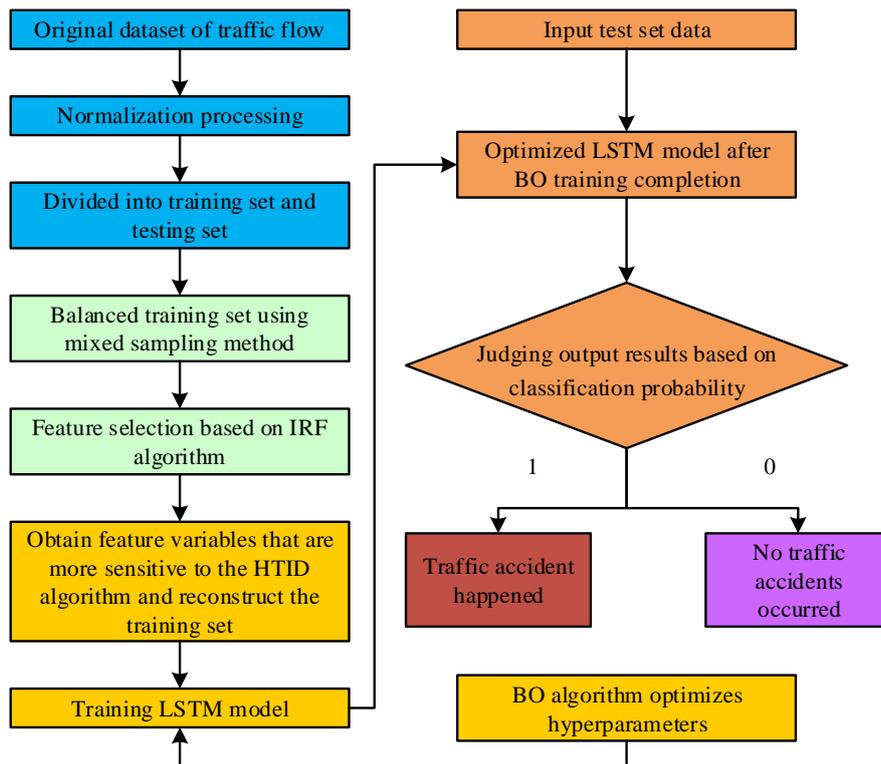


Figure 6: Schematic diagram of IRF-HTID-BO-LSTM method flow.

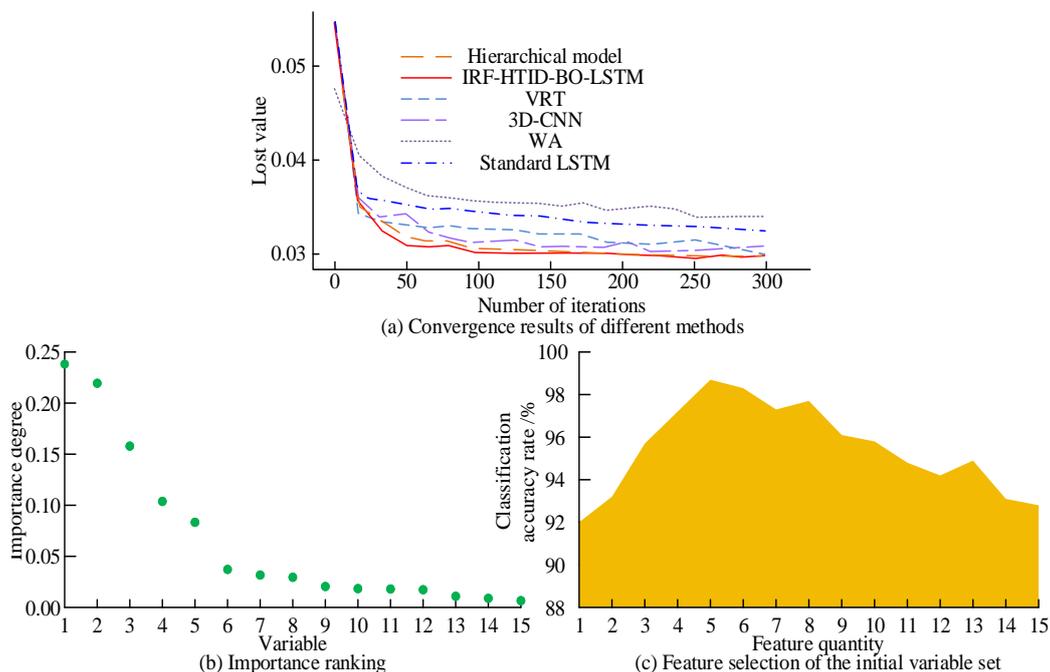


Figure 7: The convergence performance of different methods and the results of initial characteristic variable selection.

establish a surrogate model for predicting LSTM hyper-parameters. Finally, a new combination of hyper-parameters is selected through the sampling function and used to train the LSTM model. Finally, the performance is checked to check whether it meets the requirements. If it does not meet the requirements, the Gaussian process is returned. If it meets the requirements, the process can be ended. During the optimization process, the designed hyper-parameters include learning rate, batch size,

iteration count, number of hidden layer nodes, and time step size. The values of learning rate are [0.01, 0.001, and 0.0001], the range of batch size values is [32, 64, 128, 256, 512], and the settings of other hyper-parameters are based on past experience. During training, to balance the samples, a mixed sampling processing method is combined with IRF-based feature variable selection and BO-optimized LSTM to obtain the final IRF-HTID-BO-LSTM algorithm. The specific implementation process is

as follows. Firstly, the determined feature variables are applied to construct the training set and obtain the output matrix, as shown in equation (11).

$$OP = [y_i] = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{5400} \end{bmatrix} \quad (11)$$

In equation (11), $y_i \in \{0,1\}$ represents the label of the i -th input sample. 1 and 0 correspond to the event label and non-incident label, respectively. Based on the partitioning of the training set, the output corresponding to the first 2700 input samples is 1, and the remaining samples are 0. In the training and optimization part of LSTM, the input data set $\{x_s, x_s \in X_{t_set}, y_s \in Y_{t_set}\}$ is first determined. Then, the forget gate, input gate, update gate, and output gate are sequentially passed. After training each segment x_s through the highway traffic incident detection method, the final feature vector Y_s^{LSTM} output by LSTM can be obtained, as shown in equation (12).

$$Y_s^{LSTM} = L_{LSTM}(x_s; W_F, W_C, W_I, W_O) \quad (12)$$

In equation (12), $L(\square)$ represents the mapping function of LSTM. Because highway traffic incident

detection belongs to binary classification, the general loss function uses cross entropy, while the output category probability uses softmax function. The obtained probability $P_{mc}(x_s)$ and cross entropy expression are shown in equation (13).

$$\begin{cases} P_{mc}(x_s) = \text{soft max}(Y_s^{LSTM}) \\ \zeta(W, p) = -\left\{ \sum_{m=1}^n \{y_m \log[P_{mc}(x_s)] + (1-y_m) \log[1-P_{mc}(x_s)]\} \right\} \end{cases} \quad (13)$$

In equation (13), mc represents the classification label. $\zeta(W, p)$ and y_m are the objective function and the true label, respectively. The dataset is divided into training and testing sets at a ratio of 5:5. The LSTM model is set as $Z(x_s; W, p)$. The input dataset is $x_s \in X_{test}$. It is compared with the obtained classification prediction probability values to complete the classification based on the predicted label results. The flowchart of the IRF-HTID-BO-LSTM method can be obtained, as shown in Figure 6.

In Figure 6, the original traffic flow data set is normalized and partitioned. Then, the training set is balanced using a mixed sampling algorithm of Borderline SMOTE over-sampling and Tomek Links under-sampling. The IRF algorithm is used for feature selection to obtain

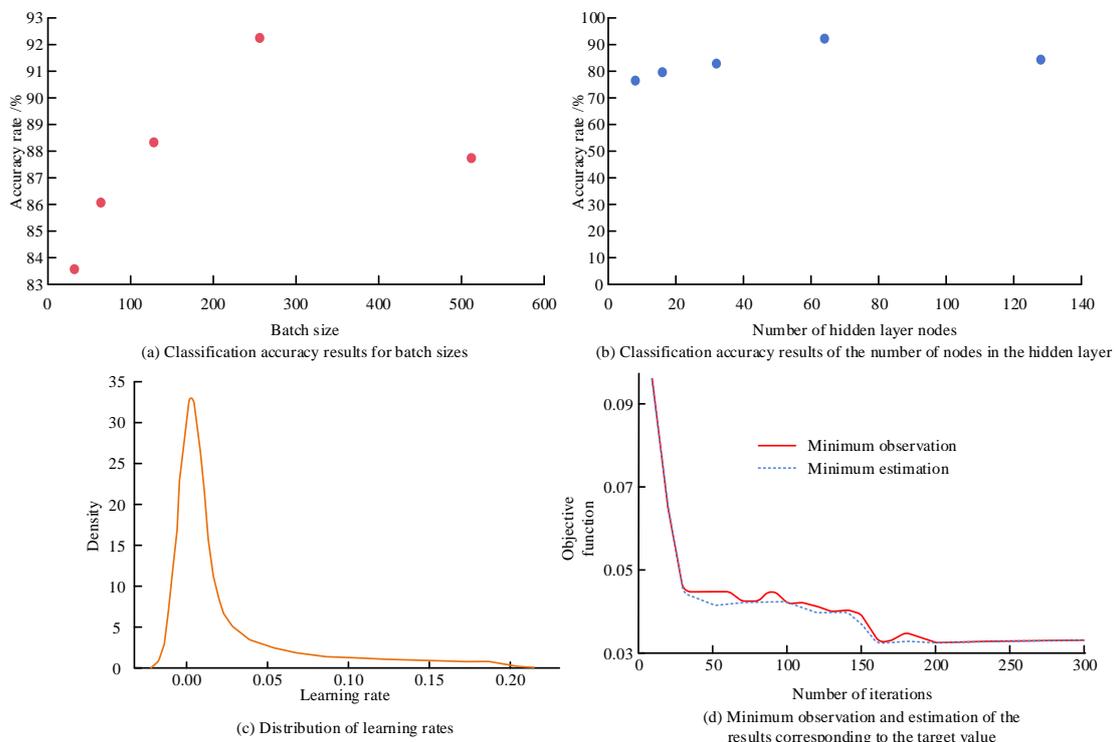


Figure 8: Training results of different hyper-parameters.

feature variables that are more sensitive to the HTID algorithm. The training set is reconstructed and input into the LSTM model for training. The BO algorithm optimizes the hyper-parameters of the LSTM, and the

output results can be determined by the classification probability.

4 Results

To test the performance and application effectiveness of the proposed IRF-HTID-BO-LSTM method, the study first examines the effectiveness of the method, and explores the correlation fit of the grading model index optimization. Then, the comparative experiments are conducted to scientifically analyze the performance. Finally, its effectiveness in practical applications is evaluated.

4.1 Performance analysis of grading model construction and traffic incident detection methods

To evaluate the performance of the research method, experiments are conducted using software Matlab. The data set used for processing is as follows. Firstly, a total of 100 sets of mountain road traffic incidents were obtained, with a total of 10000 data points. Among them, the traffic incident data set contains 2224 data points, while the rest are non traffic incident data points. Then, the traffic incident data set is divided into training and testing sets in a 1:1 ratio. Then, a mixed sampling algorithm of Borderline SMOTE oversampling and Tomek Links under-sampling is used to balance the two types of samples in the training set, resulting in 5600 samples. Finally, the samples are divided into training and testing sets in a 1:1 ratio. Afterwards, the IRF is used to screen the initial variable set constructed, obtain

feature variables that are more sensitive to the HTID algorithm, reconstruct the training set, and use it to train the LSTM network. The BO algorithm optimizer hyper-parameters are used. Finally, the output results can be used to determine whether a traffic incident has occurred on mountainous highways. In addition, the performance evaluation of the grading model is based on data collected from 487 curvature profiles of national highways, which are used for correlation degree calculation. To verify the performance of research methods more scientifically, the current mainstream methods are introduced for comparison, namely the standard LSTM method, the traffic incident detection method based on Video Recognition Technology (VRT), the detection method based on Three-Dimensional Convolutional Neural Networks (3D-CNN), and the detection method based on Wavelet Analysis (WA). In addition, commonly used indicators are used for evaluation, namely Detection Rate (DR), Mean Detection Time (MDT), False Alarm Rate (FAR), Comprehensive Performance Index (CPI), Receiver Operating Characteristic (ROC), and Area Under the Curve (AUC). MDT is the average time required to obtain test results, which is calculated by the average time required for all test times. FAR refers to the proportion of normal cases incorrectly reported as abnormal cases within a certain period of time. The lower the value, the better the specificity of the method. It is usually calculated by actually calculating the proportion of samples that are misclassified as positive. CPI is an indicator that comprehensively reflects the performance

Table 3: Statistical results of grading model structure optimization.

Model	Regression statistics	/
After improvement	Multiple R	0.873491573
	R Square	0.623749182
	Adjusted R Square	0.471134821
	Standard error	62.12849208
	Observed value	487
Before improvement	Multiple R	0.821998
	R Square	0.675681
	Adjusted R Square	0.554061
	Standard error	2526.453
	Observed value	487

of the detection method. It takes into account DR, MDT and FAR three indicators, and is calculated by the product of MDT, FAR/100 and (1-DR/100). The study first tests the convergence performance of different methods and analyzes the results of the initial feature variables. The random feature variables and decision tree are set to 4 and 1000, respectively, and the results are shown in Figure 7.

Figures 7 (a) -7 (c) respectively correspond to the convergence results of different methods, the importance ranking of the initial variable set, and their selection results. From Figure 7, the proposed grading model and traffic incident detection method did not exceed 150 iterations on the training dataset to achieve a stable convergence state. However, the VRT method, 3D-CNN method, and WA method required 265, 273, and 284

iterations, respectively, to achieve a stable convergence state and corresponding higher loss values. In addition, the standard LSTM model requires 280 iterations to reach a stable state, because the baseline model has not been specifically optimized for traffic incident detection tasks in complex mountainous highway environments, and the input features contain a large amount of redundant or irrelevant information. LSTM requires more iterations to identify and utilize effective information. After running recursive feature elimination through cross validation, the number of features with better classification accuracy can be obtained. Then, feature variables with lower rankings can be deleted, and finally the set with the highest feature score can be selected. Similarly, the detector corresponds to the difference between the measured and predicted occupancy values, the ratio of the predicted and measured

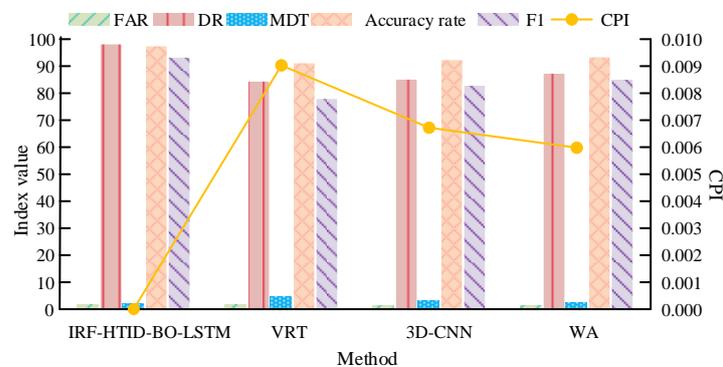
speed values, and the difference between the measured and predicted speed values. Figure 7 (b) showed that when the number of features was 5, the corresponding classification accuracy was the highest at 98.7%. Based on the importance ranking of the feature variables, the final set of feature variables can be determined as the ratio of the measured occupancy values corresponding to the upstream and downstream detectors, the product of the difference in occupancy values and the speed difference. Similarly, the difference between the measured occupancy values and the predicted values corresponding to the detectors, the ratio of the predicted speed values and the measured speed values, and the difference between the measured speed values and the predicted values can also be determined. This showed the above features contain more information about changes in traffic flow status, and can help the model better

distinguish traffic incidents and run more stably. In addition, the higher importance is due to its stronger interaction with other features. When the ratio of the measured occupancy values corresponding to the upstream and downstream detectors is the highest, it can reflect the changes in traffic flow between different detectors and is a sensitive indicator of changes in traffic flow status. The product of occupancy rate difference and speed difference contains information from both dimensions of speed and occupancy rate, which can more comprehensively describe the dynamic changes of traffic flow. To assess the effectiveness of the grading model structure optimization, the regression statistical method is used for evaluation, as displayed in Table 3.

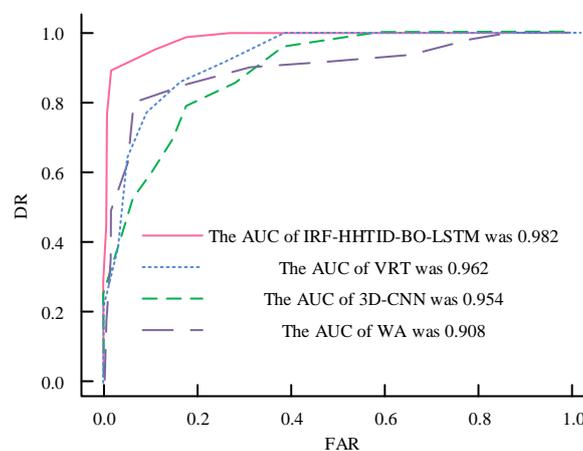
According to Table 3, the fitting degree of the grading model before optimization was only 82.2%, indicating a strong positive correlation. The R-Square

Table 4: The statistical results of the correlation degree of the optimized grading indicator of mountainous highway flat curve.

/	/	Coefficients	Standard error	t Stat	P	Lower 95%	Upper 95%
/	Intercept	0.055430518	0.098029	0.565449	0.585587	-0.16633	0.277188
Flat curve index grading	$C_{-}L^4$	0.001274648	0.001808	0.705061	0.038612	-0.00282	0.005364
	$C_{-}L^3$	3.873047452	1.19E-05	3.264968	0.009761	1.19E-05	6.56E-05
	$C_{-}L^2$	15.04555171	2273.262	6.618486	0.000166	9803.4	20287.7
	$C_{-}L^1$	27.10997887	80.05872	0.338626	0.043608	-157.506	211.7257
	$C_{-}L^0$	11.06338754	18.46679	0.599097	0.025683	-31.5211	53.64787



(a) Detection performance of different methods



(b) The combined performance of the different methods

Figure 9: Comparison of detection performance and comprehensive performance of different methods.

value showed that the model explained 67.6% of the variability of the dependent variable, indicating that the model has good explanatory power for the data. The adjusted R-Square value was 47.1%, indicating that the model has a reasonable explanatory power. The fitting degree of the optimized grading model has been significantly improved, and the fitting effect of the highway traffic condition correlation model for flat curve structure grading was improved to 87.3%, indicating a strong linear relationship between the independent and dependent variables. To further evaluate the correlation between the optimized grading indicators for mountainous highway flat curves, the specific results are displayed in Table 4.

From Table 4, C_{L^2} and C_{L^3} had a significant impact on traffic conditions when the length of the flat curve was between 300m-500m. When the length of the flat curve exceeded 500m or was less than 200m, the remaining grades exhibited marginal significance, and the intercept term was not significant, making a small contribution to the model. The above results verify that different classifications have a significant impact on the routes and traffic conditions in mountainous areas, indicating that the grading model has a good application effect on indicator grading. The indicator grading of the flat curve length of mountainous highways is reasonable. To obtain the optimal hyper-parameters for the IRF-HTID-BO-LSTM method, the study focuses on optimizing the time step and the number of hidden layer

nodes using the BO algorithm, and updates the training parameters using the adaptive moment estimation algorithm (Adam). Due to the large data set size, the fixed learning rate is 0.001, the maximum number of iterations is 300, the batch sizes are 32, 64, 128, 256, and 512, the time step range is 1-30, and the hidden layer node range is 8, 16, 32, 64, and 128. In addition, the study uses 5 cross validations for training. The obtained training results for different hyper-parameters are shown in Figure 8.

Figures 8 (a) and 8 (b) show the classification accuracy results corresponding to batch size and the number of hidden layer nodes, respectively. Figure 8 (c) and 8 (d) correspond to the distribution of learning rates and the comparison of the minimum observation and estimation corresponding to target values. Figures 8 (a) and 8 (b) showed that the accuracy results exhibited a bell-shaped curve with the change of hyper-parameters. As the number of hidden layer nodes and batch size continue to increase, the classification accuracy shows a trend of first increasing and then decreasing. Moreover, the number of hidden layer nodes has a relatively large impact on the classification accuracy. When the number of hidden layer nodes and batch size were 64 and 256, respectively, the performance of the research method was optimal, indicating that the research method has stability and robustness. Figure 8 (c) shows the distribution of learning

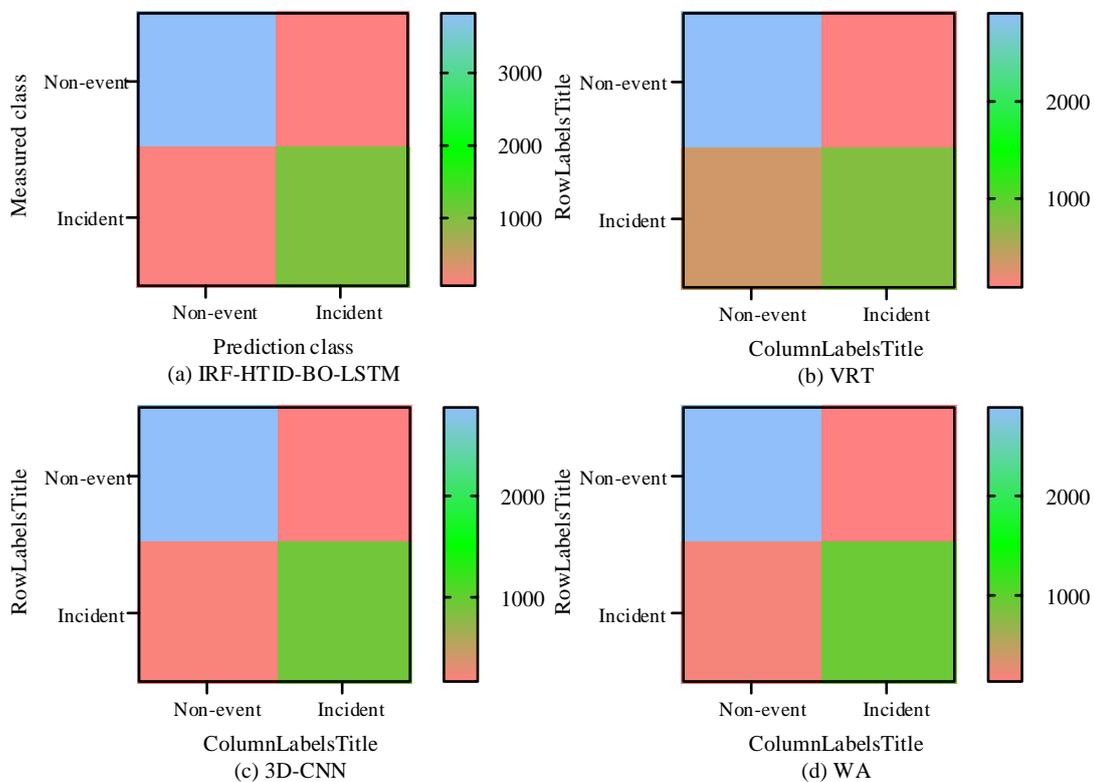


Figure 10: Comparison of classification performance results of different methods.

rates. The corresponding density reached its highest when the learning rate was 0.001. In Figure 8 (d), with the

gradual increase of iterations, the BO algorithm converged when the search reached 18 times, indicating

that the global optimal solution could be obtained. To further analyze the performance of the IRF-HTID-BO-LSTM method, experiments are conducted on different methods using DR, MDT, FAR, CPI, accuracy, F1 value, and AUC. The results are shown in Figure 9.

Figure 9 (a) and Figure 9 (b) respectively compare the detection performance and comprehensive performance results of different methods. From Figure 9, in the comparison results of detection performance, the research method had the best performance in all detection indicators except for FAR, which was 1.82%. DR, MDT, CPI, accuracy, and F1 value were 98.45%, 2.21%, 0.00159, 97.74%, and 93.52%, respectively. Compared with VRT, 3D-CNN, and WA, the detection performance

of the research method has been significantly improved. Although mainstream methods sacrifice DR values to obtain smaller FAR values, there are still many traffic incidents that have not been detected. The research method can improve other detection performance while ensuring a lower MDT. Among them, when processing data of the same scale, the WA method had the highest computational efficiency, corresponding to an MDT value of 2.6 seconds, while the research method had relatively high computational efficiency. This means that the research method can significantly improve the response capability to emergencies, thereby improving road safety and traffic flow efficiency. In the comprehensive performance results, the AUC of the

Table 5: Comparison of classification performance of different methods based on paired t-test.

Evaluating indicator	Method	df	T-score	Sig
Susceptibility	IRF-HTID-BO-LSTM	40	4.08	**
	VRT	40	5.38	**
	3D-CNN	40	7.46	**
	WA	40	8.18	**
	Standard LSTM	40	10.08	**
Specificity	IRF-HTID-BO-LSTM	40	16.30	**
	VRT	40	24.95	**
	3D-CNN	40	28.58	**
	WA	40	24.00	**
	Standard LSTM	40	28.02	**

Note: "** *" indicates $P < 0.01$.

Table 6: Comparison of model accuracy results based on K-fold cross validation.

Method	Sample size	Correlation coefficient	Root mean square error	Coefficient of variation
IRF-HTID-BO-LSTM	2000	0.65	268	0.38
VRT	2000	0.30	375	0.51
3D-CNN	2000	0.29	305	0.49
WA	2000	0.31	229	0.27
Standard LSTM	2000	0.34	243	0.32

research method was the highest, at 0.982, while the AUC values of the VRT method, 3D-CNN method, and WA method corresponded to 0.962, 0.954, and 0.908, respectively. The proposed feature variable selection method can effectively improve the performance of the final detection method, and the overall performance of the research method is also the best. The comparison of classification performance results of different methods is shown in Figure 10.

Figures 10 (a) -10 (d) show the confusion matrix results of the IRF-HTID-BO-LSTM method, VRT method, 3D-CNN method, and WA method, respectively. In Figure 10, the research method had the largest number of positive samples, which meant that the method identified the largest number of traffic incident samples. The WA method had the smallest positive class positive samples, indicating that its classification performance was the worst among mainstream methods. In summary, many indicators of the research method are superior to other mainstream algorithms, and the comprehensive performance and classification performance are the best,

indicating that the feature selection and grading model can effectively improve the performance of the method. To more accurately quantify the performance effects of different methods, statistical tests are introduced to explore the differences between each method. Paired t-tests are used to analyze the classification effects of different methods, and the results are shown in Table 5.

From Table 5, the research method had extremely significant statistical significance compared to other benchmark models on sensitivity and specificity. The research method has a certain degree of reliability and lays a good foundation for subsequent practical applications. To test the robustness of the data set used in the study, k-fold cross validation is used to analyze different methods. The specific process is as follows. The data set is randomly divided into K equally sized subsets to ensure that each subset is as similar as possible in data distribution. Then, one subset is selected as the test set, and the remaining subsets are merged as the training set. Finally, the model is trained on the training set, and its performance is evaluated on the test set. The accuracy of

the study is analyzed, and the results are shown in Table 6.

Table 6 shows that the accuracy performance from high to low is IRF-HTID-BO-LSTM, Standard LSTM, WA, VRT, and 3D-CNN. The accuracy of the research method is higher than other methods, which also indicates that the data set used in the study has stability and reliability.

4.2 Application effect analysis

To explore the effectiveness of the research method in practical applications, this study takes the traffic conditions of ordinary provincial highways as an example. The annual report data provided by Ningxia Highway Management Center in 2022 is imported into the research method. The graph analysis function in geographic information system is used to identify and mark key or abnormal road sections that exceed the

indicator grading and cause abnormal traffic volume. Taking the Haitian Highway (western section of Ningxia) as an example, the key or abnormal point road section annotations obtained are shown in Figure 11.

From Figure 11, in the practical application of the research method, there were nine abnormal points detected in the Haitian line section, and the affected range of the abnormal point road section was relatively long. In addition, when the flat curve was too short, traffic conditions were more affected by the structure of the flat curve. In practical applications, the research method can be combined with the relevant functions of geographic information systems to visually display and label abnormal sections of mountainous highways, and timely feedback to relevant departments for management. The specific data results of the abnormal points in the Haitian section are displayed in Table 7.

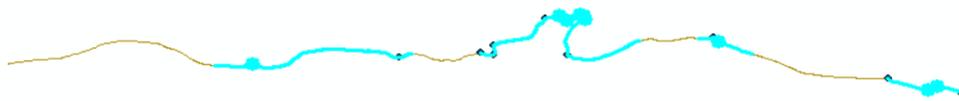


Figure 11: Mark and display map of key or abnormal points of Haitian highway (West Ningxia Section).

Table 7: The specific data results of the anomaly section of the Haitian section.

Road section name	Section number	Section length of flat curve /m	Standardized residuals
Shapotou Unity Bridge	G338640521	10.59538775	2.952860532
Gutang village, Zhongning	G338640521	76.06553674	-2.620215777
Zhongning East Hedong	G338640521	31.59364867	-2.554101242
Shapotou tourism new town	G338640502	6.762924746	-2.488856025
Shapotou tourist area	G338640502	79.81748426	2.441351606
Shapotou Yangtan east	G338640502	81.97640697	2.340711859
Kantang	G338640502	283.6514657	2.306630542
Shapo head meets the water	G338640502	173.6766195	2.305995462
Intersection of Yingshui Railway Station	G338640502	9.164988211	-2.245316332
Mengjia Bay East	G338640502	90.77954661	2.234417005

According to Table 7, in practical applications, specific road segment numbers and length information of flat curved road segments where traffic incidents occur can be obtained through the research method. The standard error range was within ± 3 . The above results suggest due to the fact that the VRT method and 3D-CNN method may miss the traffic incident in some cases, while the WA method has significant problems in feature extraction or classification decision-making, resulting in the worst performance in traffic incident recognition. In summary, the research method still demonstrates good

performance in practical applications, providing intelligent detection ideas for traffic incidents on mountainous highways and ensuring accurate detection of abnormal road sections even under complex road conditions.

5 Discussion

Mountainous highway terrain is complex, climate change, easy to appear traffic accidents. Intelligent traffic detection method can be based on real-time road

conditions, and timely detection and early warning of potential safety hazards, improving driving safety. Therefore, an index optimization method and a highway curve classification method for mountainous highway flat curve structure are proposed in this paper. The IRF algorithm is used to select characteristic variables. BO algorithm and LSTM are combined to obtain IRF-HHTID-BO-LSTM algorithm.

The research results show that when the number of features is 5, the classification accuracy of the research method is the highest 98.7%. According to the importance ranking of the feature variables, the final feature variable set can be determined as the product of the ratio of the measured value of the share corresponding to the upstream and downstream detectors, the difference of the share and the velocity difference. The same detector corresponds to the difference between the measured value of occupancy and the predicted value, the ratio of the predicted speed and the measured speed and the difference between the measured speed and the predicted value. Among them, the ratio of the measured occupancy corresponding to the upstream and downstream detectors is the highest because it can reflect the change of traffic flow between different detectors. It is a sensitive indicator of the change of traffic flow state. The product of occupancy rate difference and speed difference contains information from both dimensions of speed and occupancy rate, which can more comprehensively describe the dynamic changes of traffic flow. The above results are generated because the research method combines BO algorithm and LSTM algorithm to adjust hyper-parameters and process time series data. Therefore, the research method can learn the patterns in the data more effectively, while other mainstream algorithms lack the corresponding complexity when processing complex mountainous highway traffic incident data.

The hyper-parameter sensitivity analysis results show that the accuracy results present a bell curve with the change of hyper-parameters. With the continuous increase of the node in the hidden layer and the batch size, the classification accuracy shows an initial increase followed by a decrease. Moreover, the number of nodes in the hidden layer has a relatively large impact on the classification accuracy. The performance of the research method is the best, which indicates that the research method has stability and robustness. The optimal hyper-parameter combination is obtained as follows. The time step, batch size and number of hidden layer nodes are 5, 64 and 256, respectively. The comparison of the detection performance and comprehensive performance results of different methods shows that the performance of other detection indicators is the best except for the FAR index of 1.82%. and the DR, MDT, CPI, accuracy and F1 are 98.45%, 2.21%, 0.00159, 97.74% and 93.52%, respectively. In addition, its detection performance is obviously better than other benchmark models, but the computational efficiency of the research method is relatively high, and the corresponding MDT value is 3.1s. The research can obtain a smaller FAR by sacrificing certain DR and computational efficiency. In addition, the

best performing method in related work is the intelligent transportation service framework integrating heterogeneous data sources, whose accuracy rate is only 90%, because it is difficult to accurately analyze the road structure and traffic safety status of mountainous highways. Its subsequent management cannot get timely feedback. The comprehensive performance results show that the AUC values of the research method, VRT method, 3D-CNN method and WA method correspond to 0.982, 0.962, 0.954 and 0.908, respectively. The higher AUC value of the research method indicates because it can effectively remove redundant information through the improved feature variable selection method, improving the overall effect of the detection method.

In summary, the research method can still maintain a high comprehensive performance when facing complex mountainous highways. The effect is also excellent in practical application, but it sacrifices a certain degree of computational efficiency and has wrong classification. This may be due to the complexity and variability of mountain road traffic data. Therefore, in future, a more suitable model architecture for mountain road design can be selected according to the task characteristics, and the generalization ability of the model can be improved through data enhancement technology.

6 Conclusion

To reduce traffic accidents on mountainous highways and further enhance the ability of highways to maintain smooth traffic, an index optimization method for grading the structure of highway flat curves and a grading method for highway curve levels were designed. The IRF-HHTID-BO-LSTM method was proposed. The proposed grading model and traffic incident detection method achieved stable convergence with lower loss values on the training data set without exceeding 150 iterations. The optimal hyper-parameter combination optimized by the BO algorithm was as follows. The time step, batch size, and number of hidden layer nodes were 5, 64, and 256, respectively. The research method only had a high FAR value of 1.82%, while the performance of other detection indicators was the best. DR, MDT, CPI, accuracy, and F1 value were 98.45%, 2.21%, 0.00159, 97.74%, and 93.52%, respectively. The highest AUC value was 0.982. In practical applications, the research method accurately obtained the specific section numbers and length information of flat curved road segments where traffic incidents occurred, and the standard error range was within ± 3 . In summary, the research method can effectively classify and segment the curve line indicators of highway structures, establish the curve line grading indicator of mountainous highway structures, and ensure the rapid and accurate detection of traffic incidents on mountainous highways, providing a certain reference for the detection of other traffic incidents. However, there are still shortcomings in the research. For the analysis of highway structure, certain parameters may not be suitable for this study. For example, the results indicate that the speed variable has not reached the level of statistical significance. Therefore, in future research, the model can

be improved by introducing additional variables such as traffic intensity and lane occupancy rate.

7 Funding statement

This project was funded by National Key Research and Development Program of China (2023YFB4302701).

References

- [1] A. Haydari, and Y. Yılmaz. “Deep reinforcement learning for intelligent transportation systems: A survey,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 1, pp. 11-32, 2022, <https://doi.org/10.1109/TITS.2020.3008612>.
- [2] H. Zhao and A. Sharma. “Logistics distribution route optimization based on improved particle swarm optimization,” *Informatica*, vol. 47, no. 2, pp. 243-251, 2023, <https://doi.org/10.31449/inf.v47i2.4011>.
- [3] S. Murrar, F. M. Alhaj, and M. Qutqut. “Machine learning algorithms for transportation mode prediction: A comparative analysis,” *Informatica*, vol. 48, no. 6, pp. 117-130, 2024, <https://doi.org/10.31449/inf.v48i6.5234>.
- [4] C. Li and Z. Mu. “Analysis platform of rail transit vehicle signal system based on data mining,” *Informatica*, vol. 47, no. 3, pp. 441-450, 2023, <https://doi.org/10.31449/inf.v47i3.3942>.
- [5] J. Hawkins, and K. N. Habib. “A multi-source data fusion framework for joint population, expenditure, and time use synthesis,” *Transportation*, vol. 50, no. 4, pp. 1323-1346, 2023, <https://doi.org/10.1007/s11116-022-10279-8>.
- [6] H. Wang, L. Liao, W. Yi, and L. Zhen. “Transportation scheduling for modules used in modular integrated construction,” *International Journal of Production Research*, vol. 62, no. 11, pp. 3918-3931, 2024, <https://doi.org/10.1080/00207543.2023.2251602>.
- [7] C. Zhao, X. Chang, T. Xie, H. Fujita, and J. Wu. “Unsupervised anomaly detection-based method of risk evaluation for road traffic accident,” *Applied Intelligence*, vol. 53, no. 1, pp. 369-384, 2023, <https://doi.org/10.1007/s10489-022-03501-8>.
- [8] S. Cafiso, C. D'Agostino, and M. Kiec. “Investigating safety performance of the SAFESTAR system for route-based curve treatment,” *Reliability Engineering and System Safety*, vol. 188, pp. 125-132, 2019, <https://doi.org/10.1016/j.ress.2019.03.028>.
- [9] G. Ashley, O. A. Osman, S. Ishak, and J. Codjoe. “Investigating effect of driver-, vehicle-, and road-related factors on location-specific crashes with naturalistic driving data,” *Transportation Research Record*, vol. 2673, no. 6, pp. 46-56, 2019, <https://doi.org/10.1177/0361198119844461>.
- [10] M. R. Fatmi, and M. A. Habib. “Modeling vehicle collision injury severity involving distracted driving: Assessing the effects of land use and built environment,” *Transportation Research Record*, vol. 2673, no. 7, pp. 181-191, 2019, <https://doi.org/10.1177/0361198119849060>.
- [11] D. E. Monyo, H. J. Haule, A. E. Kitali, and T. Sando. “Are older drivers safe on interchanges? Analyzing driving errors causing crashes,” *Transportation Research Record*, vol. 2675, no. 12, pp. 635-649, 2021, <https://doi.org/10.1177/03611981211031232>.
- [12] S. B. Li, T. Sun, D. N. Cao, and L. Zhang. “Incident detection method of expressway based on traffic flow simulation model,” *Communications in Theoretical Physics*, vol. 71, no. 4, pp. 468-474, 2019, <https://doi.org/10.1088/0253-6102/71/4/468>.
- [13] X. Zhang, R. R. Souleyrette, E. Green, T. Wang, M. Chen, and P. Ross. “Collection, analysis, and reporting of kentucky traffic incident management performance,” *Transportation Research Record*, vol. 2675, no. 9, pp. 167-181, 2021, <https://doi.org/10.1177/03611981211001077>.
- [14] P. H. L. Rettore, B. P. Santos, R. R. F. Lopes, G. Maia, L. A. Villas, and A. A. F. Loureiro. “Road data enrichment framework based on heterogeneous data fusion for ITS,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 4, pp. 1751-1766, 2020, <https://doi.org/10.1109/TITS.2020.2971111>.
- [15] M. Won. “Outlier analysis to improve the performance of an incident duration estimation and incident management system,” *Transportation Research Record*, vol. 2674, no. 5, pp. 486-497, 2020, <https://doi.org/10.1177/0361198120916472>.
- [16] M. Azari, A. Paydar, FB. Eizizadeh, and V. G. Hasanlou. “A GIS-based approach for accident hotspots mapping in mountain roads using seasonal and geometric indicators,” *Applied Geomatics*, vol. 15, no. 1, pp. 127-139, 2023, <https://doi.org/10.1007/s12518-023-00490-2>.
- [17] C. Y. Lin, Y. C. Lai, S. W. Wu, F. C. Mo, and C. Y. Lin. “Assessment of potential sediment disasters and resilience management of mountain roads using environmental indicators,” *Natural Hazards*, vol. 111, no. 2, pp. 1951-1975, 2022, <https://doi.org/10.1007/s11069-021-05126-5>.
- [18] H. Duan, and Y. Song. “Grey prediction model based on Euler equations and its application in highway short-term traffic flow,” *Nonlinear Dynamics*, vol. 112, no. 12, pp. 10191-10214, 2024, <https://doi.org/10.1007/s11071-024-09611-x>.
- [19] S. T. Zheng, R. Jiang, B. Jia, J. Tian, and Z. Gao. “Impact of stochasticity on traffic flow dynamics in macroscopic continuum models,” *Transportation Research Record*, vol. 2674, no. 10, pp. 690-704, 2020, <https://doi.org/10.1177/0361198120937704>.
- [20] Y. Liu, C. Lyu, X. Liu, and Z. Liu, “Automatic feature engineering for bus passenger flow prediction based on modular convolutional neural network,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 4, pp. 2349-2358, 2021, <https://doi.org/10.1109/TITS.2020.3004254>.

- [21] A. A. Chaudhari, K. K. Srinivasan, B. R. Chilukuri, M. Treiber, and O. Okhrin. “Calibrating Wiedemann-99 model parameters to trajectory data of mixed vehicular traffic,” *Transportation Research Record*, vol. 2676, no. 1, pp. 718-735, 2022, <https://doi.org/10.1177/03611981211037543>.
- [22] T. V. Geetha, A. J. Deepa, and M. M. Linda. “Deep learning method for efficient cloud IDS utilizing combined behavior and flow-based features,” *Applied Intelligence*, vol. 54, no. 8, pp. 6738-6759, 2024, <https://doi.org/10.1007/s10489-024-05505-y>.
- [23] M. Hasanvand, M. Nooshyar, E. Moharamkhani, and A. Selyari. “Machine learning methodology for identifying vehicles using image processing,” *Artificial Intelligence and Applications*, vol. 1, no. 3, pp. 170-178, 2023, <https://doi.org/10.47852/bonviewAIA3202833>.
- [24] H. Mokayed, T. Z. Quan, L. Alkhaled, and V. Sivakumar. “Real-time human detection and counting system using deep learning computer vision techniques,” *Artificial Intelligence and Applications*, vol. 1, no. 4, pp. 221-229, 2023, <https://doi.org/10.47852/bonviewAIA2202391>.
- [25] J. Purohit, and R. Dave. “Leveraging deep learning techniques to obtain efficacious segmentation results,” *Archives of Advanced Engineering Science*, vol. 1, no. 1, pp. 11-26, 2023, <https://doi.org/10.47852/bonviewAAES32021220>.
- [26] I. Salman and J. Vomlel. “Learning the structure of bayesian networks from incomplete data using a mixture model,” *Informatica*, vol. 47, no. 1, pp. 81-94, 2023, <https://doi.org/10.31449/inf.v47i1.4497>.
- [27] G. S. Ohannesian and E. J. Harfash. “Epileptic seizures detection from EEG recordings based on a hybrid system of gaussian mixture model and random forest classifier,” *Informatica*, vol. 46, no. 6, pp. 105-116, 2022, <https://doi.org/10.31449/inf.v46i6.4203>.
- [28] C. Sivanandam, V. M. Perumal, and J. Mohan. “A novel light GBM-optimized long short-term memory for enhancing quality and security in web service recommendation system,” *Journal of Supercomputing*, vol. 80, no. 2, pp. 2428-2460, 2024, <https://doi.org/10.1007/s11227-023-05552-1>.
- [29] J. Wu, J. Tang, M. Zhang, J. Di, L. Hu, X. Wu, G. Liu, and J. Zhao. “PredictionNet: A long short-term memory-based attention network for atmospheric turbulence prediction in adaptive optics,” *Applied Optics*, vol. 61, no. 13, pp. 3687-3694, 2022, <https://doi.org/10.1364/AO.453929>.
- [30] W. Yang, W. Chang, Z. Song, F. Niu, X. Wang, and Y. Zhang. “Denoising odontocete echolocation clicks using a hybrid model with convolutional neural network and long short-term memory network,” *Journal of the Acoustical Society of America*, vol. 154, no. 2, pp. 938-947, 2023, <https://doi.org/10.1121/10.0020560>.

Cloud-Computing-Enabled Transformer Architecture for the Design of Functional Clothing Structures

Ailing Gou

Luohe Vocational Technology College, Luohe 462002, Henan, China

E-mail: ailing_gou@hotmail.com

*Corresponding author

Keywords: internet, cloud computing, functional clothing, structural graphics

Received: July 23, 2024

This paper introduces a graphical design model for smart clothing structures based on cloud computing and an integrated approach combining Transformer architecture with conditional Generative Adversarial Networks (cGANs). The model aims to revolutionize the functional clothing design industry by transforming users' diverse needs into machine-understandable vector representations using a multi-head self-attention mechanism. Subsequently, a decoder generates design elements, which are then visualized using cGAN techniques. To evaluate the model's performance, we conducted extensive computational experiments using a comprehensive dataset that includes various design styles and occupational categories, such as medical, catering, aviation, and industrial clothing. The model was trained and validated using K-fold cross-validation, ensuring robustness and generalizability. Key performance metrics were assessed, including design element similarity, layout rationality, and personalization accuracy. Experimental results show that the model achieves an average design element similarity score of over 89%, a layout rationality score of over 90%, and a personalization accuracy of nearly 92%. These performance indicators demonstrate the model's effectiveness in design accuracy, efficiency, personalization, and market adaptability, particularly for occupational clothing design in healthcare, catering, aviation, and industrial applications. The integration of Transformer and cGAN technologies significantly enhances the model's capability to generate high-quality, personalized designs while maintaining robustness and scalability. This approach provides a comprehensive solution for automating the design process, leading to improved design outcomes and enhanced user satisfaction.

Povzetek: Članek obravnava arhitekturo transformatorja, ki jo omogoča računalništvo v oblaku, za oblikovanje funkcionalnih struktur oblačil. Model, ki temelji na kombinaciji transformatorja in cGAN, pretvarja potrebe uporabnikov v vektorske predstavitve in generira oblikovalske elemente.

1 Introduction

Functional apparel, also known as occupational clothing or uniforms, are garments designed for specific occupations or work environments, which are designed not only for aesthetics and comfort, but also, more importantly, for their functionality, safety, and for brand image. Functional clothing plays a crucial role in various industries, they are not only the embodiment of dress code, but also the sign of professional identity, the booster of work efficiency, and the carrier of corporate culture [1]. In the medical industry, the white coats and operating room-specific clothing worn by healthcare workers not only create a professional image, but also have the hygienic functions of disinfection and anti-bacteria, protecting both doctors and patients from the risk of infection. In the restaurant industry, the uniforms of chefs and waiters not only keep neat and clean to meet food safety standards, but also convey the brand style of the restaurant through color and design. In the aviation and hospitality industries, the uniforms of cabin crew and receptionists are a direct reflection of the corporate image, and they convey the concept of professional and reliable service through a

uniform visual language. In the industrial field, especially in the chemical and construction industries, the protective performance of functional clothing is crucial. The application of special materials such as anti-static, fireproof, and anti-radiation provides the necessary safety for workers and reduces occupational hazards [2].

Functional apparel is equally indispensable in a variety of industries such as education, retail, and transportation, where they help to differentiate between different positions and promote teamwork while enhancing customer or public trust. The design of functional clothing needs to take into account the nature of the work, environmental factors, corporate culture, and ergonomic principles to ensure that the wearer can perform his or her job comfortably while reflecting the professionalism of the occupation and the consistency of the company [3].

Traditional functional apparel design is often limited by physical resources and geographic distances, and the design process usually requires frequent physical exchanges between designers, material suppliers, and manufacturers, which is not only time-consuming and labor-intensive, but can also lead to slow design iterations.

In addition, due to the lack of effective data management and analysis tools, design decisions often rely on experience and intuition, making it difficult to accurately capture market trends and user needs. This design approach also limits remote team collaboration and reduces design flexibility and responsiveness [4].

The advent of cloud computing has revolutionized the landscape of functional apparel design. By migrating design tools, resource libraries, and collaboration platforms to the cloud, cloud computing breaks down geographic boundaries and enables instant sharing of resources and remote team collaboration on a global scale. Designers can access the latest design software from any location, utilize cloud storage for file backup and version control, and greatly improve design efficiency [5]. The aim of this study is to explore how cloud computing can empower the structural graphic design of functional apparel, and by analyzing the application of cloud computing technology in the design process, we hope to reveal its positive impact on design efficiency, collaborative work, and innovativeness [6].

The innovation of this paper focuses on proposing a graphical design model of clothing structure based on cloud computing and Transformer architecture, which realizes efficient and precise understanding and response to diversified user needs through deep learning technology. The research mainly includes: (1) Dynamic adaptability and personalized design: relying on the elastic resources of cloud computing, the model is able to adjust in real time to respond to different user needs, ensure the uniqueness and personalization of the design, and satisfy the precise requirements for details in the design of functional clothing. (2) Multi-head self-attention mechanism: the introduced multi-head self-attention mechanism enhances the model's ability to capture the complex relationships among parts in the input sequence, and even if these parts are far away from each other in the sequence, they can still be correctly associated, thus enhancing the innovation and functionality of the design [7].

2 Literature review

2.1 Overview of the development of functional clothing

Functional apparel, i.e., functional clothing, is designed to meet the needs of the wearer in specific environments, whether it is protection from extreme climatic conditions or safety and convenience for specific occupational activities. From the initial waterproof, windproof, and breathable to the modern intelligent sensing and self-regulation, the development of functional apparel has demonstrated the deep integration of science and technology with textile innovation [8]. The development of this field has not only been driven by advances in materials science, but has also benefited from the results of ergonomic, biomechanical, and environmental adaptation research [9].

The origins of functional clothing can be traced back to the early 20th century, when the properties of natural

fibers were explored to create more durable and protective clothing [10]. However, the real revolution in functional clothing occurred in the mid-20th century with the invention of synthetic fibers such as nylon and polyester, new materials that were not only lightweight and wearable, but also had some waterproofing and warmth properties [11]. Subsequently, the emergence of waterproof and breathable membranes, such as Gore-Tex, marked a new era for functional clothing [12]. In recent years, the development of functional apparel has focused more on the smartness and responsiveness of materials. For example, phase change materials (PCMs) are capable of absorbing or releasing heat according to changes in ambient temperature to maintain the stability of the human microenvironment [13]. In addition, the application of conductive fabrics and nanotechnology enables garments to integrate sensors for health monitoring, environmental sensing, and other functions [14]. These innovations not only enhance the utility of clothing, but also open the way for personalization and customization.

The functional apparel market continues to expand as consumers demand higher levels of health, safety, and comfort, as well as an increase in outdoor sports and professional work scenarios [15]. Especially after the epidemic, the focus on personal hygiene and protection has driven the use of antimicrobial and antiviral materials in apparel [16]. Meanwhile, sustainability has become a focus of industry attention, with green materials and circular economy models being increasingly introduced into functional apparel design [17].

The future development of functional clothing will focus more on user experience and human-computer interaction. The integration of wearable technologies will make clothing part of the Internet of Things, enabling data collection and intelligent feedback [18]. In addition, with advances in artificial intelligence and machine learning, personalized design and on-demand manufacturing will become the norm, satisfying consumers' pursuit of uniqueness and adaptability [19]. Eventually, functional clothing will become more than just a piece of clothing, but an intelligent interface that connects the body to the external world.

2.2 Application of cloud computing in the field of clothing design

Cloud computing has brought unprecedented changes to the apparel design industry with its superior data processing capability and flexible service model. From design to production to supply chain management, cloud computing technology is gradually penetrating and optimizing the entire apparel industry chain, creating more value for designers, producers and consumers [20].

During the design phase, cloud computing provides powerful and easily accessible computing resources that enable designers to perform complex design simulations and 3D renderings without relying on expensive local hardware facilities [21]. For example, platforms such as the Bock Intelligent Apparel Cloud CAD System utilize cloud computing technology to allow designers and production staff to design and manage work from any

location, at any time, greatly enhancing efficiency and flexibility (. The Smart Custom Apparel Cloud CAD system even integrates advanced design tools and intelligent algorithms to help designers rapidly iterate their designs while maintaining a high level of innovation and personalization [22]. Cloud computing also facilitates the digital transformation of the apparel production process by enabling supply chain transparency and collaboration through cloud platforms, effectively reducing inventory costs and shortening the time-to-market cycle [23]. Services provided by companies such as Zeta Cloud, which utilize cloud computing and meta-universe technologies, provide a new perspective on apparel design and production, making remote collaboration and virtual presentations possible, reducing the production of physical prototypes, and saving time and resources [24]. On the consumer side, cloud computing is able to accurately capture consumer preferences and market trends through big data analysis, providing customized products and services for apparel companies [25]. Applications such as virtual fitting rooms and personalized recommendation systems allow consumers to experience the real effect of clothing before purchase, improving customer satisfaction and loyalty [26].

design efficiency, enhancing user experience, and optimizing supply chain management. For example, [27] explored how cloud technology can accelerate the 3D modeling and simulation process of apparel by providing high-performance computing power, thus shortening the product development cycle. At the user level, on the other hand, research by [28] demonstrates how cloud-driven virtual fitting technology can transform the consumer shopping experience by reducing return rates and increasing online sales through accurate body size matching. Domestic studies have also followed the international pace and are dedicated to exploring how cloud computing can empower various aspects of apparel design and manufacturing. A study by [29] revealed the application of cloud computing in the apparel supply chain, which significantly reduced operational costs by predicting market demand and optimizing inventory management through big data analysis. In addition, the project focuses on the development of intelligent design software on the cloud computing platform, which utilizes machine learning algorithms and is able to automatically generate design solutions based on fashion trends and consumer feedback, which greatly enhances the innovation and efficiency of design. The summary table of research results is specifically shown in Table 1.

2.3 Current status of domestic and international research

Overseas, research on the application of cloud computing in apparel design is quite mature, focusing on improving

Table 1: Summary of research results.

Research/Method	Accuracy	Personalization	Design Efficiency	Key Technologies/Materials	Main Findings/Contributions	Gaps/Improvement Points
Gore-Tex Waterproof & Breathable Membrane	High	Low	Medium	ePTFE	Provides excellent water resistance and breathability	Lacks personalized design and intelligent regulation capabilities
Smart PCM Clothing	Moderate	Moderate	Low	Phase Change Materials	Can regulate the microenvironment according to ambient temperature	High production cost and long design cycle
Conductive Fabric Health Monitoring	High	Moderate	Moderate	Conductive Fibers, Sensors	Achieves real-time health data monitoring	Short battery life, poor wash durability
Zeta Cloud Virtual Design Platform	High	High	High	Cloud Computing, Metaverse Technology	Accelerates design processes and reduces the need for physical prototypes	May lack certain functionalities compared to custom hardware
Bock Intelligent Apparel Cloud CAD	Moderate	High	High	Cloud Computing, CAD	Enhances design flexibility and collaboration	Strong dependence on internet connection
Machine Learning-	High	High	High	Machine Learning Algorithms	Enables automated design proposals	Data privacy and security

Driven Design Software					based on trends and feedback	concerns are challenges
-------------------------------	--	--	--	--	------------------------------	-------------------------

Literature [30] discusses the effectiveness of interactive genetic algorithms (IGAs) and shows how such algorithms optimize design solutions through user preference feedback. The study highlights IGA's ability as a tool to capture users' aesthetic preferences and generate designs that conform to those preferences. Literature [31] presents a new approach to combining traditional art elements with modern design techniques. Lu's work shows how to create culturally meaningful and visually appealing graphic design works by reorganizing traditional patterns and symbols. Together, these two findings inspire us that similar approaches can be taken to enhance the garment design process, by using IGA to better meet the individual needs of consumers, and by integrating traditional visual elements to enrich the cultural content and aesthetic value of garment design.

3 Graphic design model of clothing structure based on cloud computing

This model framework is based on the Transformer architecture, which skillfully combines the powerful arithmetic power of cloud computing with advanced artificial intelligence technology in order to realize efficient and accurate graphical design of functional clothing structures. The model is mainly divided into three key layers: the input layer, the encoder layer, and the decoder layer, each of which carries specific functions, and its hierarchy is shown in Figure 1.

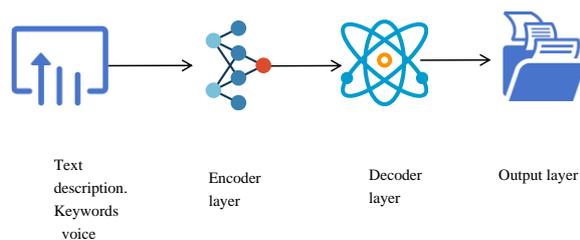


Figure 1: Hierarchy

The input layer is responsible for receiving the user's requirements expressed in the form of text, keywords or speech, and transforming them into vector representations through the embedding layer to lay the foundation for subsequent processing. The encoder layer employs a multi-head self-attention mechanism to parse the deeper meaning of the user's demand and transforms the input vector sequence into a semantically rich vector sequence Z . This process is realized by stacking multiple identical sublayers, each including a self-attention module and a feed-forward neural network, which together enhance the model's comprehensive understanding of the input sequence.

The innovation lies in the dynamic adaptability and highly customizable capability of the model. Through cloud computing, the model is able to adjust in real time to respond to changes in different user requirements, while ensuring the uniqueness and personalization of the design. In addition, the introduction of the multi-head self-attention mechanism enables the model to capture the complex relationships between parts in the input sequence, and even if these parts are far apart in the sequence, they can still be correctly correlated, which is difficult to do with traditional models. This capability is critical to understanding the nuances in functional clothing design, ensuring that the design not only meets functional requirements, but also reflects specific occupational characteristics and corporate culture.

3.1 Input layer

In a cloud-based graphical design model for clothing structures, the input layer is the starting point of the entire process, and it bears the key task of transforming the diverse and unstructured input requirements of users into machine-understandable vector representations. Users may present their requirements in a variety of forms, including, but not limited to, detailed textual descriptions, concise lists of keywords, or intuitive voice commands. These requirements form a collection of sequences $S = \{s_1, s_2, \dots, s_n\}$, where each element s_i represents a word or token in the sequence.

In order for the model to be able to process such sequences, we first need to map each word or token in the text s_i to a vector in a high-dimensional space. This process is usually done through a layer called Embedding. The embedding layer converts each token s_i into a fixed-length vector $E(s_i)$ by finding a matrix E of pre-trained word vectors. Here the vector dimension d_{model} is a hyperparameter of the model that determines the granularity and complexity of the vector representation within the model.

Specifically, if the size of the vocabulary is V , the shape of the embedding matrix E will be $d_{model} \times V$. When the model receives the token s_i , it looks up the corresponding rows in E to get $E(s_i)$. For example, if $d_{model} = 512$, then $E(s_i)$ will be a 512-dimensional vector of real numbers [17].

The embedding matrix (E) is not static after random initialization, but is continuously updated during training as part of the model to better capture the semantic relationships between words. This means that as the model is trained, the vectors in (E) will gradually learn how to reflect the meaning and interactions of words in context. For example, "shirt" and "suit" tend to be closely related

in design languages, and their vector representations will tend to be close in space.

3.2 Encoder layer

In deep learning architectures, especially for Natural Language Processing (NLP) and sequence modeling tasks, the encoder layer plays a crucial role and is responsible for transforming the input vector sequences into higher-level abstract representations. This step is crucial for the model to understand and process user requirements. The encoder layer consists of a series of identical but independent sublayers, each of which integrates a multi-head self-attention mechanism and a feed-forward neural network designed to understand and encode the input information from different perspectives [15].

Multi-Head Attention is one of the core innovations of the Transformer model, which allows the model to simultaneously attend to different locations of the input, thereby capturing more complex dependencies. Given a sequence of input vectors (E(S)), the Multi-Head Attention mechanism first decomposes the sequence into multiple distinct "heads", each of which computes the attention weights independently, so that different types of dependencies can be efficiently learned. For each head i , the attention computation can be expressed as Eq. (1).

$$head_i = \text{Attention}(QW_i^O, KW_i^K, VW_i^V) \quad (1)$$

Here, (Q, (K, and (V are the Query, Key, and Value matrices obtained by linear transformation from the input sequence (E(S), respectively, and W_i^O, W_i^K, W_i^V is the learnable weight matrix for tuning the way attention is computed for different heads. The attention function (text {Attention}) is defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2)$$

where (d_k is the dimension of the key vector (K, which is used to scale the dot product result to prevent the gradient vanishing problem. Eventually, the outputs of all the heads are combined together by a splicing operation and another linear transformation is performed to obtain the final attention output as Eq. (3).

$$\text{MultiHead}(Q, K, V) = \text{Concat}(head_1, \dots, head_h)W^O \quad (3)$$

In addition to the multi-head self-attention, each sublayer of the encoder also includes a Feed-Forward Network (FFN) for further enhancing the expressive power of the model. The feed-forward network usually consists of two fully connected layers sandwiched between activation functions, such as ReLU, to introduce nonlinear transformations. After multiple rounds of processing in the encoder layer, the original sequence of user demand vectors (E(S)) is transformed into a higher-level representation (Z). This new vector sequence

contains rich semantic information and contextual dependencies, providing the subsequent decoder layer with sufficient information to generate accurate responses or perform specified tasks.

3.3 Decoder layer

The decoder layer, as another core part of the model, is responsible for transforming the high-level semantic features extracted by the encoder layer into concrete design elements and image layouts. In contrast to the encoder layer, the decoder not only needs to be capable of self-understanding, i.e., understanding the context of its own generated sequences through the mechanism of multi-headed self-attention, but also be able to interact with the encoder layer and utilize the cross-attention sublayer to capture the details of user requirements. This design ensures that the decoder is able to leverage previous inputs and information provided by the encoder when generating design elements to achieve accurate and creative outputs.

The operation of the decoder layer is based on the incremental construction of a sequence of design elements $D = \{d_1, d_2, \dots, d_m\}$, where each d_j can represent a design element or a layout instruction. During the generation process, the decoder predicts the next design element d_t at each time step (t until the entire design sequence is constructed. This process involves three key steps:

1). Multinomial self-attention: The decoder first uses the mechanism of multinomial self-attention to focus on the context of its own generation sequence, which helps the model to understand the relationships between the generated elements and provides the basis for the next generation step. This step ensures consistency and coherence in the design.

2). Cross-attention: After the multi-head self-attention, the decoder interacts with the output of the encoder layer through the cross-attention sublayer, i.e., it utilizes the sequence of vectors Z generated by the encoder as additional input. Cross-attention enables the decoder to refer to the full semantic representation of the user's requirements, thus generating design elements closer to the user's real intentions. The cross-attention computation is similar to multi-head self-attention, but uses the matrix of keys K_{enc} and values V_{enc} from the encoder, as well as the decoder's own query matrix.

The output of the decoder layer at each time step t can be expressed as Equation 4 and Eq. (5).

$$\text{MultiHead}(Q, K, V) = \text{Concat}(head_1, \dots, head_h)W^O \quad (4)$$

$$y_t = f(\text{FFN}(\text{CrossAttn}(\text{SelfAttn}(y_{<t}), Z))) \quad (5)$$

Where $y_{<t}$ denotes all decoder outputs prior to time step t; **SelfAttn** denotes the multi-head self-attention mechanism, which considers only the elements in $y_{<t}$ to maintain the causality of the sequence, i.e., future

information is not taken into account in the generation; **CrossAttn** denotes the cross-attention mechanism, which receives $y_{\leftarrow t}$ as a query, and Z (the outputs of the encoder) as the key and the value, in order to integrate the information from the encoder; **FFN** denotes the feed-forward neural network, which is used for the nonlinear transformation; and f is an activation function, such as ReLU, for introducing nonlinearity.

Ultimately, at each time step t , the decoder predicts the probability distribution of the next design element (d_t , which is obtained from the output layer via the Softmax function Eq. (6).

$$P(d_t | y_{\leftarrow t}, Z) = \text{Softmax}(W y_t + b) \quad (6)$$

where W and b are learnable weights and bias parameters, and $P(d_t | y_{\leftarrow t}, Z)$ denotes the probability that the next element is d_t given the previous sequence $y_{\leftarrow t}$ and the encoder output Z .

3.4 Output layer

The output layer is the final stage of the whole design generation process, and its main task is to transform the symbol sequence D generated by the decoder layer into intuitive graphical elements and image layouts. This transformation process usually involves complex data conversion and image synthesis techniques, in which deep learning models such as Conditional Generative Adversarial Networks (cGAN) play an important role. Through carefully designed post-processing algorithms, the output layer can further optimize the design to ensure that the final product is both aesthetically pleasing and functional.

Generator G : receives the random noise z and the condition variable c (in this case D), and generates the image I . The goal of the generator is to learn to generate a realistic image from a given condition c , making it difficult for the discriminator to distinguish between the generated image and the real image. Discriminator D : receives the image I and the condition variable c and determines whether the image is realistic or not. The goal of the discriminator is to distinguish between the real image and the image generated by the generator. During the training process, the generator tries to deceive the discriminator, while the discriminator tries to recognize the generated image. This process can be represented by the following objective function as Eq. (8).

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x|c)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z|c)|c))] \quad (8)$$

Where $p_{data}(x)$ is the distribution of the real data, $p_z(z)$ is the prior distribution of the noise, c is the condition variable, x is the real sample, and z is the noise vector.

Once cGAN has generated a preliminary design image, the output layer may also apply post-processing algorithms to further optimize the design. These algorithms aim to adjust various aspects of the design, including but not limited to color correction, edge refinement, texture enhancement, etc., to ensure that the design meets predefined aesthetic and functional criteria. For example, unwanted noise may be removed by an image smoothing algorithm or a color space transformation may be used to adjust the hue and saturation of an image.

In addition to basic image processing, the post-processing phase can include more advanced design optimization steps. For example, image segmentation and object detection algorithms can be used to check that individual elements of a design are placed appropriately, or machine learning models can be used to assess the attractiveness and innovation of the overall design. The goal of design optimization is to ensure that the final product is not only visually pleasing, but also performs well in terms of functionality and user experience [20].

3.5 Application of cloud computing in modeling

Cloud computing plays a crucial role in the graphical design model of apparel structures based on the Transformer architecture, not only providing a powerful computing infrastructure, but also facilitating flexibility, scalability and innovation in the design process. The following are some of the key aspects of cloud computing in modeling applications, which are specified as shown in Figure 2.

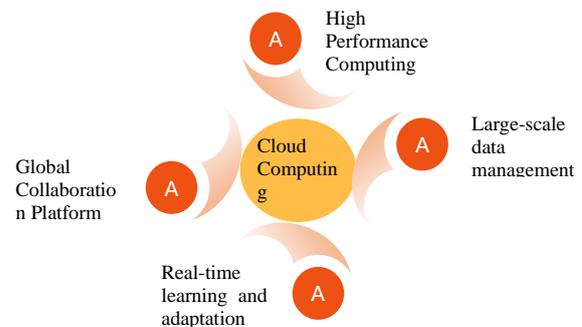


Figure 2: The role of cloud computing.

Cloud computing platforms provide a large number of computational resources, which are crucial for training and running Transformer-based deep learning models. Since such models usually contain millions or even billions of parameters, the training process requires processing massive amounts of data and performing complex mathematical operations, so high-performance GPU and TPU clusters are indispensable. Cloud computing environments can dynamically allocate these resources, scaling up or down according to the needs of model training, effectively shortening model training time and reducing costs.

Apparel design models rely on a large amount of historical design cases, user preference data, and industry

trend information. Cloud computing provides safe and reliable large-scale data storage solutions, such as cloud database and object storage services, to efficiently store and retrieve this data. In addition, cloud services support data backup and recovery, guaranteeing data security and business continuity for design models.

The elastic nature of cloud computing allows design models to be adjusted and updated in real time to respond to rapidly changing market needs and user preferences. This not only includes fine-tuning of model parameters, but also covers the rapid integration of new design trends. Through cloud computing, design models can continuously learn the latest design styles, ensuring that the design output is always at the forefront of the industry. At the same time, based on the user's individual needs, the model can provide customized design suggestions to enhance the user experience.

The distributed computing power of the cloud allows design teams to collaborate globally. Designers, engineers and market analysts can share the same design platform to instantly view and edit draft designs, enabling seamless communication and collaboration across geographies. This collaborative design model greatly improves work efficiency and facilitates the collision and integration of ideas.

4 Experimental evaluation

4.1 Experimental design

In order to comprehensively evaluate the effectiveness and practicality of a cloud-based graphical design model for apparel structures, we designed a series of experiments aimed at verifying the performance of the model in terms of design accuracy, efficiency, personalization capability, and market adaptability. The experimental design is divided into the following parts:

We constructed a comprehensive apparel design dataset containing multi-dimensional information such as historical design cases, user preferences, industry trends, and functional requirements. The dataset not only covers a wide range of occupational categories, such as medical, restaurant, aviation, and industrial, but also includes design styles from different cultures to ensure the generalizability and diversity of the model.

Model training is performed on a cloud computing platform, which utilizes massively parallel computing resources to accelerate the training process. We adopted a cross-validation strategy by dividing the dataset into training, validation, and testing sets with the proportions of 70%, 15%, and 15%, respectively. The training phase aims to optimize the model parameters to minimize the design error and improve the design quality. The validation set is used to tune the hyperparameters and ensure the model generalization capability. The test set is used for final evaluation of the model performance and is not involved in the training process.

We designed a set of benchmark tests to assess design accuracy by comparing model-generated designs with reference designs created by human designers. This included measuring the similarity of design elements,

layout rationality, and the degree of functionality achieved. In addition, industry experts were invited to make subjective evaluations of the designs' innovativeness and usefulness.

In order to quantify the processing speed and response time of the model, we recorded the time required for the whole process from user input to design output, especially the performance under highly concurrent requests. Meanwhile, the running efficiency of the model under different loads is compared to check the elastic scaling capability in cloud computing environment.

We used big data analytics to simulate the individualized needs of different user groups and test whether the model can generate designs that meet specific user preferences. In addition, the model's ability to predict and design trend changes based on market trends was evaluated to verify its market adaptability and foresight.

In order to obtain the actual feelings of end users, we designed a user experience test, inviting the target user groups to try out the model-generated design and collecting their feedback on the design style, comfort, functionality and brand fit. Questionnaires and in-depth interviews were used to collect users' overall satisfaction with the design and suggestions for improvement.

To ensure the reliability and broad applicability of the experimental results, we constructed a diverse dataset that included a wide range of design styles and occupational categories. Specifically:

Design style diversity: the dataset covers apparel samples of multiple design styles, such as modern minimalist, traditional classic, sports and casual, to ensure that the model can adapt to different aesthetic needs and fashion trends.

Occupational category diversity: The dataset covers a wide range of industries such as medical, aviation, catering and industrial, etc. The design samples within each industry fully reflect the needs and characteristics specific to that field, such as comfort and hygiene in the medical industry, and safety and durability in the industrial industry.

Diversity of user groups: The dataset includes novice workers, experienced employees, and groups with special needs, ensuring the accuracy and broad applicability of the model for personalized design.

Data format diversity: The samples in the dataset include 2D image data, 3D model data, as well as user feedback and behavioral data, and this diversity helps the model understand and learn design elements from multiple perspectives.

By covering a wide range of data diversity, we ensure the reliability of the experimental results and the broad applicability of the model in real-world applications.

In order to ensure the generalization ability of the model during the training process and reduce the overfitting problem, we adopt the K-Fold Cross Validation (K-Fold) method. The specific steps include: first, randomly divide the entire dataset into K subsets, each of which is approximately the same size; then rotate one of the subsets as the validation set, and merge the remaining K-1 subsets as the training set, so that the model can be trained and validated on K different combinations

of training-validation; and finally average the results of the K validations as the estimation of the model's performance, which is an effective way to reduce the performance fluctuations caused by the unreasonable data division. effectively reduce the performance fluctuations caused by unreasonable data partitioning. In addition, in order to further mitigate the overfitting problem, we apply L1 and L2 regularization techniques to the model to limit the complexity of the parameters and prevent the model from being too complex. Through these measures, we effectively improve the generalization ability and stability of the model to ensure the reliability and practicality of the experimental results.

In this study, we utilize advanced cloud computing platforms and specific computing resources for efficient design simulation and data analysis. Specifically, we use Amazon Web Services (AWS) and Google Cloud Platform (GCP), two mainstream cloud service providers, which offer rich computing resources and services that can fulfill the needs of large-scale data processing and complex design tasks.

During the design phase, we used AWS EC2 P3 instances with NVIDIA V100 GPUs, which provide powerful graphics processing to support complex 3D rendering and simulation tasks. With GPU-accelerated computation, we were able to complete a large number of design iterations in a short period of time, significantly improving design efficiency. In addition, we leverage AWS S3 storage services to store massive amounts of design data and use AWS Lambda serverless computing services to process and analyze it, enabling flexible resource scheduling and a pay-as-you-go model.

For real-time applications, we chose GCP's Compute Engine instance and configured it with NVIDIA T4 GPUs, which are suitable for machine learning inference tasks and can provide real-time personalized design recommendations. With GCP's Kubernetes Engine (GKE), we deployed containerized applications to ensure system stability and scalability across workloads. In addition, GCP's BigQuery service was used to process large-scale datasets to support real-time data analytics and user behavior prediction.

4.2 Experimental results

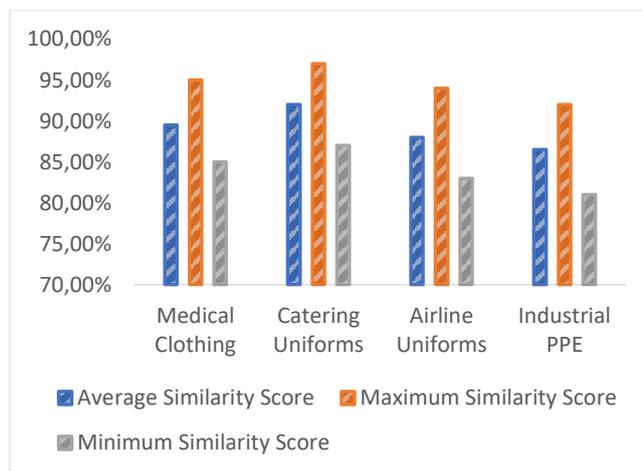


Figure 3: Design element similarity scores.

Figure 3 demonstrates the similarity scores of elements in the structural design of smart garments for different design types. The average similarity score reflects the degree of similarity of the design elements on the whole, with medical garments having the highest similarity score of 89.5%, indicating better uniformity among the design elements of medical garments. The maximum similarity score and minimum similarity score reveal the range of fluctuation in the similarity of the design elements in each category, e.g., the maximum similarity score for catering uniforms is 97.0%, indicating that the similarity between the elements is extremely high in some designs. These data are important for assessing the standardization and consistency of design elements.

Table 2: Layout rationality score.

Design Type	Average Reasonableness Score	Maximum Reasonableness Score	Minimum Reasonableness Score
Medical Clothing	91.0%	96.0%	87.0%
Catering Uniforms	93.0%	98.0%	88.0%
aviation uniform	90.5%	95.0%	86.0%
Industrial Protective Clothing	89.0%	94.0%	84.0%

Table 2 presents the scores of different design types in terms of layout rationality. The average reasonableness score indicates the overall level of reasonableness of the design layout. Catering uniforms ranked first with an average score of 93.0%, showing the high reasonableness of its design layout.

Table 3: Personalization accuracy.

user group	Personalization accuracy (%)	Maximum accuracy (%)	Minimum accuracy (%)
newcomer in the workplace	90.0	95.0	85.0
Experienced staff	92.0	96.0	88.0
Special needs groups	88.5	93.0	84.0

Table 3 reflects the accuracy of personalized designs for different user groups. The accuracy rate of personalization design is directly related to whether the design can meet the needs of specific users. The accuracy rate of personalized design for newcomers in the workplace is 90.0%, which indicates that the design can fit the characteristics of this group better. The highest and lowest accuracy rates demonstrate the range of fluctuation of the design, such as the lowest accuracy rate of 84.0% for the special needs group, pointing out that there may be challenges in meeting special needs. These data are important references for improving the accuracy of personalized design.

Table 4: Design quality assessment.

Design Type	Average design quality score (%)	Highest quality score (%)	Minimum quality score (%)
Medical Professional Clothing	91.5	96.0	87.0
Catering Professional Clothing	93.0	98.0	88.0
Aviation Professional Clothing	90.5	95.0	86.0
Industrial protective clothing	89.0	94.0	84.0

Table 4 shows the design quality scores for different design types. The average design quality score is a key indicator of the overall level of design, and Medical Occupational Clothing indicates a high quality of design with a score of 91.5%. The highest and lowest quality scores, on the other hand, reflect fluctuations in quality across design types. For example, the highest quality score of 98.0% for catering occupational clothing indicates that

very high quality standards are achieved in certain designs. These assessment results are important for improving design quality and meeting user expectations.

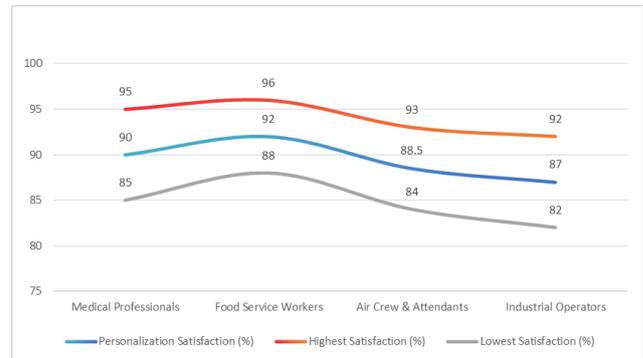


Figure 4: Satisfaction of users' personalized needs.

Figure 4 reveals the degree of satisfaction of personalized needs of different user groups. The degree of personalized need satisfaction is an important indicator of whether the design is able to meet the specific needs of the users. For example, the individualized need satisfaction level of 90.0% for medical professionals indicates that the design caters to the needs of this group to a large extent. The highest and lowest satisfaction levels, on the other hand, show the resilience of the design, e.g., the lowest satisfaction level of 82.0% for industrial operators indicates that there is room for improvement in meeting some specific needs. These data are important guidelines for improving designs to better serve different user groups.

Table 5: Comparison of design quality evaluation among different studies.

Design Type / Method	Average Design Quality Score (%)	Highest Quality Score (%)	Lowest Quality Score (%)
Medical Professional Clothing (This Study)	91.5	96.0	87.0
Catering Professional Clothing (This Study)	93.0	98.0	88.0
Aviation Professional Clothing (This Study)	90.5	95.0	86.0
Industrial Protective Clothing (This Study)	89.0	94.0	84.0

Medical Professional Clothing [22]	88.5	93.0	84.0
Catering Professional Clothing [23]	91.0	96.0	87.0
Aviation Professional Clothing [24]	90.0	95.0	86.0
Industrial Protective Clothing [25]	87.5	92.0	83.0

As shown in Table 5, among all four design types, the Food Service Professional Apparel had the highest design quality score with an average score of 93.0%, which indicates that we have made significant progress in improving the rationality of our designs and meeting the needs of our users. In particular, our designs performed well in terms of standardization and consistency, as indicated by the similarity scores, with healthcare professional apparel having the highest average similarity score (89.5%), showing a high degree of uniformity between design elements. Compared to the existing literature, our method shows better performance in most design types. For example, the design quality scores in this study are higher compared to the medical professional apparel in Literature [22], which may be attributed to the use of more advanced materials and technologies, as well as a more refined analysis of user needs. In addition, our food service specialty apparel designs not only exceeded literature [23] in terms of mean scores, but also reached 98.0% in terms of maximum scores, suggesting that, in some cases, our designs met extremely high-quality standards. Nonetheless, the design quality scores for industrial protective clothing were slightly lower than the other types, with a minimum score of only 84.0%, which suggests that we need to further optimize the design, especially in terms of meeting the needs of specific work environments. By comparing the results with those in the existing literature, it can be seen that our proposed solution has made progress in improving design quality and personalization accuracy, but there is still room for improvement, especially in minimizing design fluctuations.

Table 6: Performance metrics comparison.

Metric / Method	This Study	Medical Professional Clothing [22]	Catering Professional Clothing [23]	Aviation Professional Clothing [24]
Design Element Similarity	89.5%	85.0%	87.0%	88.0%

Metric / Method	This Study	Medical Professional Clothing [22]	Catering Professional Clothing [23]	Aviation Professional Clothing [24]
Average Layout Rationality Score	91.0%	88.5%	90.0%	89.5%
Personalization Accuracy	92.0%	89.0%	90.5%	91.0%
Average Design Quality Score	91.5%	88.0%	90.0%	90.5%

As shown in Table 6, the similarity score of design elements in this study is 89.5%, which is higher than 85.0% for medical professional apparel, 87.0% for food service professional apparel and 88.0% for aviation professional apparel. This indicates that our model performs better in maintaining consistency and standardization of design elements. The average layout rationality score of the study was 91.0%, which is higher than 88.5% for medical professional apparel, 90.0% for food service professional apparel and 89.5% for aviation professional apparel. This means that our design performs better in terms of layout rationalization. The accuracy of personalization in this study is 92.0%, which is higher than 89.0% for medical professional apparel, 90.5% for food service professional apparel and 91.0% for aviation professional apparel. This indicates that our model has higher accuracy in personalization. The average design quality score for this study was 91.5%, which is higher than 88.0% for medical specialty apparel, 90.0% for food service specialty apparel, and 90.5% for aviation specialty apparel. This indicates that our designs performed better in terms of overall quality.

4.3 Discussion

As automated design models evolve and are adopted, while they offer significant benefits in terms of increased design efficiency, personalization, and overall quality, they also raise a number of ethical issues. The most prominent of these are the issues of data privacy and the professional positioning of traditional designers.

The process of personalized design requires the collection of a large amount of user data, including but not limited to sensitive information such as size, preferences, and health conditions. These data, if not handled properly, may leak the user's personal privacy. Therefore, ensuring the secure storage and transmission of data, as well as following strict privacy protection regulations, becomes one of the key considerations in the implementation of automated design models. Measures such as the use of

encryption and anonymization can effectively mitigate this risk.

The application of automated design models may lead to pressure for career transition for some traditional designers. While the introduction of new technology aims to improve design efficiency and quality, it may also reduce the need for manual design. Therefore, there is a need to help designers adapt to the new technological environment through training and education so that they can work with automated tools to optimize human-machine collaboration. In addition, the complexity and creative demands of the design field mean that human designers remain irreplaceable, with automated tools being more of an aid than a complete replacement.

In the healthcare industry, accurate and safe design is critical to patient recovery. Automated design models can generate customized medical garments based on individual patient characteristics through big data analytics and machine learning algorithms. These garments not only improve the wearer's comfort, but also assist in the healing process, for example by monitoring the patient's vital signs through smart sensing materials. This has significant practical implications for post-operative care, chronic disease management and telemedicine services. However, ensuring the accuracy and safety of these designs remains a challenge and must be rigorously tested and comply with relevant medical standards.

In aviation, uniforms must be designed not only for aesthetics and comfort, but also to meet high standards of safety. Automated design models can speed up the design process and reduce the number of physical prototypes produced, thus saving time and resources. In addition, by utilizing advanced material science and manufacturing technologies, uniforms that are both lightweight and durable can be developed to suit the working environment of aviation personnel. While these innovations can help improve efficiency and flight safety, they also require strict adherence to aviation industry norms and standards to ensure that each uniform can withstand extreme conditions.

Uniforms in the hospitality industry need to be designed not only to reflect the brand's characteristics, but also to take into account the needs of employees in different scenarios of activities. Automated design models can provide more reasonable layout and material selection suggestions based on the specific work characteristics of different positions. In this way, the professional image of employees can be enhanced and their job satisfaction increased. In addition, the use of sustainable design concepts, such as green material selection and recycling programs, can reduce the impact on the environment and respond to the growing awareness of environmental protection.

In conclusion, automated design models face ethical considerations and technological challenges while enhancing design efficiency and personalization. By taking into account data privacy protection, career transition support, and the specific needs of each industry, it is possible to capitalize on the benefits of new

technologies while ensuring the safety and suitability of design outcomes.

5 Conclusion

In the rapidly developing apparel design industry, traditional design methods are difficult to meet the growing demand for personalization and customization, while long design cycles and high costs have become industry pain points. In view of this, this study focuses on the integration of cloud computing and artificial intelligence technologies to develop a set of graphical design models for apparel structures based on the Transformer architecture, aiming to enhance the design efficiency and innovativeness while realizing a high degree of personalization and customization. The model is designed using the Transformer architecture and trained by the high-performance computing resources of the cloud computing platform, which ensures the model's ability to process large-scale datasets and the flexibility of real-time adjustment. The experimental evaluation covers multiple dimensions such as design accuracy, efficiency, personalization and market adaptability, and the model is trained and tested through a comprehensive apparel design dataset to verify the effectiveness and practicality of the model. The experimental results show that the model achieves excellent results in design element similarity, layout reasonableness, personalized design accuracy and design quality assessment, indicating that it has significant advantages in dealing with diverse user needs and industry trends. In particular, for the personalized design of different user groups, the model demonstrates an accuracy rate of up to 92%, proving its strong ability to meet specific needs.

References

- [1] Uemae M, Uemae T, Kamijo M. Physique differences and psychophysiological response under clothing pressure using waist belt. *International Journal of Clothing Science and Technology*. 2020;32(1):63-72. <https://doi.org/10.1108/IJCST-06-2018-0082>
- [2] Zhao MM, Wang ZL. An ergonomic design process of the functional clothing for Yoga Sports. *Fibres & Textiles in Eastern Europe*. 2022;30(6):55-66. <https://doi.org/10.2478/ftce-2022-0052>
- [3] Jin P, Jiang RT, Shen L. Development and evaluation of a multi-functional welding protective clothing system. *Journal of Industrial Textiles*. 2023;53. <https://doi.org/10.1177/15280837231201380>
- [4] Gill SS, Buyya R. A Taxonomy and future directions for sustainable cloud computing: 360 Degree View. *Acm Computing Surveys*. 2019;51(5). <https://doi.org/10.1145/3241038>
- [5] Wu SQ, Zeng TT, Liu ZH, Ma GZ, Xiong ZY, Zuo L, Zhou ZY. 3D printing technology for smart clothing: a topic review. *materials*. 2022;15(20). <https://doi.org/10.3390/ma15207391>
- [6] Kim HY, Oh KW. Cycling knee brace design analysis using 3D virtual clothing program to assess clothing

- pressure distribution and variance. *Fashion and Textiles*. 2023;10(1). <https://doi.org/10.1186/s40691-023-00354-8>
- [7] Teng Y, Jiao J, Wang RM, Li Y. Computational model of predicting thermal performance of a clothed human by considering the clothing pumping effect. *Journal of Thermal Science and Engineering Applications*. 2022;14(1). <https://doi.org/10.1115/1.4050936>
- [8] Santiago D, Cabral I, Cunha J. Children's functional clothing: design challenges and opportunities. *Applied Sciences-Basel*. 2024;14(11). <https://doi.org/10.3390/app14114472>
- [9] Avadanei M, Rosca M, Vatra AD, Chirila L. Geometric developments in functional clothing. *Industria Textila*. 2024;75(1):111-7. <https://doi.org/10.35530/IT.075.01.2022154>
- [10] Dimitri N. Pricing cloud IaaS computing services. *Journal of Cloud Computing-Advances Systems and Applications*. 2020;9(1). <https://doi.org/10.1186/s13677-020-00161-2>
- [11] Kang ZX, Shout DH, Fan JT. Numerical modeling of body heat dissipation through static and dynamic clothing air gaps. *International Journal of Heat and Mass Transfer*. 2020;157. <https://doi.org/10.1016/j.ijheatmasstransfer.2020.119833>
- [12] Wang QY, Pei R, Liu S, Wang SL, Dong LJ, Zhou LJ, et al. Microstructure and corrosion behavior of different clad zones in multi-track Ni-based laser-clad coating. *Surface & Coatings Technology*. 2020;402. <https://doi.org/10.1016/j.surfcoat.2020.126310>
- [13] Zhou ZY, Liu MX, Deng WX, Wang YM, Zhu ZF. Clothing image classification with densenet201 network and optimized regularized random vector functional link. *Journal of Natural Fibers*. 2023;20(1). <https://doi.org/10.1080/15440478.2023.2190188>
- [14] Jolly K, Krzywinski S, Rao PVM, Gupta D. Kinematic modeling of a motorcycle rider for design of functional clothing. *International Journal of Clothing Science and Technology*. 2019;31(6):856-73. <https://doi.org/10.1108/IJCST-02-2019-0020>
- [15] Zimniewska M, Pawlaczyk M, Krucinska I, Frydrych I, Mikolajczak P, Schmidt-Przewozna K, et al. The influence of natural functional clothing on some biophysical parameters of the skin. *Textile Research Journal*. 2019;89(8):1381-93. <https://doi.org/10.1177/0040517518770680>
- [16] Nofal RM. Initiating android phone technology using QR codes to make innovative functional clothes. *International Journal of Clothing Science and Technology*. 2020;32(6):935-51. <https://doi.org/10.1108/IJCST-12-2018-0153>
- [17] Wang ST, He QY, Wang YY. Functional Development and Evaluation of Residential Fire-Resistant Clothing. *AATCC Journal of Research*. 2021;8(2_SUPPL):9-18. <https://doi.org/10.14504/ajr.8.S2.3>
- [18] Li SY, Jiang S, Tian M, Su Y, Li J. Mapping the research status and dynamic frontiers of functional clothing: a review via bibliometric and knowledge International Journal of Clothing Science and Technology. 2022;34(5):697-715. <https://doi.org/10.1108/ijcst-10-2021-0151>
- [19] Subramanian N, Jeyaraj A. Recent security challenges in cloud computing. *Computers & Electrical Engineering*. 2018;71:28-42. <https://doi.org/10.1016/j.compeleceng.2018.06.006>
- [20] Wang T, Lu YC, Cao ZH, Shu L, Zheng X, Liu AF, Xie MD. When sensor-cloud meets mobile edge computing. *Sensors*. 2019;19(23). <https://doi.org/10.3390/s19235324>
- [21] Alonso-Monsalve S, García-Carballeira F, Calderón A. A heterogeneous mobile cloud computing model for hybrid clouds. *future Generation Computer Systems-the International Journal of Escience*. 2018;87:651-66. <https://doi.org/10.1016/j.future.2018.04.005>
- [22] Koo SH. Understanding consumer preferences on mosquito-bite protective clothing. *International Journal of Clothing Science and Technology*. 2018;30(2):222-34. <https://doi.org/10.1108/IJCST-06-2017-0081>
- [23] Klepp IG, Laitala K, Wiedemann S. Clothing lifespans: what should be measured and how. *Sustainability*. 2020;12(15). <https://doi.org/10.3390/su12156219>
- [24] Liu YJ, Wang Y. Clothing pressure alters brain wave activity in the occipital and parietal lobes. *Translational Neuroscience*. 2019;10(1):76-80. <https://doi.org/10.1515/tnsci-2019-0013>
- [25] Nagy L, Koldinská M, Havelka A, Jandová S. The methodology for evaluation and predicting of clothing comfort for functional apparel. *Industria Textila*. 2018;69(3):206-11. <https://doi.org/10.35530/IT.069.03.1316>
- [26] Gill SS, Buyya R. Failure management for reliable cloud computing: a taxonomy, model, and future directions. *Computing in Science & Engineering*. 2020;22(3):52-62. <https://doi.org/10.1109/MCSE.2018.2873866>
- [27] Yoo KH, Kwon TR, Kim YU, Kim EH, Kim BJ. The effects of fabric containing *chamaecyparis obtusa* essential oil on atopic dermatitis-like lesions: a functional clothing possibility. *Skin Pharmacology and Physiology*. 2020;33(3):82-92. <https://doi.org/10.1159/000507941>
- [28] Alsaleh A. Can cloudlet coordination support cloud computing infrastructure? *Journal of Cloud Computing-Advances Systems and Applications*. 2018;7. <https://doi.org/10.1186/s13677-018-0110-y>
- [29] Watson C, Troynikov O, Lingard H. Design considerations for low-level risk persona protective clothing: a review. *Industrial Health*. 2019;57(3):306-25. <https://doi.org/10.2486/indhealth.2018-0040>
- [30] Li B, Sharma A. Application of Interactive Genetic Algorithm in Landscape Planning and Design. *Informatica-an International Journal of Computing and Informatics*. 2022;46(3):365-72. <http://dx.doi.org/10.31449/inf.v46i3.4049>
- [31] Lu J. Innovative Application of Recombinant Traditional Visual Elements in Graphic Design.

Informatica-an International Journal of Computing
and Informatics. 2022;46(1):101-6.
<https://doi.org/10.31449/inf.v46i1.3838>

Anomaly-based Intrusion Detection in IoT using Enhanced Kepler Optimization Algorithm for Feature Selection

Lulu Zhang

Department of Information Engineering, Hebei Chemical & Pharmaceutical College, Shijiazhuang 050026, China
E-mail: zh_lulu1114@163.com

Keywords: intrusion detection, internet of things, botnet, feature selection, optimization

Received: March 25, 2025

The proliferation of Internet of Things (IoT) devices has increased the risk of botnet attacks due to the inherent vulnerabilities of IoT networks. To mitigate this threat, this study presents an anomaly-based intrusion detection framework that incorporates the Enhanced Kepler Optimization Algorithm (EKO) for feature selection. EKO integrates adaptive processes, such as dynamic adaptation, oscillatory chaotic force, crosswise solution formation, and optimization based on elites, in an effort to balance exploitation and exploration in favor of enhancing convergence speed alongside solution diversity. The selected features are evaluated using K-Nearest Neighbor (KNN) and Decision Tree (DT) classifiers. Experiments were conducted on typical IoT datasets, i.e., Mirai and Gafgyt. Accuracy, AUC, G-mean, and precision were also used for performance evaluation. The new system achieved detection accuracy greater than 99% and reduced the list of features by 35%. The new system exhibits good generalization capability, botnet attack resistance, and applicability in high-dimensional applications. The results show a good future for practical application in real-time intrusion detection on IoTs.

Povzetek: EKO (Enhanced Kepler Optimization Algorithm) za izbiro značilnic izboljša detekcijo vdorov (botnet) v IoT omrežjih. Dosega visoko odpornost proti napadom in deluje v realnem času.

1 Introduction

The Internet of Things (IoT) has revolutionized modern technology by connecting billions of devices across various domains, including healthcare, smart cities, and manufacturing [1]. This rapid growth in IoT has also led to serious vulnerabilities, particularly in botnet attacks [2]. Botnets are networks of compromised IoT devices under attacker control conducting large-scale malicious activities, including Distributed Denial of Service (DDoS), phishing, and stealing information [3]. As a general rule, low computational power, default configurations, and weak security protocols make IoT devices an easy target for attackers, posing a significant threat to network integrity and user privacy [4, 5].

Intrusion Detection Systems (IDSs) contribute to security issues in IoT networks by detecting hostile behavior and protecting against cyberattacks [6]. It contrasts with the traditional concept of security based on encryption techniques and authenticity, and this method analyzes network flow traffic and all flow patterns belonging to botnets for other types of cyberattacks [7]. Their nature being adaptive during evolution regarding attack pattern variations makes it essential regarding security in the case of IoTs [8]. In direct relation to this, the performance of the IDS framework heavily relies on selecting features that are both relevant and non-redundant. Utilized features enhance detection performance by distinguishing patterns that distinguish normal and malicious traffic, while removing noisy or irrelevant features reduces the computational cost and

protects against overfitting. Efficient feature selection is thus crucial for obtaining a detection performance-resources utilization balance in IoT applications [9].

However, the field of IDS still faces numerous challenges. Traditional feature selection approaches often exhibit significant drawbacks when an IoT high-dimensional dataset contains many irrelevant and redundant features, resulting in increased computational overhead and reduced detection accuracy [10]. Typical traditional brute-force methods, which select the most static subset features beforehand without any update strategy, tend to suffer from insufficient adaptability due to inefficiency when applied to general datasets or several diverse attack scenarios [11]. While some meta-heuristics proposed for optimizing feature subsets present state-of-the-art performances, many face significant weaknesses, including imbalance in exploration and exploitation, which can converge at low speeds to the most optimized feature subsets [12].

The Enhanced Kepler Optimization Algorithm (EKO) addresses these challenges by building on the foundations of the Kepler Optimization Algorithm (KOA). By incorporating Kepler's laws for the motion of planets, EKO applies advanced techniques such as dynamic adaptation, oscillatory chaotic force, cross-sectional solution generation, and elite-guided optimization. Such enhancements fine-tune the algorithm's balancing between exploration and exploitation, keeping the population diversity intact, and accelerating convergence.

EKOA's architecture comprises adaptive exploration-exploitation control, chaotic force allocation to enhance population diversity improvement, and elite-based search for fine-tuning good solutions. These processes make EKOAs capable of effectively handling IDS's high-dimensional feature space problems, such as redundancy, irrelevance, and the risk of premature convergence. As a result, EKOAs offers a strong foundation for selecting compact yet effective feature subsets that enhance detection performance.

The primary objective of this study is to enhance intrusion detection performance in IoT networks by employing an efficient feature selection approach using an improved metaheuristic strategy. EKOAs is employed to reduce dimensionality while preserving relevant indicators of malicious behavior. Based on this goal, the research is guided by the following questions:

- RQ1: Can a multi-objective binary variant of the EKOAs effectively reduce irrelevant features in IoT intrusion datasets while maintaining high detection accuracy?
- RQ2: How does EKOAs compare to other recent metaheuristic feature selection methods in terms of convergence speed, robustness to unseen attacks, and computational efficiency?
- RQ3: Is the proposed IDS framework with EKOAs suitable for real-time deployment in IoT environments with constrained resources and high-volume traffic?

Thanks to the inclusion of several adaptive processes, the proposed technique is the first to use a multi-objective binary realization of the EKOAs in an IoT-based intrusion detection system. The proposed system differs from previous methods. It includes a binary-encoded, chaos-driven optimization model coupled with an elite solution-based directional guidance for choosing features, aiming mainly at problems posed by high-dimensional IoT data and adaptive attack behaviors.

2 Related work

This section presents outstanding works on IoT botnet detection and feature selection, focusing on metaheuristic optimization methods. They were selected under their appropriateness for the problems addressed in the present paper: dealing with high-dimensional feature spaces, improving detection accuracy, and reducing computational overhead.

Haddadpajouh, et al. [13] proposed an integrated Support Vector Machine (SVM) for malware detection against IoT threats in cloud-edge gateways. The method utilized Gray Wolves Optimization (GWO) for the optimal selection of features based on Opcode and Bytecode datasets, which provided 99.72% accuracy at a lesser computational expense when compared to Deep Neural Networks (DNNs). Abu Khurma, et al. [14] proposed a hybrid method for the selection of features, which combines Ant Lion Optimization (ALO) and Salp Swarm Algorithm (SSA). Based on the N-BaIoT dataset, the method reported a 99.9% actual positive rate, besides resolving the issue of high-dimensional feature space.

Hosseini, et al. [15] suggested botnet detection with a hybrid Slime Mold Algorithm (SMA) and SSA for choosing features. The algorithm utilized chaos theory to balance exploration and exploitation, achieving higher detection in UCI datasets. Gharehchopogh, et al. [16] suggested a binary Multi-Objective Dynamic Harris Hawks Optimization (MODHHO) algorithm for choosing features. The algorithm utilized mutation operators and different classifiers (KNN, SVM, MLP, DT), which showed higher speed and accuracy on five datasets.

Alkhamash [17] offered a metaheuristic-based, blockchain-integrated model for DDoS attack detection (MHADMA-BCIDL). The model utilized Arctic Tern Optimization (ATO) for attribute selection and CNN-BiLSTM for classification, achieving 99.32% accuracy on the BoT-IoT dataset. Maghrabi, et al. [18] proposed a hybrid deep learning-based model (BESO-HDLBD) that incorporated Bald Eagle Search Optimization (BESO) for selecting attributes and a CNN-BiLSTM-Attention for bot identification. The model worked best on benchmarked datasets, outperforming existing methods in speed and accuracy.

Maazalahi and Hosseini [19] proposed a hybrid algorithm as a fusion between Whale Optimization Algorithm (WOA), Particle Swarm Optimization (PSO), and Sailfish Optimizer (SFO). The algorithm was tested on BoT-IoT and UNSW-NB15 datasets and achieved a detection accuracy of 99.8% in less execution time. Elsedimy and AboHashish [20] proposed FCM-SWA, an integration between fuzzy C-means clustering and Sperm Whale Algorithm (SWA), for IoT-driven innovative systems. The algorithm outperformed existing methods on BoT-IoT, NSL-KDD, and AWID datasets using adaptive threshold techniques in accuracy and precision.

Despite the good performance encompassed in existing feature selection and classification techniques, as indicated in Table 1, typical weaknesses still hold. Some models incur an inefficient exploration-exploitation balance, leading to convergence in the latter parts or locally optimal sets of features. Others attain good detection accuracy but are computationally costly, especially on high-dimensional or IoT streaming datasets. Most research works provide limited assessment on the impact of the selected feature on the robustness and generalizability of models to unknown attacks. The paper bridges these gaps by introducing EKOAs for feature selection, aiming at achieving computational effectiveness, convergence speed, and high detection accuracy.

3 Materials

3.1 Multi-objective optimization

Multi-objective optimization aims to optimize multiple conflicting objectives, often resulting in trade-offs among them. Improving one objective can result in the deterioration of another, requiring solutions that balance these trade-offs.

Table 1: An overview of related works

Reference	Contribution	Accuracy	TPR	Shortcoming
[13]	Suggested a multi-kernel SVM using GWO for feature selection, achieving 99.72% accuracy with reduced training time.	99.72%	-	Limited evaluation datasets and focus on specific malware types (Cortex A9 samples).
[14]	Developed SSA–ALO hybrid for feature selection, achieving 99.9% TPR on N-BaIoT datasets with superior efficiency.	-	99.9%	High computational complexity for large-scale datasets.
[15]	Introduced SMA + SSA with chaos theory for balanced exploration and exploitation in feature selection.	-	-	Results lack comprehensive comparison with advanced optimization algorithms.
[16]	Presented MODHHO for multi-objective feature selection and versatile classification across multiple datasets.	98.1%	-	Moderate accuracy improvement compared to existing approaches.
[17]	Proposed MHADMA-BCIDL with blockchain integration and CNN-BiLSTM for DDoS detection, achieving 99.32% accuracy.	99.32%	-	Dependence on blockchain may introduce overhead in real-time systems.
[18]	Designed BESO-HDLBD with hybrid deep learning for spatial-temporal feature extraction and botnet detection.	99.4%	-	The computational cost is due to the BiLSTM and attention mechanisms in large datasets.
[19]	Proposed SFO-WOA-PSO-K-means hybrid with 99.8% accuracy and low execution time for botnet detection.	99.8%	-	Limited scalability for highly dynamic IoT environments.
[20]	Introduced FCM-SWA with enhanced clustering and global optimization for IoT-based innovative systems.	98.9%	-	Lacks evaluation on diverse IoT network scenarios and attack types.

The solutions to such problems are termed Pareto-optimal solutions, also known as the Pareto front, in which no objective can be enhanced without compromising at least one other objective. The mathematical formulation of a multi-objective optimization problem can be stated as follows:

$$\min F = \{f_1(X), f_2(X), \dots, f_M(X)\} \quad (1)$$

Subject to:

$$\begin{aligned} g_i(X) &\leq 0, \quad i = 1, 2, \dots, q \\ h_i(X) &\leq 0, \quad j = 1, 2, \dots, p \\ X &\in \Omega \end{aligned} \quad (2)$$

Where $X = \{x_1, x_2, \dots, x_D\}$ represents a decision vector in a D -dimensional space, $g_i(X)$ and $h_i(X)$ represent constraints of inequality and equality, respectively, and F is the set of M objective functions to optimize, Ω defines the feasible decision space.

In multi-objective optimization, a solution $U = \{u_1, u_2, \dots, u_D\}$ is supposed to dominate another solution $V = \{v_1, v_2, \dots, v_D\}$, denoted as $U < V$, if the following conditions are satisfied:

$$\begin{aligned} f_i(U) &\leq f_i(V), \quad \forall i \in \{1, 2, \dots, M\} \\ f_i(U) &< f_i(V), \quad \exists i \in \{1, 2, \dots, M\} \end{aligned} \quad (3)$$

Non-dominated or Pareto-optimal solutions are the ones not dominated by another. They constitute the Pareto front of the problem, which is said to be the best set of conflicting objective trade-off solutions. The feasible solution is said to be satisfying all the constraints and is in the set of non-dominated solutions if and only if it qualifies for the criteria outlined above. The Pareto front, thus, shows all the Pareto-optimal solutions for a problem.

3.2 Feature selection

Feature selection is a crucial step in data classification, where the target is to select a subset of features from the total feature set F_{et} , consisting of D features and N samples, to maximize classification performance while minimizing computational cost [21]. The process can be formulated mathematically as follows:

A feature subset X is represented as a binary vector $X = (x_1, x_2, \dots, x_D)$, where $x_j \in \{0, 1\}$ specifies whether the j^{th} feature is selected ($x_j = 1$) or not ($x_j = 0$). The following equation can then describe the task of feature selection:

$$\max H(X) \quad (4)$$

$H(X)$ represents the objective function that evaluates the classification accuracy of the selected feature subset X .

The classifier's running time is directly proportional to the selected number of features. With a larger set of features, the classifier is computationally costlier and runs slower, but classification accuracy is potentially lower for a smaller set. The compromise between optimizing accuracy and keeping the selected number of features small is thus required. The compromise can be framed as a bi-objective optimization problem:

$$\min F = (ERR(X), |X|) \quad (5)$$

Where $ERR(X) = 1 - H(X)$ is the classification error for the selected feature set X and $|X|$ is the number of selected features.

3.3 Disruption operator

The disruption operator is taken from astrophysical phenomena, which tries to enhance population diversity for optimization methods. Including variation in the population expands the search area and manages exploration and exploitation effectively. The disruption operator successfully enhances the performance of optimization methods to avoid premature convergence. The disruption operator is mathematically represented as:

$$\begin{aligned} D_{op} &= \begin{cases} D_{i,j} \times \delta(-2, 2), & \text{if } D_{i,j,best} \geq 1 \\ 1 + D_{i,j,best} \times \delta\left(-\frac{10^{-4}}{2}, \frac{10^{-4}}{2}\right), & \text{otherwise} \end{cases} \end{aligned} \quad (6)$$

Where $D_{i,j}$ signifies the Euclidean distance between the i^{th} and j^{th} solutions in the population, $D_{i,j,best}$ denotes the Euclidean distance between the i^{th} solution and the best solution identified so far, and $\delta(x, y)$ is a random value generated within the interval $[x, y]$.

The operator dynamically adjusts its impact based on the proximity of solutions to the best-known solution. If $D_{i,j,best} \geq 1$, a larger variation is introduced, allowing for greater exploration in the search space. Otherwise, a minor variation is applied, encouraging fine-tuned exploitation around the best solution. This design ensures that the algorithm strikes a balance between discovering new areas in the search space and refining existing solutions, thereby enhancing overall optimization performance.

4 Methodology

KOA is a physics-inspired metaheuristic algorithm based on Kepler’s laws of planetary motion. These laws define the motion of planets around the sun in elliptical orbits, the relationship between areas swept by the planets, and the proportionality between the square of their orbital period and the cube of their semi-major axis [22]. KOA applies these ideas to mimic optimization such that the planets are treated as potential solutions, while the sun is treated as the best. The algorithm starts by using an initial population of planets characterized by a given orbital eccentricity and spin period. The initialization is specified as below.

$$X_i^j = X_{i,lb}^j + rand \times (X_{i,ub}^j - X_{i,lb}^j), \tag{7}$$

$$i=1, 2, \dots, N; \quad j=1, 2, \dots, D \tag{8}$$

$$e_i = rand, i=1, 2, \dots, N \tag{8}$$

$$OP_i = |rand|, i=1, 2, \dots, N \tag{9}$$

Where D represents the problem's dimensionality, N is the population size, $X_{i,lb}^j$ and $X_{i,ub}^j$ are the lower and upper bounds for the j^{th} variable, and $rand$ is a random number in the interval $[0, 1]$. This reflects the binary nature of the feature selection problem, where each feature can either be included (1) or excluded (0) from the subset. These normalized bounds ensure that the optimization begins within a valid real-valued range before binary conversion via the sigmoid-based transformation.”

The planets rotate around the sun in elliptical orbits, undergoing two phases: moving closer to the sun and moving away. The gravitational force between the sun and a planet, which governs the planet's motion, is calculated as:

$$F_{gi}(t) = e_i \cdot \mu(t) \cdot \frac{M_s \cdot m_i}{R_i^2 + \epsilon} + r_1 \tag{10}$$

Where M_s and m_i represent the normalized masses of the sun and the planet, calculated as follows:

$$M_s = r_2 \cdot \frac{fit_s(t) - worst(t)}{\sum_{k=1}^N (fit_k(t) - worst(t))} \tag{11}$$

$$m_i = \frac{fit_i(t) - worst(t)}{\sum_{k=1}^N (fit_k(t) - worst(t))} \tag{12}$$

Where $\mu(t) = \mu_0 \cdot exp(-\gamma \cdot t/T)$ is the gravitational constant, $R_i = \sqrt{\sum_{j=1}^d (X_{sj}(t) - X_{ij}(t))^2}$ is the distance between the planet and the sun, and r_1 and r_2 are random values, and $\epsilon \in$ is a small constant.

The velocity of a planet, influenced by its distance from the sun, is updated as follows:

$$\vec{v}_i(t) = \begin{cases} \delta \cdot (2r_4 \cdot \vec{X}_i - \vec{X}_b) + \delta' \cdot (\vec{X}_a - \vec{X}_b) + (1 - R_{norm}(t)) \cdot \sigma \cdot \vec{U}_1 \cdot r_5 \cdot (\vec{X}_{i,ub} - \vec{X}_{i,lb}), & R_{norm}(t) \leq 0.5 \\ r_4 \cdot \kappa \cdot (\vec{X}_a - \vec{X}_i) + (1 - R_{norm}(t)) \cdot \sigma \cdot \vec{U}_2 \cdot r_5 \cdot (r_3 \cdot \vec{X}_{i,ub} - \vec{X}_{i,lb}), & \text{otherwise} \end{cases} \tag{13}$$

The position is then updated as follows:

$$\vec{X}_i(t+1) = \vec{X}_i(t) + \sigma \cdot \vec{v}_i(t) + (F_{gi}(t) + |r|) \cdot \vec{U} \cdot (\vec{X}_s(t) - \vec{X}_i(t)) \tag{14}$$

In the second stage, KOA refines the planetary positions around the sun using an adaptive factor h and the exploration formula:

$$\vec{X}_i(t+1) = \vec{X}_i(t) \cdot \vec{U}_1 + (1 - \vec{U}_1) \cdot \left(\frac{\vec{X}_i(t) + \vec{X}_j(t) + \vec{X}_a(t)}{3} + h \cdot \left(\frac{\vec{X}_i(t) + \vec{X}_j(t) + \vec{X}_a(t)}{3} - \vec{X}_b(t) \right) \right) \tag{15}$$

$$h = \frac{1}{e^{\eta r}}, \quad \eta = (l - 1) \cdot r_4 + 1, \quad l = -1 - 1 \cdot \left(\frac{t \% T}{T} \cdot T \right)$$

Balancing exploration (broad search) and exploitation (fine-grained tuning), KOA effectively finds global optimum solutions in high-dimensional search spaces. Its physical inspiration maintains a balanced optimization process that can be applied to various applications.

EKOA is an improvement on the traditional KOA. EKOA addresses the weakness in the first algorithm, i.e., poor convergence for high-dimensional issues, an insufficient balance between exploration and exploitation, and sub-standard handling of complex solution spaces. EKOA achieves higher accuracy, rapid convergence, and solution diversity by utilizing new strategies, i.e., dynamic fine-tuning, oscillatory chaotic force, cross-direction solution creation, and elite-based optimization.

The adaptive adjusting policy dynamically updates the weight between exploration and exploitation in every iteration. At the initial phases, EKOA emphasizes exploration to thoroughly explore the search area. With increasing iterations, the algorithm transfers the focus step by step toward exploitation, adjusting the promising area for the optimum solution. The weighing is mathematically represented as:

$$w = w_{min} + (w_{max} - w_{min}) \cdot \frac{t}{T} \tag{16}$$

Where w_{min} and w_{max} are the minimum and maximum weights, respectively, t stands for the ongoing iteration, and T is the total number of iterations. This adaptability prevents the algorithm from prematurely converging to

local optima, ensuring a more robust search across the solution space.

In traditional KOA, the gravitational constant $\mu(t)$ gradually decreases to focus the search around promising areas. EKOA improves this process by introducing an oscillatory chaotic force constant, which dynamically modulates gravitational force to increase diversity in solutions and prevent stagnation. The gravitational constant is updated as:

$$\mu(t) = s_{map}(t) + \mu_0 \cdot \exp\left(-\frac{\gamma t}{T}\right) \quad (17)$$

Where $s_{map}(t)$ is an oscillatory chaotic function calculated as follows.

$$s_{map}(t+1) = \alpha \cdot \sin(\pi \cdot s_{map}(t)) \quad (18)$$

Where μ_0 is the initial gravitational constant, γ is a decay factor, and T represents the total cycle count. This chaotic mechanism ensures greater randomness in gravitational influence, allowing EKOA to escape local optima and maintain diverse solutions.

The crosswise solution generation strategy accelerates convergence by improving population diversity and generating new candidate solutions. Based on their current positions, two "satellite" solutions are created around existing solutions. The crossover equations are:

$$KX_{a,j}(t) = r_1 \cdot X_{a,j}(t) + (1 - r_1) \cdot X_{b,j}(t) + c_1 \cdot (X_{a,j}(t) - X_{b,j}(t)) \quad (19)$$

$$KX_{b,j}(t) = r_1 \cdot X_{b,j}(t) + (1 - r_2) \cdot X_{a,j}(t) + c_2 \cdot (X_{b,j}(t) - X_{a,j}(t)) \quad (20)$$

Where r_1 and r_2 are random values between $[0, 1]$, and c_1 and c_2 are constants controlling the influence of each parent solution. If the satellite solutions improve the fitness value, they replace their parent solutions, ensuring that the population progressively improves over iterations.

The elite-driven optimization strategy focuses on refining the best solutions to enhance convergence accuracy and precision. It combines three sub-strategies: elite movement, elite cooperation, and elite-driven optimization. Elite movement refines elite solutions by adding perturbations based on their distance from non-elite solutions:

$$\begin{aligned} GX_i(t) &= X_i(t) + A_1 \cdot D_1 + \Delta h \\ A_1 &= l_1 \cdot (r_1 - 0.5) + 1 \\ D_1 &= 2 \cdot r_2 \cdot X_{best}(t) - X_i(t) \end{aligned} \quad (21)$$

Where Δh is a refinement term derived from the Levy flight mechanism.

Elite solutions collaborate by sharing information to improve population diversity:

$$GX_i(t) = X_r(t) + r_1 \cdot D_3 + \Delta h \quad (22)$$

Where $D_3 = X_a(t) - X_b(t)$ and $X_r(t)$ is a randomly selected elite solution.

Elite-driven optimization focuses on aggressively refining elite solutions:

$$GX_i(t) = X_r(t) + A_2 \cdot D_2 + \Delta h \quad (23)$$

Where $D_2 = X_i(t) - X_r(t)$ and A_2 is a randomly selected elite solution.

As shown in Figures 1 and 2, the EKOA workflow begins with the initialization of the population, where the positions, velocities, and orbital parameters of the solutions are set within the defined bounds. The algorithm evaluates the fitness of each solution in light of the optimization objective. In subsequent iterations, the lateral crossover generates new candidate solutions, the oscillatory chaotic force changes the search region, and the adaptive weight allows for smooth movement between exploration and exploitation. Finally, the elite-driven optimization refines the top-performing solutions, consistently enhancing the population's quality. The loop terminates when the terminating criteria are met, e.g., a fixed number of iterations or a convergence point.

The enhanced algorithm is a multi-target anomaly-based IDS for IoT. The algorithm employs Pareto dominance to attain a practical compromise between conflicting objectives, e.g., enhancing classification accuracy but reducing the number of features chosen to be analyzed. Non-dominated solutions are preserved in each iteration through a repository-based structure, which presents different Pareto-optimal solutions. EKOA optimizes leader solutions in each iteration, chosen from the repository through a roulette-wheel selection algorithm in conjunction with hypercube scores and the Boltzmann function. The leader is therefore guaranteed to be a good choice for optimizing the process.

The repository is divided into two components: the grid and the controller. The grid organizes solutions for easier assessment and diversity, and the controller decides whether new solutions are to be added to the repository. To improve the repository's quality, dominated solutions are purged at fixed intervals so that high-quality solutions are retained. This architecture promotes a well-distributed Pareto front for the opposing accuracy and feature reduction objectives.

The algorithm initializes a randomly started population set of solutions and their positions, speeds, and gravity constants. Non-dominated solutions are moved apart in a repository, while dominated members are preserved in the base population. The leader solution is chosen in the initial step of every iteration from the repository through the roulette-wheel technique, under the governance of the Boltzmann function. The EKOA framework applies its mechanisms, including adaptive weight adjustment, oscillatory chaotic force, crosswise solution generation, and elite-driven optimization strategies, to effectively explore and exploit the solution space. Since feature selection is a binary problem, solutions are converted from the discrete domain to the binary domain using Eq. 24.

$$y_i^{j+1} = f(x) \begin{cases} 1, & \text{if } w(x_i^{j+1}) \geq r_{rand} \\ 0, & \text{otherwise} \end{cases} \quad (24)$$

$$w(a) = \frac{1}{1 + e^{10(a-0.5)}}$$

This change ensures the algorithm yields a binary code for the feature subset. The optimization step is followed by adding new non-dominated points to the

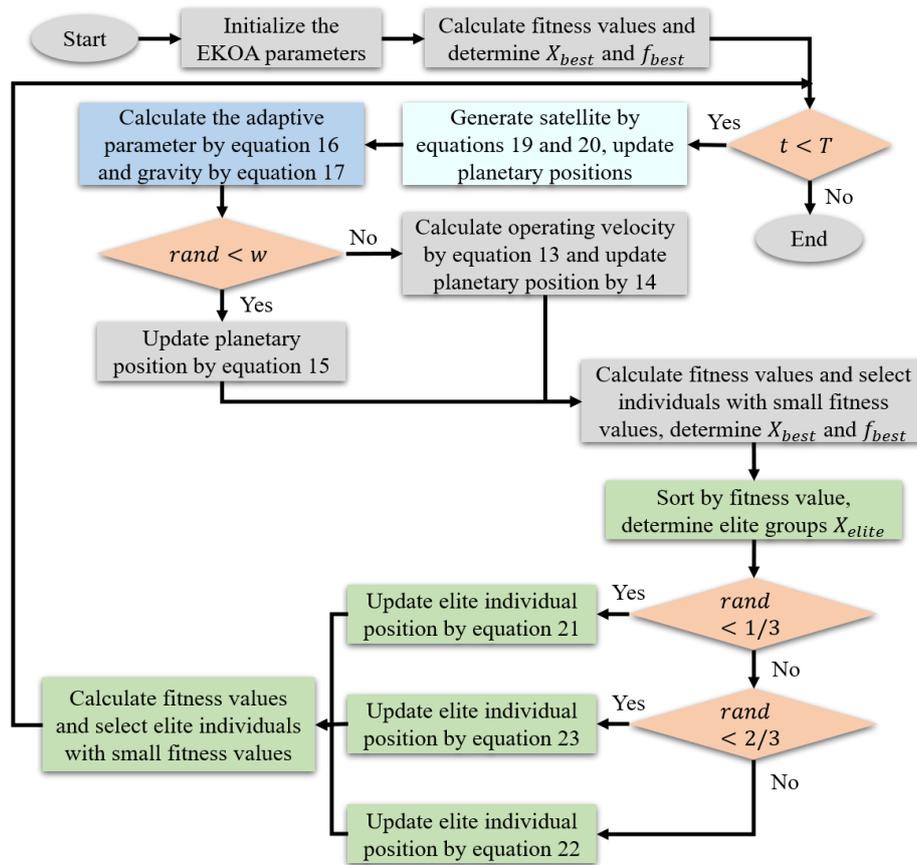


Figure 1: Flowchart of EKO

repository and eliminating dominated ones. If the repository size exceeds, the lesser-quality points are eliminated in favor of higher-quality points. The disruption operator is called at regular intervals to introduce controlled randomness to prevent stagnation. The algorithm terminates when the maximum iterations

are exceeded or a predetermined convergence criteria are met.

5 Results

To evaluate the performance of the proposed EKO-based intrusion detection system, experiments were

Input:

- Population size N
- Maximum number of iterations T
- Search space dimensionality D

Output:

- Optimal solution X_{best}
- Best fitness value f_{best}

Initialize parameters: gravitational constant μ_0 , decay factor γ , and total iterations T

Generate initial population using Eq. 7; assign orbital eccentricity using Eq. 8, and compute orbital period via Eq. 9

Evaluate the initial fitness for each candidate; identify the best individual, and set its fitness; initialize iteration count

Repeat until $t \geq T$:

For each individual $i = 1$ to N :

- Compute gravitational attraction using Eq. 10
- Measure distance to the current best solution using Eq. 12
- Determine velocity using Eq. 13

If a random number $<$ threshold:

Update position with Eq. 14

Else:

Use Eq. 15 for position update

Recalculate fitness; if this new fitness improves upon f_{best} , update both f_{best} and X_{best}

Increment iteration count: $t = t + 1$

End

Figure 2: Pseudo-code

conducted on three datasets: Mirai and Gafgyt. The datasets were normalized and encoded according to standard preprocessing and labeling methods. Feature selection was performed according to the proposed binary multi-objective EKO algorithm. Table 2 shows the complete set of algorithms and classifier hyperparameters, like population size, number of iterations, and crossover coefficients. The levels for these variables were selected after preliminary tuning for stable convergence and desirable search behavior. The 5-fold cross-validation method was chosen for statistical robustness. The experiments were repeated five times, and the means were calculated across performance metrics.

EKO was evaluated on ten feature selection datasets and nine botnet detection datasets, as indicated in Tables 3 and 4. The datasets had at least 100,000 samples, and some even exceeded a million. Table 5 provides the nature of the datasets, e.g., normal/abnormal class ratios.

Mirai botnet attacks were utilized for the training set (70%), while Gafgyt botnet attacks comprised the test set (30%). This was to keep the model robust, as we are training on attacks, we are aware of, but testing on the ability to identify new patterns previously unseen. The test was performed in MATLAB on an Intel Core i5-8400 processor computer, running 8 GB of RAM.

Feature selection experiments compared EKO against five multi-objective algorithms: MOHFOFA [23], NSGA-IIFS [24], B-MOABCFS [25], and MOPSOFS [26]. The Hyper-Volume (HV) and several feature subsets (FN) metrics were used to evaluate

Table 2: Hyperparameter settings for algorithms and classifiers

Component	Parameter	Range	Description
EKO	μ_0	0.1	Initial gravitational constant for attraction force
	γ	15	Gravitational decay factor controlling convergence speed
	w_{min}	0.4	Minimum adaptive weight for exploration
	w_{max}	0.9	Maximum adaptive weight for exploitation
	c_1 and c_2	Uniform [-1, 1]	Crossover coefficients in lateral crossover mechanism
	Population size	30	Number of individuals in the population
	Max iterations	100	Maximum number of optimization iterations
KNN	k-value	3	Number of neighbors used for classification
Decision tree	Max depth	None (default)	Tree expansion continues until full purity or constraint
SVM	Kernel	RBF	Radial basis function kernel for non-linear classification

Table 3: Summary of datasets used for feature selection

Dataset	No. of features	No. of classes	No. of samples
Yale_64	1024	15	165
CNAE-9	857	9	540
LSVT	309	2	126
Musk	167	2	476
Urban land cover	148	9	507
Hill-valley	100	2	606
Sonar	60	2	208
Ionosphere	34	2	351
Vehicle	18	4	846
Vowel	10	11	990

performance. Table 6 presents the HV results for the ten datasets, which measure solution convergence and diversity. Table 7 reports the FN values (average and standard deviation) to assess the effectiveness of dimensionality reduction. EKO's classification relies on K-Nearest Neighbors (KNN) and Leave-One-Out Correlation (LOOCV) scores to measure classification errors. Experiments demonstrated EKO's ability to effectively optimize high classification accuracy feature subsets and outperform traditional multi-objective algorithms.

EKO optimized both the anomaly detection and feature selection for botnet detection. Non-dominated solutions with lowest error rates in each iteration were saved in a second external archive. Table 8 is a comparison between EKO and other algorithms, which indicates that EKO performs better than all the algorithms in all metrics: True Positive Rate (TPR), True Negative Rate (TNR), False Alarm Rate (FAR), accuracy, Area Under the Curve (AUC), and Geometric Mean (G-mean)

Accuracy measures the proportion of correctly labeled records, combining True Negatives (TNs) and True Positives (TPs) over the total population:

$$A = \frac{TN + TP}{FN + FP + TN + TP} \quad (25)$$

FAR evaluates the proportion of False Positives (FPs) among standard samples:

$$FAR = \frac{FP}{TN + FP} \quad (26)$$

TPR or sensitivity quantifies the percent of true positives identified successfully:

$$TPR = \frac{TP}{TP + FN} \quad (27)$$

TNR or specificity determines the percentage of true negatives recorded:

$$TNR = \frac{TN}{TN + FP} \quad (28)$$

G-mean balances sensitivity and specificity, providing a harmonic mean between TPR and TNR. AUC measures the relationship between TPR and FAR over a range of classification thresholds:

$$AUC = \frac{TPR \cdot FAR}{2} + \frac{(1 + TPR) \cdot (1 - FAR)}{2} \quad (29)$$

The consistently superior performance of EKO across datasets is primarily due to its hybrid optimization

Table 4: Summary of datasets used for botnet detection

Dataset	Bashlite (%)	Mirai (%)	Anomaly (%)	Normal (%)	No. of records	No. of features
Ennio doorbell	89	0	89	11	3,55,506	115
Samsung webcam	86	0	86	14	3,75,228	115
Monitoring equipment XC1003	39	59	98	3	8,15,237	115
Monitoring equipment XC1002	37	58	94	6	8,29,079	115
Ecobee thermostat	38	61	98	2	8,35,887	115
Monitoring equipment PT838	37	51	88	12	8,36,902	115
Monitoring equipment PT737	40	52	92	7	8,28,271	115
Danmini doorbell	31	64	95	5	10,18,309	115
Baby monitor	28	55	84	16	10,98,688	115

structure, which is well-aligned with the nature of IoT botnet detection. EKOAs adaptive strategy dynamically shifts the focus from exploration to exploitation, improving convergence without overfitting. Sinusoidal chaotic force addition introduces controlled randomness to enhance population diversity, which is crucial in avoiding a local optimum because of redundancy or noise that is common in high-dimensional data from IoT. The elite-guided aspect also introduces localized optimization for possibly good candidates, in such a way that compact and effective sub-sets of features are chosen, leading to higher accuracy in the classifier.

To assess the proposed framework’s ability to generalize across unseen botnet types, a cross-family evaluation was conducted. Specifically, two scenarios were tested:

- Scenario 1: Training on Mirai samples and testing on Gafgyt samples
- Scenario 2: Training on Gafgyt samples and testing on Mirai samples

These setups simulate real-world IoT environments where the intrusion detection system must detect novel attack variants without prior exposure during training. The results for both scenarios using KNN and DT classifiers are summarized in Table 9.

These results consolidate that the resultant EKOAs-based feature selection method facilitates successful generalizability in new attack patterns. Surprisingly, the

performance is marginally higher for the KNN classifier under domain shift scenarios. The reason is that EKOAs can weed out noise in the datasets and emphasize behavior-centric patterns usable for different families of botnets.

6 Discussion

The proposed EKOAs-based intrusion detection system exhibits clear comparative benefits compared to the diversity of state-of-the-art methods in Table 1. The different methods all contribute to metaheuristic-based feature choice or hybrid detection methods. Nevertheless, EKOAs presents clear performance benefits in various dimensions, such as generalizability, convergence speed, and deployability.

In experiments on typical test datasets such as Mirai and Gafgyt, the EKOAs framework always achieved detection accuracy greater than 99% and reduced the set of features by 35%. This is on par with methods such as GWO-SVM and MHADMA-BCIDL, which performed with high accuracy in narrow-use cases but were evaluated on less inclusive datasets or a few malware types. EKOAs’s consistent performance on diverse attack types, e.g., DDoS, data exfiltration, and command-and-control traffic, shows higher generalizability to new threats.

From an algorithmic point of view, EKOAs addresses several weaknesses characteristic of metaheuristic-based detection systems. Such approaches as SSA-ALO and

Table 5: Distribution of botnet-related classes in training and testing sets

Dataset	Testing set (%)		Training set (%)	
	First class	Second class	First class	Second class
Monitoring equipment XC1003	Gafgyt (95)	Normal (5)	Mirai (96)	Normal (4)
Monitoring equipment XC1002	Gafgyt (85)	Normal (15)	Mirai (92)	Normal (8)
Ecobee thermostat	Gafgyt (94)	Normal (6)	Mirai (96)	Normal (4)
Monitoring equipment PT838	Gafgyt (76)	Normal (14)	Mirai (82)	Normal (18)
Monitoring equipment PT737	Gafgyt (84)	Normal (16)	Mirai (86)	Normal (14)
Danmini doorbell	Gafgyt (85)	Normal (15)	Mirai (91)	Normal (9)
Baby monitor	Gafgyt (65)	Normal (35)	Mirai (78)	Normal (22)

Table 6: HV results for feature selection experiments

HV	MOHHOFOA	B-MOABCFs	NSGA-IIFS	MOPSOFS	EKOAs
	Std/average	Std/average	Std/average	Std/average	Std/average
Yale 64	0.009/0.687	0.003/0.752	0.0135/0.448	0.005/0.645	0.002/0.771
CNAE-9	0.011/0.823	0.011/0.834	0.019/0.487	0.008/0.755	0.006/0.852
LSVT	0.029/0.768	0.072/0.822	0.006/0.404	0.004/0.752	0.025/0.881
Musk	0.005/0.936	0.008/0.942	0.014/0.618	0.022/0.894	0.003/0.957
Urban land cover	0.005/0.884	0.007/0.871	0.024/0.596	0.011/0.836	0.003/0.892
Hill-valley	0.004/0.649	0.021/0.631	0.017/0.531	0.007/0.652	0.002/0.932
Sonar	0.004/0.91	0.005/0.913	0.021/0.699	0.016/0.891	0.003/0.922
Ionosphere	0.002/0.93	0.003/0.927	0.091/0.841	0.003/0.922	0.003/0.944
Vehicle	0.006/0.684	0.006/0.693	0.033/0.613	0.007/0.692	0.003/0.724
Vowel	0.001/0.826	0.001/0.828	0.03/0.815	0.009/0.826	0.009/0.839

Table 7. Feature subset results for feature selection experiments

FN	MOHHOFOA	B-MOABCFS	NSGA-IIFS	MOPSOFS	EKOA
	Std/average	Std/average	Std/average	Std/average	Std/average
Yale_64	1.27/8.45	2.46/12.56	1.91/5.55	1.66/3.45	1.11/10.28
CNAE-9	2.15/8.92	3.66/9.29	2.44/7.41	0.44/5.53	2.66/10.49
LSVT	1.04/5.23	1.25/4.51	0.83/3.39	3.45/5.02	1.26/6.23
Musk	1.73/14.03	2.75/10.86	3.55/6.8	1.42/10.05	1.99/15.74
Urban land cover	2.01/14.05	2.53/11.42	2.56/9.22	1.88/10.55	1.13/14.53
Hill-valley	1.21/9.15	2.76/9.26	1.11/7.02	2.71/8.25	0.45/9.21
Sonar	1.42/11.1	1.76/11.02	1.63/5.81	1.96/10.23	2.28/12.18
Ionosphere	0.71/7.22	0.88/7.26	1.15/5.22	0.71/6.41	0.41/7.56
Vehicle	0.47/5.23	0.36/5.42	0.22/4.13	0.41/5.35	0.31/5.91
Vowel	0/9.01	0/9.01	0.19/8.34	0.44/8.44	0/9

Table 8: Comparative performance analysis of EKOA and other algorithms for botnet detection

Datasets	Algorithms	AUC	G-mean	TPR	TNR	FAR	Accuracy
Monitoring equipment XC1003	MOHHOFOA	0.88	0.87	0.94	0.82	0.18	0.89
	NSGA-IIFS	0.68	0.67	0.84	0.54	0.47	0.69
	B-MOABCFS	0.83	0.82	0.91	0.74	0.26	0.84
	MOPSOFS	0.67	0.65	0.83	0.52	0.48	0.68
	EKOA	0.98	0.98	0.97	0.99	0.08	0.98
Monitoring equipment XC1002	MOHHOFOA	0.89	0.87	0.92	0.86	0.14	0.89
	NSGA-IIFS	0.68	0.67	0.73	0.62	0.39	0.68
	B-MOABCFS	0.74	0.73	0.61	0.87	0.14	0.69
	MOPSOFS	0.62	0.61	0.78	0.48	0.53	0.64
	EKOA	0.98	0.98	0.98	0.98	0.02	0.97
Ecobee thermostat	MOHHOFOA	0.89	0.88	0.94	0.85	0.17	0.9
	NSGA-IIFS	0.72	0.71	0.87	0.57	0.44	0.72
	B-MOABCFS	0.85	0.86	0.91	0.82	0.18	0.87
	MOPSOFS	0.77	0.77	0.78	0.73	0.28	0.78
	EKOA	0.99	0.99	0.99	0.98	0.05	0.98
Monitoring equipment PT838	MOHHOFOA	0.9	0.9	0.95	0.85	0.16	0.91
	NSGA-IIFS	0.76	0.74	0.92	0.58	0.41	0.78
	B-MOABCFS	0.81	0.81	0.9	0.7	0.28	0.82
	MOPSOFS	0.77	0.77	0.86	0.67	0.34	0.78
	EKOA	0.98	0.98	0.96	0.98	0.009	0.98
Monitoring equipment PT737	MOHHOFOA	0.79	0.79	0.92	0.68	0.33	0.82
	NSGA-IIFS	0.65	0.63	0.82	0.49	0.51	0.66
	B-MOABCFS	0.76	0.75	0.88	0.64	0.36	0.78
	MOPSOFS	0.64	0.61	0.84	0.44	0.55	0.65
	EKOA	0.97	0.97	0.98	0.95	0.08	0.97
Danmini doorbell	MOHHOFOA	0.87	0.86	0.93	0.82	0.19	0.88
	NSGA-IIFS	0.66	0.63	0.88	0.44	0.56	0.69
	B-MOABCFS	0.77	0.74	0.93	0.59	0.41	0.84
	MOPSOFS	0.69	0.67	0.87	0.52	0.48	0.71
	EKOA	0.98	0.98	0.97	0.91	0.04	0.98
Baby monitor	MOHHOFOA	0.92	0.92	0.96	0.88	0.12	0.92
	NSGA-IIFS	0.68	0.68	0.76	0.62	0.39	0.67
	B-MOABCFS	0.83	0.83	0.79	0.89	0.11	0.83
	MOPSOFS	0.71	0.71	0.61	0.82	0.18	0.71
	EKOA	0.97	0.97	0.99	0.94	0.06	0.97

MODHHO are typically susceptible to premature convergence or limited diversity in solution space, potentially inhibiting precision or becoming unstable. EKOA integrates four strategic enhancements to overcome such shortcomings:

- Dynamic adjustment strategy: adapts parameters dynamically based on search progress, balancing stably between exploration and exploitation.
- Oscillatory chaotic force: introduces controlled randomness to prevent stagnation and enhance escape from local optima.
- Crosswise Solution Generation: enhances diversity in candidate solutions in later iterations.
- Elite-driven optimization: ensures that high-performance solutions control the evolutionary

process, increasing the likelihood of a global optimum.

Such processes cause EKOA to converge faster than classical evolutionary techniques but without a loss in solution quality. For example, in comparison with BESO-HDLBD and SFO-WOA-PSO, which involve the use of high-level neural structures or multi-level optimization steps, EKOA discovers optimal or near-optimal ensembles of features in fewer iterations and with much less computational expenditure.

Feasibility in practical deployments is of concern for IoT-driven intrusion detection systems, which are usually resource-limited and need a low-latency response. EKOA is suitable for such an environment because:

Table 9: Cross-attack generalization results using EKOA-selected features

Scenario	Classifier	Accuracy	TPR (Recall)	FAR	TNR	G-Mean	AUC
Train: Mirai → Test: Gafgyt	KNN	96.7%	95.6%	4.3%	95.7%	95.6%	95.5%
Train: Mirai → Test: Gafgyt	DT	95.4%	93.7%	5.1%	94.9%	94.3%	94.1%
Train: Gafgyt → Test: Mirai	KNN	97.2%	96.0%	3.8%	96.2%	96.1%	95.9%
Train: Gafgyt → Test: Mirai	DT	95.8%	94.3%	4.7%	95.3%	94.8%	94.6%

- Low execution time: Feature selection using EKOA is computationally lightweight and does not depend on deep learning backbones or large ensemble models.
- Classifier compatibility: The system leverages efficient classifiers (KNN and decision tree), which are known for fast inference times and ease of integration on edge devices.
- Scalability: The modular design allows the framework to be deployed on distributed or hierarchical architectures such as cloud-edge systems, IoT gateways, and embedded devices.

These advantages make EKOA a compelling and high-performance substitute for more advanced or specialized intrusion detection methods. It perfectly balances speed, accuracy, and scalability, the essential properties for real-time IoT network security in high-speed applications.

7 Conclusion

This paper proposed an EKOA-driven optimal IoT security feature selection intrusion detection system. EKOA incorporates adaptive control, chaotic force modulation, cross-sectional solution construction, and elite-based fine-tuning to promote convergence and robustness. Experimental verifications demonstrated higher detection accuracy and reduced feature dimensionality against state-of-the-art contemporary multi-objective methods on standard benchmark sets. Future work will extend the system for real-time intrusion detection based on online learning models. Secondly, realization in realistic-edge scenarios and exploring transfer learning methods between IoT applications will be attempted to enhance adaptability and scalability.

Conflict of interest

The authors declare that they have no conflicts of interest.

References

- [1] K. Halimi, A. Hadjadj, Z. Kouahla, and B. Farou, "A Fuzzy Logic-Driven Semantic and Binary Tree-Based Indexing Framework for Scalable IoT Data Storage and Retrieval," *Informatica*, vol. 49, no. 24, 2025, doi: <https://doi.org/10.31449/inf.v49i24.8039>.
- [2] R. M. Ghadban, H. Z. Neima, and H. A. Jasim, "A Blockchain-Based Security Framework for IoT Networks: Design, Implementation, and Evaluation," *Informatica*, vol. 49, no. 24, 2025, doi: <https://doi.org/10.31449/inf.v49i24.8122>.
- [3] B. Bala and S. Behal, "AI techniques for IoT-based DDoS attack detection: Taxonomies, comprehensive review and research challenges," *Computer science review*, vol. 52, p. 100631, 2024, doi: <https://doi.org/10.1016/j.cosrev.2024.100631>.
- [4] T. Al-Shurbaji *et al.*, "Deep Learning-Based Intrusion Detection System For Detecting IoT Botnet Attacks: A Review," *IEEE Access*, 2025, doi: <https://doi.org/10.1109/ACCESS.2025.3526711>.
- [5] O. Malkawi, N. Obaid, and W. Almobaideen, "Intrusion Detection System for 5G Device-to-Device Communication Technology in Internet of Things," *Informatica*, vol. 48, no. 15, 2024, doi: <https://doi.org/10.31449/inf.v48i15.4646>.
- [6] E. Rivandi and R. Jamili Oskouie, "A Novel Approach for Developing Intrusion Detection Systems in Mobile Social Networks," *Available at SSRN 5174811*, 2024, doi: <https://dx.doi.org/10.2139/ssrn.5174811>.
- [7] A. Heidari and M. A. Jabrael Jamali, "Internet of Things intrusion detection systems: a comprehensive review and future directions," *Cluster Computing*, vol. 26, no. 6, pp. 3753-3780, 2023, doi: <https://doi.org/10.1007/s10586-022-03776-z>.
- [8] N. Mohamed, "Artificial intelligence and machine learning in cybersecurity: a deep dive into state-of-the-art techniques and future paradigms," *Knowledge and Information Systems*, pp. 1-87, 2025, doi: <https://doi.org/10.1007/s10115-025-02429-y>.
- [9] J. Azimjonov and T. Kim, "Stochastic gradient descent classifier-based lightweight intrusion detection systems using the efficient feature subsets of datasets," *Expert Systems with Applications*, vol. 237, p. 121493, 2024, doi: <https://doi.org/10.1016/j.eswa.2023.121493>.
- [10] J. Li, M. S. Othman, H. Chen, and L. M. Yusuf, "Optimizing IoT intrusion detection system: feature selection versus feature extraction in machine learning," *Journal of Big Data*, vol. 11, no. 1, p. 36, 2024, doi: <https://doi.org/10.1186/s40537-024-00892-y>.
- [11] K. Harahsheh, R. Al-Naimat, and C.-H. Chen, "Using Feature Selection Enhancement to Evaluate Attack Detection in the Internet of Things Environment," *Electronics*, vol. 13, no. 9, p. 1678, 2024, doi: <https://doi.org/10.3390/electronics13091678>.
- [12] M. Ahmadi *et al.*, "Optimal allocation of EVs parking lots and DG in micro grid using two-stage GA-PSO," *The Journal of Engineering*, vol. 2023, no. 2, p. e12237, 2023, doi: <https://doi.org/10.1049/tje2.12237>.

- [13] H. Haddadjpajouh, A. Mohtadi, A. Dehghantanaha, H. Karimipour, X. Lin, and K.-K. R. Choo, "A multikernel and metaheuristic feature selection approach for IoT malware threat hunting in the edge layer," *IEEE Internet of Things Journal*, vol. 8, no. 6, pp. 4540-4547, 2020, doi: <https://doi.org/10.1109/JIOT.2020.3026660>.
- [14] R. Abu Khurma, I. Almomani, and I. Aljarah, "IoT botnet detection using salp swarm and ant lion hybrid optimization model," *Symmetry*, vol. 13, no. 8, p. 1377, 2021, doi: <https://doi.org/10.3390/sym13081377>.
- [15] F. Hosseini, F. S. Gharehchopogh, and M. Masdari, "A botnet detection in IoT using a hybrid multi-objective optimization algorithm," *New Generation Computing*, vol. 40, no. 3, pp. 809-843, 2022, doi: <https://doi.org/10.1007/s00354-022-00188-w>.
- [16] F. S. Gharehchopogh, B. Abdollahzadeh, S. Barshandeh, and B. Arasteh, "A multi-objective mutation-based dynamic Harris Hawks optimization for botnet detection in IoT," *Internet of Things*, vol. 24, p. 100952, 2023, doi: <https://doi.org/10.1016/j.iot.2023.100952>.
- [17] M. Alkhamash, "A Metaheuristic Approach to Detecting and Mitigating DDoS Attacks in Blockchain-Integrated Deep Learning Models for IoT Applications," *IEEE Access*, 2024, doi: <https://doi.org/10.1109/ACCESS.2024.3519132>.
- [18] L. A. Maghrabi *et al.*, "Enhancing cybersecurity in the internet of things environment using bald eagle search optimization with hybrid deep learning," *IEEE Access*, vol. 12, pp. 8337-8345, 2024, doi: <https://doi.org/10.1109/ACCESS.2024.3352568>.
- [19] M. Maazalahi and S. Hosseini, "Machine learning and metaheuristic optimization algorithms for feature selection and botnet attack detection," *Knowledge and Information Systems*, pp. 1-49, 2025, doi: <https://doi.org/10.1007/s10115-024-02322-0>.
- [20] E. Elsedimy and S. M. AboHashish, "An intelligent hybrid approach combining fuzzy C-means and the sperm whale algorithm for cyber attack detection in IoT networks," *Scientific Reports*, vol. 15, no. 1, p. 1005, 2025, doi: <https://doi.org/10.1038/s41598-024-79230-4>.
- [21] M. B. Bagherabad, E. Rivandi, and M. J. Mehr, "Machine Learning for Analyzing Effects of Various Factors on Business Economic," *Authorea Preprints*, 2025, doi: <https://doi.org/10.36227/techrxiv.174429010.09842200/v1>.
- [22] M. Abdel-Basset, R. Mohamed, S. A. A. Azeem, M. Jameel, and M. Abouhawwash, "Kepler optimization algorithm: A new metaheuristic algorithm inspired by Kepler's laws of planetary motion," *Knowledge-based systems*, vol. 268, p. 110454, 2023, doi: <https://doi.org/10.1016/j.knosys.2023.110454>.
- [23] B. Abdollahzadeh and F. S. Gharehchopogh, "A multi-objective optimization algorithm for feature selection problems," *Engineering with Computers*, vol. 38, no. Suppl 3, pp. 1845-1863, 2022, doi: <https://doi.org/10.1007/s00366-021-01369-9>.
- [24] T. M. Hamdani, J.-M. Won, A. M. Alimi, and F. Karray, "Multi-objective feature selection with NSGA II," in *Adaptive and Natural Computing Algorithms: 8th International Conference, ICANNGA 2007, Warsaw, Poland, April 11-14, 2007, Proceedings, Part I 8*, 2007: Springer, pp. 240-247, doi: https://doi.org/10.1007/978-3-540-71618-1_27.
- [25] E. Hancer, B. Xue, M. Zhang, D. Karaboga, and B. Akay, "Pareto front feature selection based on artificial bee colony optimization," *Information Sciences*, vol. 422, pp. 462-479, 2018, doi: <https://doi.org/10.1016/j.ins.2017.09.028>.
- [26] B. Xue, M. Zhang, and W. N. Browne, "Particle swarm optimization for feature selection in classification: A multi-objective approach," *IEEE transactions on cybernetics*, vol. 43, no. 6, pp. 1656-1671, 2012, doi: <https://doi.org/10.1109/TSMCB.2012.2227469>.

PSO-JSO: A Hybrid Metaheuristic for Load Balancing in Cloud Computing

Na Li

Department of Information Technology, ZhengZhou Vocational College of Finance and Taxation, Zhengzhou, Henan, 450048, China

E-mail: xinxizhengzhou@126.com

Keywords: cloud computing, load balancing, resource utilization, optimization

Received: March 22, 2025

Cloud computing platforms face growing challenges in efficiently allocating resources and balancing loads due to the dynamic and heterogeneous nature of workloads. This work introduces PSOJSO, an innovative hybrid optimization algorithm that fuses Particle Swarm Optimization (PSO) and Jellyfish Search Optimization (JSO). It presents a dynamic time-control mechanism and adaptive coefficients to balance global exploration and exploitation. Experiments were simulated under a time-sharing scheduling policy using CloudSim 3.0.2 for 2 data centers, eight virtual machines, and 10–100 cloudlets. PSOJSO is compared against five baseline algorithms: ACO, ABC, BA, CSA, and PSO alone. PSOJSO achieved a reduction of up to 25.3% in makespan, a 19.8% reduction in energy consumption, and a 17.5% enhancement in resource utilization, proving its validity for dynamic cloud environments.

Povzetek: Hibridni algoritem PSOJSO (PSO-JSO) je namenjen uravnoteženju obremenitve v računalništvu v oblaku. Z dinamičnim časovnim nadzorom zmanjša makespan, porabo energije in izboljša izkoriščenost virov, kar potrjuje njegovo učinkovitost v dinamičnih oblačnih okoljih.

1 Introduction

Cloud computing offers on-demand computational resources over the Internet [1]. It also introduced scalability and a pay-as-you-use model, allowing for its use across various industries [2]. Large dynamic workloads and configurable Virtual Machine (VM) setups are significant challenges for resource and workload management [3]. This can lead to inefficient resource utilization, bottlenecks, energy waste, and increased operating costs [4]. Hence, creative techniques are necessary to optimize the system's performance and energy efficiency. Correlation research into creative grid energy management indicates adaptive control is essential under uncertain conditions [5]. Moreover, machine learning techniques have been promising for investigating and optimizing dynamic business economics-like systems [6].

Optimal performance in cloud environments is achieved by proper load balancing. Even workload distribution through load balancing avoids the underutilization and overloading of VMs, optimizes response times, reduces energy consumption, and improves resource utilization [7]. An efficient load balancing strategy is essential while maintaining Quality of Service (QoS) and adhering to Service-Level Agreements (SLAs) in dynamic and complex cloud systems [8].

Recently, metaheuristics have emerged as one of the most effective methodologies for solving NP-hard problems, such as load balancing. Nature-inspired metaheuristics provide flexible and robust frameworks for exploring complex solution spaces [9, 10]. Algorithms inspired by nature, such as Particle Swarm Optimization

(PSO) and Jellyfish Search Optimization (JSO), have already demonstrated considerable promise in minimizing makespan and enhancing resource efficiency [11–13]. However, combining global exploration with local exploitation remains an open challenge.

The hybridization of PSO with JSO is intuitively inspired by their complementary strengths. While PSO has been highly efficient in exploring the global solution space, JSO possesses a strong capability of local exploitation to refine the solutions. Such a combination provides an unprecedented way to balance exploration and exploitation. It mitigates some of the drawbacks of standalone implementations, particularly in terms of enhanced performance in complex load-balancing scenarios.

This paper presents the PSOJSO algorithm, combining PSO and JSO for load-balancing in cloud computing. The novelty of this work lies in three main contributions. First, we propose a dynamic hybrid framework that alternates between PSO and JSO phases based on a nonlinear time control function, rather than combining the two heuristics statically. This allows the algorithm to leverage PSO's global exploration in early iterations and JSO's local exploitation in later stages. Second, we incorporate nonlinear, time-varying inertia weights and adaptive cognitive/social coefficients to enhance convergence and prevent premature stagnation. Third, we apply this hybrid PSOJSO formulation to a cloud computing context with multiple conflicting objectives, including makespan, energy efficiency, and resource utilization. To the best of our knowledge, this is the first implementation of such a hybrid algorithm within this specific performance-driven context.

2 Related work

Annie Poornima Princess and Radhamani [14] presented an efficient cloud workload balancing problem with a hybrid metaheuristics optimization approach, using Pigeon-inspired Optimization (PO) and Harries Hawks Optimization (HHO) algorithms. Kruekaew and Kimpan [15] suggested a Multi-objective Artificial Bee Colony Algorithm with Q-learning (MOABCQ) for task scheduling and resource utilization in cloud computing. Thakur and Goraya [16] suggested a hybrid optimization approach using Phasor PSO and Dragonfly algorithms, PPSO-DA, for load balancing and resource management. Ramya and Ayothi [17] developed HDWOA-LBM, a hybrid method combining dingo and whale optimization algorithms to maximize resource utilization and reliability while minimizing makespan.

Narwal [18] presented a credit-based scheduling method integrating Walrus and Lyrebird optimization algorithms called HO-CB-RALB-SA. Gabhane, et al. [19] proposed a hybrid algorithm combining Ant Colony Optimization (ACO) and Tabu Search (TS), called ACOTS, for scalable load distribution. Singhal, et al. [20] suggested the Rock Hyrax (RH)-based algorithm, focusing on dynamic load balancing and energy efficiency using QoS parameters. Karuppan and Bhalaji [21] proposed the African Vultures Algorithm (AVA) that balances server workloads while reducing makespan and energy usage.

In various reviewed studies, as represented in Table 1, several hybrid algorithms are proposed to perform load balancing with significant improvement in performance indicators for makespan, resource utilization, and energy efficiency. Nevertheless, most methods concentrate on specific aspects, such as throughput optimization or resource utilization, and very few provide a holistic balance between exploration and exploitation.

While prior hybrid metaheuristics, such as HDWOA-LBM, MOABCQ, and ACOTS, have achieved improvements in specific metrics, they often fall short in addressing all critical QoS parameters holistically, particularly under dynamic and heterogeneous workloads. Furthermore, many studies lack adaptability mechanisms or focus solely on either exploration or exploitation. Our proposed PSOJSO algorithm addresses these gaps through a combination of techniques: (1) a time-adaptive switching mechanism to balance global and local search phases, (2) nonlinear, time-varying inertia and acceleration coefficients to dynamically modulate the search behavior, and (3) chaotic logistic initialization to improve population diversity and prevent premature convergence. These components, integrated within a single framework, provide a novel and scalable solution for robust load balancing in cloud computing environments.

3 Cloud load balancing

This section provides a two-part explanation of load balancing. The first part presents an overview of cloud load balancing. The primary metrics related to load balancing are discussed in the second part.

3.1 Problem statement

Balancing the load will be effective in the cloud computing environment, leading to better resource utilization and sustained high system performance. This is generally concerned with distributing the arriving tasks among the available VMs to avoid underutilization or overloading of resources. Several benefits can be ensured through load balancing, including minimized energy consumption, shorter response times, and excellent system reliability. However, designing effective load balancing strategies has become very challenging due to the changing and heterogeneous structure of cloud workloads.

Load balancing is also concerned with route selection, minimizing bottlenecks, and preventing overflows. Without a well-thought-out strategy, poor cloud environments are often plagued with performance bottlenecks and even SLA violations. The dynamic assignment of tasks to heterogeneous VMs in cloud data centers presents formidable challenges toward workload imbalance avoidance; inefficient workload allocation will lead to server underperformance and increased operational costs.

Network topology greatly influences how resources are allocated and balanced. The unbalanced system, shown in Figure 1(a), illustrates the outcome of an ineffective load-balancing strategy, where some PMs are overloaded and others underutilized. In contrast, a balanced system, Figure 1 (b), illustrates the benefit of equal workload distribution to ensure all resources are utilized efficiently.

Therefore, the challenge lies in applying the load-balancing strategy dynamically in cloud environments to accommodate variations in workloads and heterogeneity, while aiming to maximize resource utilization, minimize delays, and reduce energy consumption. The paper explores a hybrid algorithm to optimally balance these loads and achieve high system efficiency.

3.2 Problem formulation

Load balancing in cloud computing facilitates even distribution of tasks among available VMs while meeting specific performance targets, including better response time, resource utilization, and energy usage. The formulation of this problem involves defining tasks, VMs, dependencies, and performance metrics that guide the optimization process.

Table 1: An overview of relevant works

Study	Algorithm	Main contribution	Tool/testbed	Metrics	Drawbacks
[14]	HHO+PO	Optimized load balancing and resource utilization	JAVA	Makespan, cost, throughput, and latency	Limited focus on energy consumption and dynamic workload handling
[15]	MOABCQ	Improved task scheduling using reinforcement learning	CloudSim	Makespan, throughput, cost, and imbalance degree	Does not consider energy efficiency or real-time adaptability
[16]	PPSO-DA	Balanced resource allocation across physical machines	CloudSim	Resource usage and load imbalance	High computational complexity and limited energy optimization
[17]	HDWOA-LBM	Enhanced resource utilization and reliability	CloudSim	Throughput, reliability, and makespan	Limited scalability and lack of energy consumption focus
[18]	HO-CB-RALB-SA	Credit-based scheduling for equitable task distribution	CloudSim	Makespan, cost, execution time, and resource usage	Complex implementation and limited metrics evaluation
[19]	ACOTS	Faster file delivery and better multi-resource load balancing	CloudSim	Throughput and task completion time	A narrow focus on throughput, lacking other critical metrics
[20]	RH	Addressed local maxima and improved energy efficiency	CloudSim	Makespan and energy usage	Limited consideration of workload distribution
[21]	AVA	Load balancing with minimized makespan and energy usage	CloudSim	Resource utilization and makespan	Overhead in evaluating selection factors
[22]	DRABC-LB	Dynamic task scheduling based on resource awareness		Resource utilization, response time, makespan, and throughput	Limited adaptability to highly dynamic cloud environments

The tasks in the system are denoted as $T = \{T_1, T_2, \dots, T_n\}$, where each T_i (for $1 \leq i \leq n$) represents a non-preemptive task with a specific set of instructions. These tasks are executed across a set of VMs $VM = \{VM_1, VM_2, \dots, VM_m\}$, where VM_j (for $1 \leq j \leq m$) refers to the j^{th} VM in the system. The VMs may have heterogeneous configurations, reflecting variations in computational power, memory, and storage capacities.

The dependencies between tasks are represented using a Directed Acyclic Graph (DAG), $G = (T, E)$, where T denotes the set of tasks and E signifies the edges between them. An edge $e(T_i, T_j)$ indicates that T_j cannot start execution until T_i is completed. These dependencies ensure the logical execution of tasks and impact the scheduling process. To evaluate the efficiency of a load-balancing mechanism, multiple performance metrics are used, each addressing a critical aspect of system performance:

Energy consumption (E): Energy consumption is modeled based on the power usage of Complementary Metal Oxide–Semiconductor (CMOS)-based microprocessors. The dynamic power is given by:

$$P_e = \alpha cv^2 f \tag{1}$$

Where ac is the dynamic power parameter, v is the supply voltage, and f is the processor frequency. The total energy consumption for all tasks is calculated using Eq. 2.

$$E = \sum_{i=1}^n TCT_i \cdot acv_i^2 f_i \tag{2}$$

Where TCT_i stands for the completion time of i^{th} task, v_i refers to the voltage supplied to the processor on which i^{th} task is executed, and n denotes the number of VMs.

Imbalance degree (I): Eq. 3 quantifies the imbalance in task distribution among VMs.

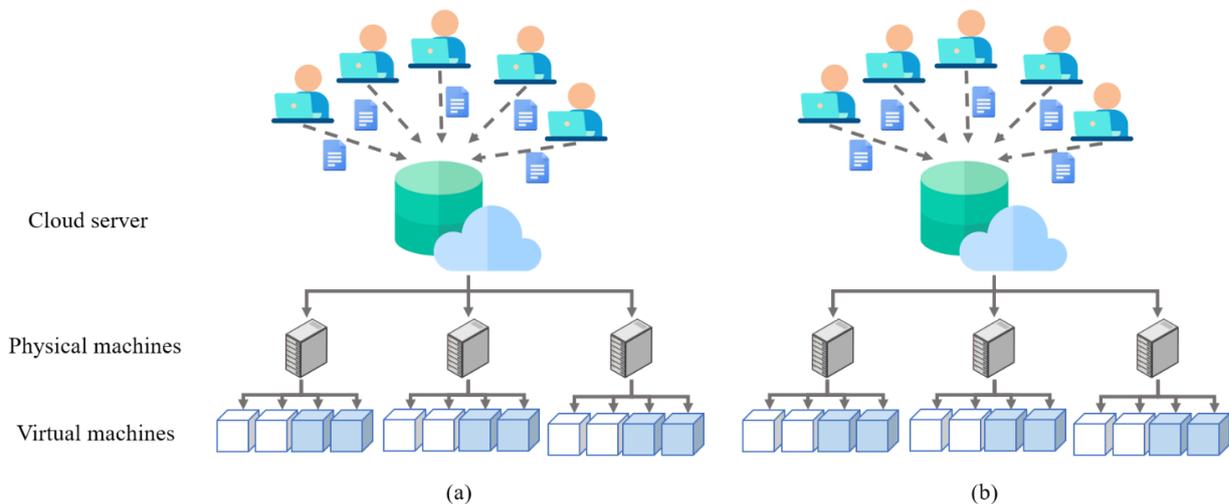


Figure 1: Network topology

$$D = \frac{T_{max} - T_{min}}{T_{avg}} \quad (3)$$

T_{max} and T_{min} are the maximum and minimum completion times across VMs. T_{avg} is the average task length calculated by Eq. 4, where m signifies the VM count.

$$T_{avg} = \frac{\text{Total task length}}{\text{Total VM capacity} \times m} \quad (4)$$

Resource Utilization (RU): Resource utilization measures the percentage of resources used across all VMs and is given by:

$$RU_{VM} = \frac{\sum_{i=1}^n \sum_{j=1}^m TCT_{ij}}{MS \times m} \times 100 \quad (5)$$

Higher resource utilization indicates more efficient use of VM resources.

Response Time (RT): The total response time for all tasks in the system is calculated as:

$$RT = \sum_{i=1}^n \sum_{j=1}^m TCT_{ij} \quad (6)$$

This metric evaluates the cumulative time spent on task execution across all VMs.

Makespan (MS): Makespan is the maximum completion time across all VMs, representing the time required to finish all tasks.

$$MS = \max \{TCT_{VM_j} | j = 1, 2, \dots, m\} \quad (7)$$

Task Completion Time (TCT): The time taken to execute a task T_i on the j^{th} VM is calculated by E. 8.

$$TCT_{ij} = FT(T_i) - ST(T_i) \quad (8)$$

Where $ST(T_i)$ and $FT(T_i)$ are the start and finish times of T_i , respectively. For a specific VM VM_j , the total completion time of all its tasks is calculated as follows:

$$TCT_{VM_j} = \sum_{i=1}^{N_{VM_j}} TCT_{ij} \quad (9)$$

The optimization objective is to diminish the makespan, energy usage, and degree of imbalance while improving resource utilization and minimizing response time. This multi-objective problem is addressed using a fitness function that aggregates all metrics into a single-objective optimization framework as follows:

$$F = \text{minimize} \left(\beta_1 \frac{1}{RU} + \beta_2 MS + \beta_3 E \right) \quad (10)$$

Where β_1 , β_2 , and β_3 are weights assigned to each metric, reflecting their relative importance. The weights are normalized to ensure their sum equals 1. The following principles are applied to simplify the multi-target problem to a single-target one.

- **Scalarization validity:** It ensures linear trade-offs among objectives, making the single-target function properly represent the multi-target problem.
- **Objective alignment:** Maximization metrics, like resource utilization, are converted into minimization objectives to unify the fitness function.
- **Normalization:** Metrics are scaled to a standard range for meaningful aggregation.
- **Weight assignment:** Equal weights ($\beta_1 = \beta_2 = \beta_3 = 1/3$) are assigned initially, ensuring no metric is prioritized.

To guide the design and evaluation of the proposed PSOJSO algorithm, the following research questions (RQs) are formulated:

- RQ1: Can a hybrid PSO-JSO algorithm with dynamic time-control outperform existing metaheuristic methods in reducing makespan under variable cloud workloads?
- RQ2: Does PSOJSO lead to improved energy efficiency and resource utilization compared to standalone and hybrid baselines?
- RQ3: How does the use of nonlinear time-varying coefficients and chaotic initialization impact task scheduling accuracy and load balancing quality?
- RQ4: Is the PSOJSO framework scalable and adaptable to dynamic cloud environments with heterogeneous VM configurations and fluctuating task demands?

Although the scalarization method used in Eq. 10 simplifies multi-objective optimization by combining objectives into a single fitness function, it inherently assumes linear trade-offs between competing goals such as energy consumption and response time. This simplification is suitable for the current study's scope, which emphasizes algorithmic efficiency and comparative evaluation. Nonetheless, we acknowledge that Pareto-based multi-objective optimization could provide a more expressive treatment of conflicting objectives. Integrating a Pareto-dominance approach into PSOJSO is identified as a promising direction for future enhancement of this work.

4 System model

This part discusses the suggested load balancing scheme and the strategy for resolving the problem.

4.1 Proposed model

The cloud load balancing system in this approach utilizes a framework with multiple entities based on CloudSim. A multi-component task scheduling and load distribution mechanism optimizes energy consumption. The entire system processes tasks submitted by users through a centralized control mechanism, comprising several sub-entities, such as a Cloud Information Service (CIS), brokers, and VM managers.

The proposed model integrates multilevel components for cloud load balancing, ensuring efficiency in task scheduling, dynamic redistribution of loads, and energy-aware management. The proposed solution will be implemented using the CloudSim framework, allowing user-submitted tasks to be effectively allocated to the best resources.

The broker treats the tasks submitted by any user. The broker uses a CIS that maintains a continuously updated cloud data center database. The CIS records resource availability, configuration, and performance characteristics. Once the broker receives tasks, it queries CIS to find appropriate data centers. After selecting data centers, it finally allocates the tasks to them. The CIS plays a vital role in this process, ensuring that resource utilization is optimized while directing tasks to data centers with available capacity.

The arriving tasks are handled hierarchically within each data center, as shown by the multiple layers in Figure 2. The data center controller is responsible for executing the task. The scheduler executes tasks on a specific virtual machine according to predefined algorithms, such as Round Robin. After this, the load balancer monitors the dynamic task distribution and detects situations of under- or overloading. It finally controls the VM manager to handle the operational states, whether active or asleep, of the VMs according to the workload.

The load balancer plays an essential role in the system, distributing the workload evenly among the VMs. If the workload of a VM exceeds the upper threshold, the load balancer will remove tasks from it and redistribute them to underloaded VMs for optimal task allocation. If the workload of any VM exceeds the upper threshold, the load balancer removes tasks from it and reallocates them to underloaded VMs for optimal task allocation.

The hierarchical model ensures dynamic adaptability and scalability, effectively managing workloads and

operational states of VMs while maintaining low energy consumption and optimizing resource utilization. This proposed model balances task execution at data centers effectively, providing a strong and energy-efficient solution for modern cloud environments.

4.2 Proposed method

The PSOJSO approach aims to combine the strengths of PSO and JSO by addressing their respective shortcomings. For example, PSO features global exploration capability and performs well in the early stages of the exploration process. In contrast, the JSO relies heavily on local exploitation to refine any solution. The PSOJSO approach balances the two algorithms, enabling the exploration and exploitation of load balancing in cloud computing. Novel elements are nonlinear time-varying weights, adaptive coefficients, and a time control mechanism that dynamically switches between PSO and JSO phases to ensure robust and efficient optimization. Detailed explanations of the components and mathematical formulations are given below.

PSO draws inspiration from the swarming tendency of birds and fishes. It was proposed as an iterative optimization method. In PSO, candidate solutions have been represented as a population of particles in an algorithm. Every particle updates its position in the search space concerning its best-known previous position and the global best position of the swarm. The update for the velocity of every particle should be made using Eq. 11.

$$V_i^{t+1} = \omega \cdot V_i^t + c_1 \cdot r_1 \cdot (P_{best}^t - X_i^t) + c_2 \cdot r_2 \cdot (G_{best}^t - X_i^t) \tag{11}$$

Where w refers to the inertia weight controlling the balance between exploration and exploitation, c_1 and c_2 are cognitive and social coefficients, representing the particle's own experience and the swarm's collective knowledge, and r_1 and r_2 are random values in the range $[0, 1]$.

The positions of particles are updated using Eq. 12.

$$X_i^{t+1} = X_i^t + V_i^{t+1} \tag{12}$$

JSO models the dynamics of jellyfish within an ocean, including motions related to ocean currents or local swimming motions. The position of each jellyfish is updated based on ocean currents or movements within the swarm [23]. Initial positions of the individual jellyfish can be generated using a chaotic logistic map as follows:

$$X_i = LB + (UB - LB) \cdot L_i \tag{13}$$

Where L_i is a logistic map value, and UB and LB are the upper and lower bounds of the search space, respectively.

The position update based on the ocean current is defined as follows:

$$X_i^{t+1} = X_i^t + r_1 \cdot (X^* - 3 \cdot r_2 \cdot X_i^t) \tag{14}$$

Where X^* refers to the current global best position.

Movements within the swarm include passive motion and active motion, calculated as follows:

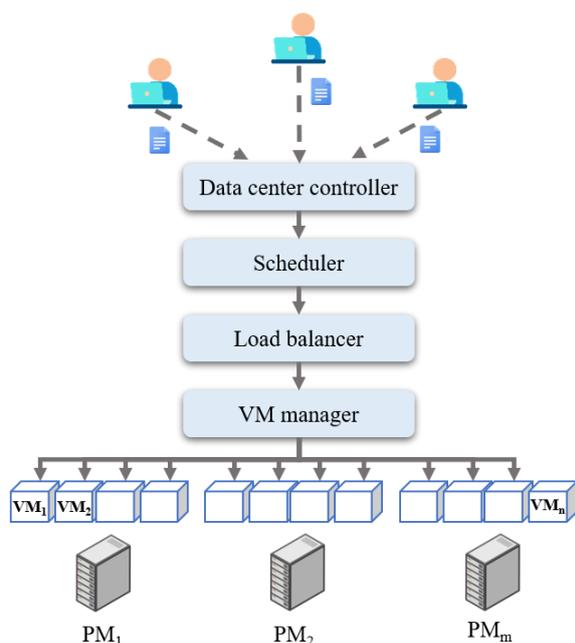


Figure 2: System model

$$X_i^{t+1} = X_i^t + \gamma \cdot r_1 \cdot (UB - LB) \quad (15)$$

$$X_i^{t+1} = X_i^t + \omega \cdot r_1 \cdot \overrightarrow{Step} \quad (16)$$

Where γ is the motion coefficient and \overrightarrow{Step} is determined by the direction of movement, depending on the quality of solutions.

PSOJSO alternates between PSO and JSO phases using a dynamic time control function $c(t)$, determining whether the algorithm follows the PSO or JSO phases.

$$c(t) = \left(1 - \frac{t}{T}\right) \cdot (2r - 1) \quad (17)$$

Where t refers to the current iteration, and T refers to the total number of iterations. If $c(t) \geq 0.5$, the algorithm executes the PSO phase, favoring exploration. Otherwise, it switches to the JSO phase for exploitation.

To enhance exploration in the early stages and exploitation in the later stages, nonlinear time-varying inertia weights (ω), cognitive coefficients (c_1), and social coefficients (c_2) are introduced:

$$\omega = \omega_{min} + (\omega_{max} - \omega_{min}) \left(1 - \frac{t}{T}\right)^\beta \quad (18)$$

$$c_1 = c_{min} + (c_{max} - c_{min}) \sin\left(\frac{\pi}{2} \cdot \left(1 - \frac{t}{T}\right)\right) \quad (19)$$

$$c_2 = c_{min} + (c_{max} - c_{min}) \cos\left(\frac{\pi}{2} \cdot \left(1 - \frac{t}{T}\right)\right) \quad (20)$$

Where β controls the rate of inertia weight decay, and ω_{min} , ω_{max} , c_{min} , and c_{max} define the lower and upper bounds of the weights.

In implementation, the control function $c(t)$ governs dynamic switching between phases. When $c(t) \geq 0.5$, the PSO phase is executed; when $c(t) < 0.5$, the algorithm enters the JSO phase. Thus, the first half of the iterations generally emphasize exploration, and the second half focuses on exploitation. The parameters were set as follows: $\omega_{max} = 0.9$, $\omega_{min} = 0.4$, $c_1 = c_2 = 2$, and $\beta = 1.5$. These values were selected based on recommendations from prior hybrid metaheuristic literature and validated through empirical testing. A limited sensitivity analysis, varying each parameter within $\pm 20\%$, demonstrated that the selected configuration maintained consistent performance and yielded superior results in terms of makespan and energy consumption.

PSOJSO uses the global search capability of PSO during the early stages and local exploitation by JSO during the later stages. Switching between PSO and JSO maintains a balance between exploration and exploitation, thus preventing early convergence. Such hybridization allows PSOJSO to adapt dynamically to the search space and optimizes cloud computing load balancing.

The algorithm initializes the position of all individuals in the population by using a chaotic logistic map to introduce randomness, thereby improving diversity and preventing convergence. Meanwhile, it sets the velocity of each particle or jellyfish to zero. The execution of the PSO

or JSO phases is decided based on the time control mechanism. $c(t)$ denotes a time control function indicating whether the PSO or JSO phases should be executed.

The computational complexity of the PSOJSO algorithm can be estimated as $O(N \times D \times T)$, where N is the number of particles, D is the dimensionality of the problem, and T is the number of iterations. This is consistent with the standard complexity of PSO and JSO, although PSOJSO introduces a slight overhead due to its dynamic switching mechanism and time-varying coefficient updates. However, this additional cost is marginal relative to the performance gains.

5 Results

The proposed load-balancing algorithm draws inspiration from nature-inspired optimization techniques. It involves identifying overloaded servers, selecting appropriate servers for task migration, and optimizing the fitness function to balance workloads across cloud environments.

The proposed algorithm was implemented in CloudSim version 3.0.2 on Windows 10. CloudSim is a well-known simulation tool that models and simulates cloud computing environments. This work utilized the tool to create virtual machines, data centers, and cloudlets, thereby realistically modeling a cloud computing scenario. The experimental environment has four data centers and workload resources with increasing task counts. Table 2 presents the details of the experimental environment configuration. The time-sharing policy was used for resource allocation in this architecture by the VMs across the four data centers.

While the overall system model is built upon CloudSim's default architecture, key components have been extended to support dynamic and energy-aware load balancing. The default round-robin policy was replaced with a custom task allocation policy governed by the PSOJSO fitness function. Additionally, the broker and load balancer were modified to dynamically monitor VM utilization and trigger task migration when workload thresholds are breached. These enhancements allow the system to respond adaptively to fluctuating workloads and to optimize key QoS metrics in real time.

Table 2: Simulation variables

Entities	Variables	Values
Data center	Count	2
VM	CPU	4
	Policy	Time-sharing
	MIPS	1000
Host	Operating system	Windows
	RAM	2 GB
	VMs	8
	Bandwidth	512
	Storage	40 GB
Cloudlets	RAM	4 GB
	Hosts	4
	Length	500-10000
User	Cloudlets	10-100
PSOJSO	Iterations	100
	Population size	50
Simulation	Random seed	42
	Runs per scenario	30

The parameters used in the PSOJJO algorithm were selected through empirical calibration and informed by prior literature on PSO and JJO hybrids. For the PSO phase, the inertia weight ω was varied nonlinearly from 0.9 to 0.4 using a decay exponent $\beta = 0.1$, while cognitive (c_1) and social (c_2) coefficients varied within [0.5, 2.5] using complementary sinusoidal schedules to encourage early exploration and late exploitation. These bounds were chosen based on commonly accepted best practices in swarm optimization. The motion coefficient γ for JJO was fixed at 0.1, ensuring a controlled range of swarm movement. The time control function $c(t)$ plays a critical role by governing the phase transition; its design promotes exploration during early iterations and gradual transition to exploitation.

To analyze its effectiveness, the proposed hybrid PSO-JSO algorithm was considered for various key QoS parameters, including makespan, energy efficiency, and response time. Makespan measures the total time to complete a set of tasks and reflects the scheduler’s efficiency in allocating resources under dynamic workloads. It directly correlates with throughput and infrastructure utilization. Energy consumption is crucial due to the environmental and operational cost implications of large-scale data centers, where inefficient load distribution results in unnecessary power usage. Response time represents the delay experienced by users from task submission to completion, making it vital for meeting SLA requirements and user satisfaction.

Furthermore, the proposed method was tested against several prevalent optimization algorithms, including the Cuckoo Search Algorithm (CSA), Bat Algorithm (BA), ABC, PSO, and ACO. Results related to energy consumption, depicted in Figures 3 and 4, demonstrate that PSOJJO can reduce power consumption by dynamically adjusting the number of active VMs in response to workload demands. In this regard, resources are utilized efficiently during high workload periods, while idle servers are transitioned into a low-power state to minimize unnecessary energy consumption. This methodology not only enhances energy efficiency but also reduces operational costs.

Figures 5 and 6 depict that PSOJJO achieves lower and more consistent response times compared to the baseline methods, as observed across varying VM counts. In the proposed algorithm, the scheduler will be continuously updated in real-time by the broker about server availability to ensure that tasks are scheduled optimally. This reduces waiting time in server queues and maximizes the accuracy of estimated response times, thereby increasing user satisfaction and enabling the handling of more critical tasks.

Figures 7 and 8 illustrate the relationship between the number of cloudlets and the makespan. PSOJJO maintains a relatively stable makespan with increasing tasks, whereas methods such as ACO and PSO degrade their effectiveness with increased workloads. This is because the algorithm optimizes resource utilization and balances the workload effectively, thereby preventing excessive delays in task execution.

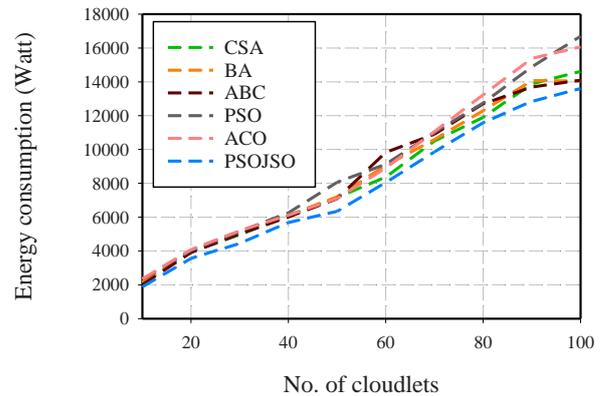


Figure 3: Energy consumption results for 10 VMs

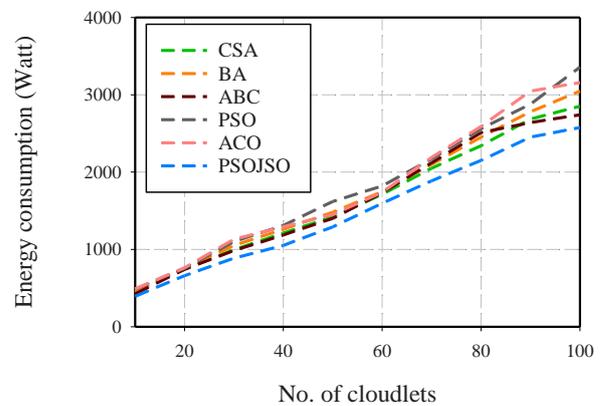


Figure 4: Energy consumption results for 50 VMs

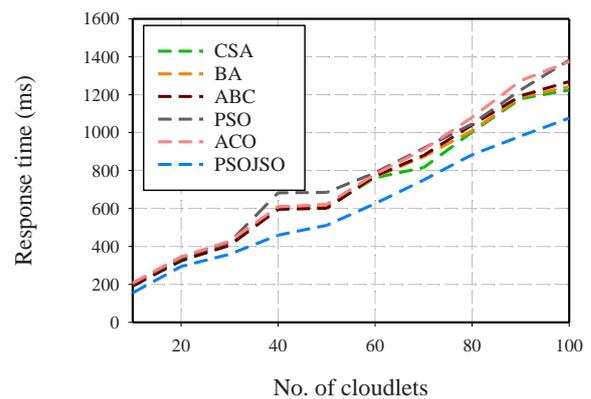


Figure 5: Response time results for 10 VMs

PSOJJO aims to enhance system performance by dynamically scaling the number of VMs up or down according to workload demand. This makes the model adaptable, ensuring sufficient resources to cope with high demands, thereby avoiding bottlenecks and slowdowns. During low times, it deactivates unnecessary VMs, minimizing energy wastage. These measures contribute to

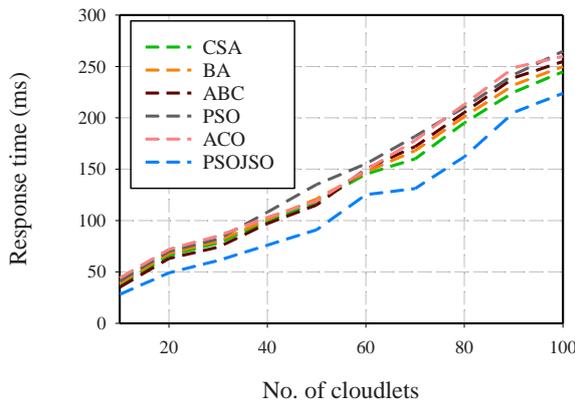


Figure 6: Response time results for 50 VMs

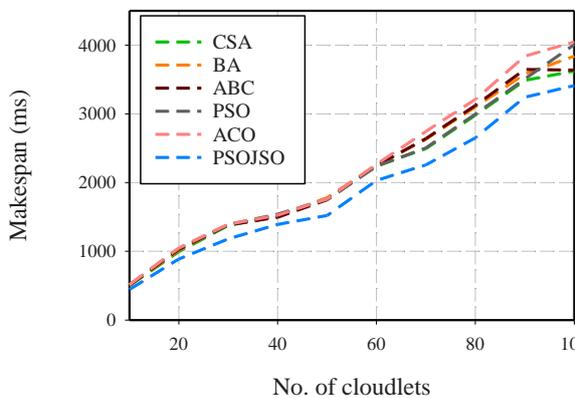


Figure 7: Makespan results for 10 VMs

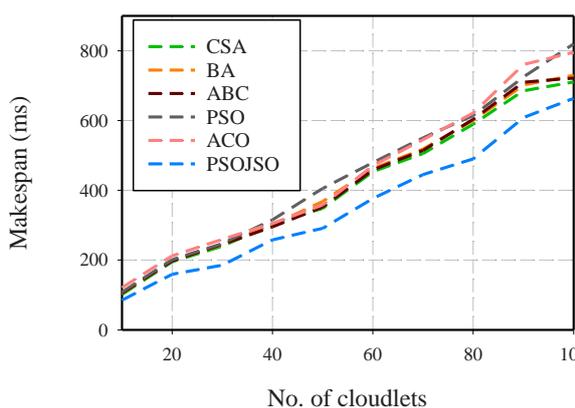


Figure 8: Makespan results for 50 VMs

a greener computing philosophy by saving energy while optimizing resource utilization.

6 Discussion

The proposed PSOJSO algorithm demonstrated superior performance in critical metrics, such as makespan, energy efficiency, resource utilization, and response time, when

compared to widely used metaheuristics, including ACO, PSO, ABC, BA, and CSA. These improvements can be directly attributed to three key features.

- **Dynamic time-control mechanism:** This component enables the algorithm to dynamically shift between exploration (PSO phase) and exploitation (JSO phase) phases adaptively, thereby preventing early convergence, a common issue in standalone PSO, and improving solution diversity throughout the optimization process.
- **Nonlinear time-varying coefficients:** The inertia weight and acceleration coefficients are modulated over time using nonlinear decay functions. This ensures that early iterations favor global search, while later stages intensify local search, enhancing convergence stability and precision.
- **Chaotic logistic initialization:** By seeding the population using a chaotic map, the algorithm ensures a well-distributed initial search space, reducing the risk of stagnation and improving coverage in high-dimensional problems.

Compared to the state-of-the-art hybrid algorithms, such as HDWOA-LBM, MOABCQ, PPSO-DA, and ACOTS, the developed PSOJSO method shows better exploration and exploitation balance, mainly under dynamic and diverse cloud workload conditions. For one, although MOABCQ uses reinforcement learning to promote scheduling, there is no adaptive parameter control, and energy efficiency is not tackled directly. HDWOA-LBM prioritizes reliability and throughput but is unsuitable for scalability and optimizing energy. PPSO-DA and ACOTS have balanced resource allocation but high computational complexity or restricted metric consideration (e.g., throughput only). PSOJSO, however, introduces a nonlinear time-varying design of the coefficients, a chaotic initialization method, and dynamic switching between PSO and JSO phases so that PSOJSO can achieve better results along several QoS dimensions, including makespan, energy consumption, response time, and imbalance degree. These innovations enable PSOJSO to generalize and scale better under multi-resource settings and complete the gaps of other techniques.

PSOJSO maintains robust performance as the number of tasks increases, with makespan remaining stable even as workload intensity grows. This suggests good scalability in multi-VM, multi-data center environments. Additionally, since the method does not rely on problem-specific heuristics, it is generalizable to other cloud computing scenarios such as fog/edge computing or task offloading in IoT.

7 Conclusion

This paper presented the PSOJSO hybrid load-balancing algorithm for resource allocation problems and workload management for cloud computing environments. This integrated solution utilizes the exploration of PSO and exploitation of JSO to achieve an optimal trade-off, thereby preventing premature convergence and enhancing solution quality. It also utilized dynamic scaling, adaptive coefficients, and nonlinear inertia weights for improved

resource utilization and efficiency. PSOJSO obtained stable makespan without increasing workloads, maximized throughput by reducing VM idle time, and decreased energy consumption due to adaptive resource scaling. A detailed performance evaluation based on CloudSim confirmed that the designed PSOJSO algorithm outperformed existing algorithms, including BA, PSO, ABC, ACO, and CSA, in all key QoS metrics, namely, makespan, response time, and energy efficiency. These enhancements would result in lower operational costs, higher user satisfaction, and improved system reliability.

References

- [1] A. Gou, "Cloud-Computing-Enabled Transformer Architecture for the Design of Functional Clothing Structures," *Informatica*, vol. 49, no. 11, 2025, doi: <https://doi.org/10.31449/inf.v49i11.6763>.
- [2] S. Sun, J. Dong, Z. Wang, X. Liu, and L. Han, "An on-demand collaborative edge caching strategy for edge-fog-cloud environment," *Computer Communications*, vol. 228, p. 107967, 2024, doi: <https://doi.org/10.1016/j.comcom.2024.107967>.
- [3] K. Saidi and D. Bardou, "Task scheduling and VM placement to resource allocation in Cloud computing: challenges and opportunities," *Cluster Computing*, vol. 26, no. 5, pp. 3069-3087, 2023, doi: <https://doi.org/10.1007/s10586-023-04098-4>.
- [4] W. Yao, Z. Wang, Y. Hou, X. Zhu, X. Li, and Y. Xia, "An energy-efficient load balance strategy based on virtual machine consolidation in cloud environment," *Future Generation Computer Systems*, vol. 146, pp. 222-233, 2023, doi: <https://doi.org/10.1016/j.future.2023.04.014>.
- [5] A. Kermani *et al.*, "Energy management system for smart grid in the presence of energy storage and photovoltaic systems," *International Journal of Photoenergy*, vol. 2023, no. 1, p. 5749756, 2023, doi: <https://doi.org/10.1155/2023/5749756>.
- [6] M. B. Bagherabad, E. Rivandi, and M. J. Mehr, "Machine Learning for Analyzing Effects of Various Factors on Business Economic," *Authorea Preprints*, 2025, doi: <https://doi.org/10.36227/techrxiv.174429010.09842200/v1>.
- [7] C. Vijaya and P. Srinivasan, "Multi-objective meta-heuristic technique for energy efficient virtual machine placement in cloud data centers," *Informatica*, vol. 48, no. 6, 2024, doi: <https://doi.org/10.31449/inf.v48i6.5263>.
- [8] S. Ghafir, M. A. Alam, F. Siddiqui, and S. Naaz, "Load balancing in cloud computing via intelligent PSO-based feedback controller," *Sustainable Computing: Informatics and Systems*, vol. 41, p. 100948, 2024, doi: <https://doi.org/10.1016/j.suscom.2023.100948>.
- [9] B. Pourghebleh, A. Aghaei Anvigh, A. R. Ramtin, and B. Mohammadi, "The importance of nature-inspired meta-heuristic algorithms for solving virtual machine consolidation problem in cloud environments," *Cluster Computing*, vol. 24, no. 3, pp. 2673-2696, 2021, doi: <https://doi.org/10.1007/s10586-021-03294-4>.
- [10] M. Ahmadi *et al.*, "Optimal allocation of EVs parking lots and DG in micro grid using two-stage GA-PSO," *The Journal of Engineering*, vol. 2023, no. 2, p. e12237, 2023, doi: <https://doi.org/10.1049/tje2.12237>.
- [11] Z. Jafari, A. Habibizad Navin, and A. Zamanifar, "Task scheduling approach in fog and cloud computing using Jellyfish Search (JS) optimizer and Improved Harris Hawks optimization (IHHO) algorithm enhanced by deep learning," *Cluster Computing*, pp. 1-25, 2024, doi: <https://doi.org/10.1007/s10586-024-04347-0>.
- [12] Y. Gao, B. Yang, S. Wang, G. Fu, and P. Zhou, "A multi-objective service composition method considering the interests of tri-stakeholders in cloud manufacturing based on an enhanced jellyfish search optimizer," *Journal of Computational Science*, vol. 67, p. 101934, 2023, doi: <https://doi.org/10.1016/j.jocs.2022.101934>.
- [13] K. Shao, H. Fu, and B. Wang, "An efficient combination of genetic algorithm and particle swarm optimization for scheduling data-intensive tasks in heterogeneous cloud computing," *Electronics*, vol. 12, no. 16, p. 3450, 2023, doi: <https://doi.org/10.3390/electronics12163450>.
- [14] G. Annie Poornima Princess and A. Radhamani, "A hybrid meta-heuristic for optimal load balancing in cloud computing," *Journal of grid computing*, vol. 19, no. 2, p. 21, 2021, doi: <https://doi.org/10.1007/s10723-021-09560-4>.
- [15] B. Kruekaew and W. Kimpan, "Multi-objective task scheduling optimization for load balancing in cloud computing environment using hybrid artificial bee colony algorithm with reinforcement learning," *IEEE Access*, vol. 10, pp. 17803-17818, 2022, doi: <http://dx.doi.org/10.1109/ACCESS.2022.3149955>.
- [16] A. Thakur and M. S. Goraya, "RAFL: A hybrid metaheuristic based resource allocation framework for load balancing in cloud computing environment," *Simulation Modelling Practice and Theory*, vol. 116, p. 102485, 2022, doi: <https://doi.org/10.1016/j.simpat.2021.102485>.
- [17] K. Ramya and S. Ayothi, "Hybrid dingo and whale optimization algorithm-based optimal load balancing for cloud computing environment," *Transactions on Emerging Telecommunications Technologies*, vol. 34, no. 5, p. e4760, 2023, doi: <https://doi.org/10.1002/ett.4760>.
- [18] A. Narwal, "Resource Utilization Based on Hybrid WOA-LOA Optimization with Credit

- Based Resource Aware Load Balancing and Scheduling Algorithm for Cloud Computing," *Journal of Grid Computing*, vol. 22, no. 3, p. 61, 2024, doi: <https://doi.org/10.1007/s10723-024-09776-0>.
- [19] J. P. Gabhane, S. Pathak, and N. M. Thakare, "A novel hybrid multi-resource load balancing approach using ant colony optimization with Tabu search for cloud computing," *Innovations in Systems and Software Engineering*, vol. 19, no. 1, pp. 81-90, 2023, doi: <https://doi.org/10.1007/s11334-022-00508-9>.
- [20] S. Singhal *et al.*, "Energy Efficient Load Balancing Algorithm for Cloud Computing Using Rock Hyrax Optimization," *IEEE Access*, 2024, doi: <http://dx.doi.org/10.1109/ACCESS.2024.3380159>.
- [21] A. S. Karuppan and N. Bhalaji, "Efficient load balancing strategy for cloud computing environment with African vultures algorithm," *Wireless Networks*, pp. 1-17, 2024, doi: <https://doi.org/10.1007/s11276-024-03810-5>.
- [22] R. Mishra and M. Gupta, "DRABC-LB: A Novel Resource-Aware Load Balancing Algorithm Based on Dynamic Artificial Bee Colony for Dynamic Resource Allocation in Cloud," *SN Computer Science*, vol. 5, no. 2, p. 233, 2024, doi: <https://doi.org/10.1007/s42979-023-02570-x>.
- [23] J.-S. Chou and D.-N. Truong, "A novel metaheuristic optimizer inspired by behavior of jellyfish in ocean," *Applied Mathematics and Computation*, vol. 389, p. 125535, 2021, doi: <https://doi.org/10.1016/j.amc.2020.125535>.

A Personalized Music Intervention Framework for Elderly Mental Health Using SWPSI-KNN and Neural Collaborative Filtering Based on EEG Signals

Minyong Zhang

Art and Sports Department, Henan College of Transportation, No.259 Tonghui Road, Zhengzhou City, Henan Province, China

E-mail: awb2023@126.com

Keywords: mental health of elderly people, EEG signals, K-nearest neighbor algorithm, neural collaborative filtering, music intervention

Received: April 11, 2025

Under the trend of global aging, psychological problems such as depression and anxiety are becoming increasingly prevalent among the elderly population. Traditional intervention methods suffer from lagging emotional recognition and insufficient personalization. To improve the mental health problems of the elderly, this study innovatively combines the optimization of the temporal characteristics of EEG signals (power spectral density, time-frequency analysis, dynamic time regularization) with a collaborative recommendation mechanism. A music electrical signal data acquisition system for the psychological health of elderly people based on dynamic analysis of EEG signals has been developed. The system employs real-time EEG acquisition (1000Hz sampling rate), preprocessing (1-50Hz bandpass filtering, ICA-based noise removal), and feature extraction, utilizing an enhanced K-Nearest Neighbor (KNN) algorithm (with sliding windowing and dynamic weight adjustment) to predict EEG responses under music intervention. Experiments involved 80 elderly participants from a nursing home, with datasets including baseline anxiety scales and EEG recordings, validated through randomized controlled trials. The results indicated that the model reduced the EEG tracking error (MAE) from the traditional KNN of 3.24 μV to 0.07 μV . The NCF mechanism achieved 93.2% accuracy in anxiety state classification. In practical applications, the anxiety relief efficiency reached 96.21%, compared to 72.5% in the control group, and the user satisfaction score was 9.5/10. By dynamically optimizing temporal features through dynamic time warping and real-time EEG feedback-driven music adjustment, the system enables personalized intervention, offering an innovative solution combining real-time monitoring and precision adjustment for elderly mental health.

Povzetek: Raziskava predlaga hibridni okvir SWPSI-KNN-NCF za personalizirano glasbeno intervencijo in izboljšanje duševnega zdravja starejših, ki združuje EEG signale in globoko učenje. Model SWPSI-KNN izboljša sledenje EEG (MAE 0,07 μV), medtem ko NCF omogoča razvrščanje anksioznosti in učinkovito lajša anksioznost.

1 Introduction

Against the backdrop of accelerating global population aging, the issue of Mental Health of the Elderly (MHOE) is becoming increasingly prominent [1]. According to the Blue Book of "China Aging Development Report 2024 - Mental Health Status of Chinese Elderly", 26.4% of the elderly in China have varying degrees of depression symptoms, of which 6.2% have moderate to severe depression symptoms. In addition, 23.76% of the elderly experience varying degrees of loneliness, with 4.75% of them frequently feeling lonely. These psychological problems not only seriously affect the quality of life of the elderly but may also exacerbate the development of chronic diseases and increase the medical burden [2]. Traditional intervention methods rely heavily on questionnaire surveys and scale evaluations, which have limitations such as lagging emotion recognition and insufficient dynamic monitoring, making it difficult to meet personalized needs [3]. In this context, non-

pharmacological interventions such as music therapy have attracted much attention due to their low-risk and high acceptance characteristics [4]. Music can improve cognitive function and alleviate negative emotions by activating the limbic system and dopamine secretion, especially in the elderly population [5]. In recent years, the application of music intervention in MHOE has gradually deepened. Wang et al. used partial least squares structural equation modeling technique to quantitatively analyze the data to explore the impact of music intervention on the mental health of college students. Emotional intelligence, as a regulatory factor, significantly and positively regulates the relationship between music education and students' mental health [6]. De Witte et al. conducted a multilevel meta-analysis to evaluate the strength of the impact of music therapy on physiological and psychological stress-related outcomes. They found that music therapy has a significant overall impact on stress-related outcomes [7]. Vajpeyee et al.

proposed a music therapy combined with yoga training to alleviate the mental health stress of medical staff. This study has demonstrated that this method can effectively improve the mental health of medical staff and reduce work pressure [8].

The K-Nearest Neighbor (KNN) algorithm, as a data instance analysis-based algorithm, is often used to analyze mental health data. However, the traditional KNN algorithm has problems such as high computational complexity, so many scholars have made improvements to it. Pamungkas et al. proposed a KNN-based mental health disorder diagnosis model optimized by case-based reasoning to provide more accurate and effective mental health treatment plans. The model achieved an accuracy of 84.62% on the test data [9]. Wibowo et al. proposed a voting classifier that integrates KNN, Gaussian Naive Bayes, and the Random Forest algorithm for effective diagnosis of bipolar disorder and depression. The accuracy of this classifier ranged from 66.67% to 91.67%, which was superior to traditional psychiatric diagnosis methods [10]. Cheng et al. proposed a mental illness prediction model using ensemble logistic regression, KNN, and random forest to address the issues of over-detection or under-detection in traditional mental illness prediction methods. The model exhibited good generalization ability on the prediction test set data, with an accuracy of 83.23%, a recall rate of 89.87%, and a precision of 78.02% [11].

In other aspects, Kolenik T et al. proposed a computational psychotherapy system that simulates Theory of Mind to address the lack of advanced prediction and behavior change mechanisms in existing computational psychotherapy systems. The experiment showed that the system outperformed the current optimal system in terms of prediction accuracy and intervention effect [12]. Kolenik T proposed a digital mental health technology framework based on the Internet of Things (IoT) to address the worsening of mental health issues and insufficient professional resources. This method achieved precise evaluation and personalized intervention through multi-modal data monitoring of physiology, behavior, etc. Experiments have shown that it can effectively improve intervention effectiveness [13].

However, the above methods still have some shortcomings, as some studies rely on subjective evaluations and lack objective physiological data support, resulting in incomplete and inaccurate results. Some studies, although using KNN for data analysis, still have limitations in capturing emotional changes in real-time and cannot reflect the effects of music intervention promptly. To address these challenges, this study aims to achieve two main objectives. One is to improve the prediction accuracy of Electroencephalogram (EEG) signals by developing a Stacked Window Power Spectral Intensity KNN (SWPSI-KNN) algorithm that enhances temporal resolution and noise robustness. The second is to enhance personalized music intervention through Neural Collaborative Filtering (NCF) integration, achieving dynamic alignment between neural responses and music attributes. Based on this, this study proposes a hybrid psychological health music intervention model

(SWPSI-KNN-NCF). This study innovatively introduces the SWPSI feature extraction method to improve the accuracy of emotion prediction, and combines multi-scale time-frequency domain feature matrices with adaptive denoising processing methods to solve the noise interference of dynamic signals. Furthermore, a dual-channel decision framework driven by NCF is constructed to accurately recommend suitable music, thereby enhancing the precision and sustainability of music intervention effects.

2 Method

2.1 Improved KNN algorithm for MHOE detection

The global aging population is intensifying, and the issue of MHOE is becoming increasingly prominent. Emotional disorders such as depression and anxiety are more common in the elderly population, seriously affecting their quality of life and physical and mental health [14]. Music intervention, as a non-pharmacological therapy, can provide psychological support and emotional comfort to the elderly by regulating their emotions. The music intervention MHOE and traditional emotion detection methods are shown in Figure 1.

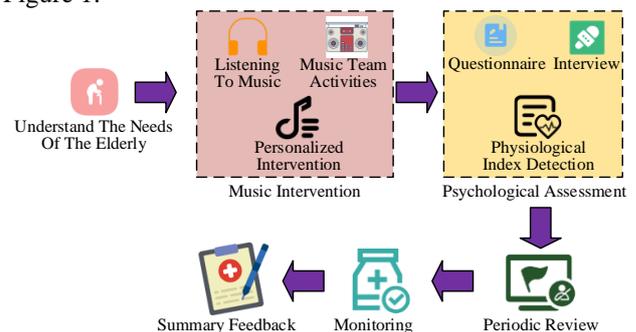


Figure 1: Music intervention for the MHOE and traditional emotion detection methods

In Figure 1, the first step is to communicate with the elderly and their families to understand their interests, music preferences, lifestyle habits, etc., to choose suitable music and activity forms. Music intervention is based on the preferences and needs of the elderly. Music intervention methods include passive listening, structured group activities, and personalized music intervention methods [15]. Traditional emotion detection methods mainly rely on questionnaire surveys and scale evaluations, collecting data through self-reports by subjects to indirectly reflect emotional states. After music intervention, it is necessary to evaluate the effectiveness of music intervention through methods such as scales, observations, and interviews, to understand the changes and degree of improvement in emotions, cognition, socialization, and other aspects of the elderly. Finally, it is required to summarize the entire music intervention process, provide feedback on the intervention effect to the elderly and their families,

collect their opinions and suggestions, and provide a reference for subsequent intervention work.

However, traditional emotion detection methods are relatively limited, mainly collecting data through questionnaires, which makes it difficult to reflect real emotions in real time. To improve the effectiveness of emotion detection, achieve more effective music intervention, and enhance MHOE level, it is urgent to use advanced algorithms for rapid emotion recognition. The KNN algorithm is a classic non-parametric machine learning method that learns based on instances and does not require strict assumptions about the probability distribution of emotional data [16]. The KNN algorithm, with its non-parametric assumption advantage in data distribution, can accurately classify the psychological health status of elderly people based on their physiological and behavioral data. This algorithm has been widely used in the field of MHOE detection. However, traditional KNN algorithms have weak capabilities in feature extraction and selection, and cannot automatically extract discriminative features from raw data, requiring additional preprocessing and feature engineering steps [17]. The traditional KNN algorithm has shortcomings in feature extraction and high-dimensional data processing. Therefore, this study proposes two key improvements, namely the SWPSI-KNN algorithm. Firstly, a sliding window mechanism is introduced to dynamically update the neighborhood range and enhance real-time performance. Secondly, dynamic weight adjustment is adopted to allocate weights based on feature importance, thereby improving prediction accuracy. The SWPSI feature extraction method captures the time-frequency characteristics and complexity of EEG signals by calculating their power spectral intensity and fractal dimension, forming high-dimensional feature vectors. SWPSI can reflect the dynamic changes of EEG signals, while fractal dimension quantifies the complexity of signals. These feature vectors provide richer information for the KNN algorithm, enabling it to more accurately distinguish different emotional states. The structure of SWPSI-KNN is shown in Figure 2.

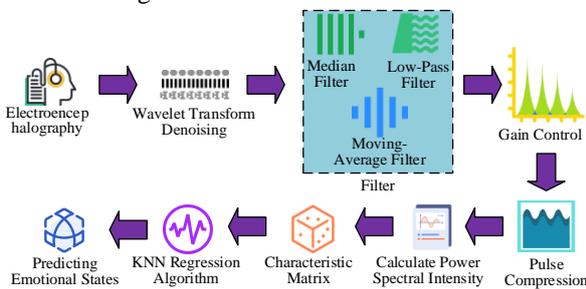


Figure 2: Structure of SWPSI-KNN algorithm

In Figure 2, the algorithm first preprocesses the EEG signal, including denoising and filtering operations, to improve signal quality. By using a heatmap to visualize the feature matrix, the energy distribution and signal complexity changes of different frequency components can be clearly observed. Signal denoising uses wavelet transform method to decompose the signal into different

scales and positions, and the calculation formula is shown in equation (1) [18].

$$y_{after} = wavelet_denoise(x, threshold) \quad (1)$$

In equation (1), y_{after} is the denoised signal and x is the original signal. $threshold$ is a threshold used to control the denoising parameters, while $wavelet_denoise$ is the wavelet transform denoising coefficient. The filtering operation adopts sliding average filtering, and the formula is shown in equation (2) [19].

$$y(n) = \frac{1}{N} \sum_{k=0}^{N-1} x(n-k) \quad (2)$$

In equation (2), $y(n)$ is the value of the filtered signal at the n -th sampling point. k is half the length of the window. $x(n-k)$ is the value of the original signal at the $n-k$ -th sampling point. N is the window length. For salt and pepper noise and pulse noise, this study uses median filtering method for calculation, as shown in equation (3).

$$y(n) = median(x(n-k), x(n-k+1), \dots, x(n+k)) \quad (3)$$

In equation (3), $median$ is the median filtering parameter. k is half the length of the window. $x(n-k)$ to $x(n+k)$ are the sampling point values within the window. High frequency noise is calculated using low-pass filtering, as shown in equation (4).

$$y(n) = \sum_{k=0}^N h(k) \times x(n-k) \quad (4)$$

In equation (4), $h(k)$ is the impulse response of the filter. To ensure that the signal is within the dynamic range, this study uses automatic gain control technology to adjust the gain of the received signal, as shown in equation (5).

$$y[n] = \frac{x[n]}{mean(x[n])} target_gain \quad (5)$$

In equation (5), $y[n]$ is the signal after gain control. $x[n]$ is the input signal before gain. $mean(x[n])$ is the mean of the signal. $target_gain$ is the target gain. To improve signal resolution, this study applies matched filters to the echo signals of each distance unit using pulse compression, as shown in equation (6).

$$y_z[n] = match_filter(x_z[n], h_z[n]) \quad (6)$$

In equation (6), $y_z[n]$ represents the compressed signal parameters. $x_z[n]$ represents the signal parameters before compression. $h_z[n]$ represents the impulse response of the matched filter. $match_filter$ represents pulse compression. The signal received by pulse compression is the processed signal, whose amplitude is within the dynamic range of the system, ensuring that the signal is not distorted and can be

effectively processed. Next, the algorithm divides the processed EEG signal into multiple windows and calculates the power spectral intensity of different frequency bands for each window, forming a time-frequency domain feature matrix. The formula for power spectral intensity is shown in equation (7).

$$PSI_f = \sum_{t=1}^T |X(f,t)|^2 \quad (7)$$

In equation (7), PSI_f is the power spectral intensity at frequency f . $X(f,t)$ is a frequency domain signal obtained through fast Fourier transform. T is the length of the time window. These feature matrices are stacked to form the input feature sequence. The calculation of introducing SWPSI feature extraction method to improve the accuracy of emotion prediction is shown in equation (8).

$$PSI_{stacked} = \begin{pmatrix} PSI_1 \\ PSI_2 \\ \vdots \\ M_i \\ \vdots \\ PSI_n \end{pmatrix} \quad (8)$$

In equation (8), $PSI_{stacked}$ represents the stacked power spectral intensity matrix. M_i is the power spectral intensity vector of the i -th window. n is the number of windows. The SWPSI feature extraction method traverses EEG signals by sliding fixed length windows and calculates the power spectral intensity of each window using fast Fourier transform. Then, these intensity vectors are superimposed into a comprehensive feature matrix to capture the temporal variation of EEG signal activity. This matrix serves as a key input for emotion prediction, as different emotional states are associated with different EEG activity patterns on each frequency band, enabling the KNN algorithm to distinguish different emotional responses based on the extracted SWPSI features. Subsequently, the KNN algorithm is used to predict emotions in the feature sequence. By calculating the distance between the input instance and all instances in the training set, KNN instances are found. The emotional state is predicted based on the output values of these neighbors. By visualizing the distribution and classification boundaries of sample points in the feature space through scatter plots, the classification logic and performance of the model can be demonstrated. The distance formula is shown in equation (9).

$$distance(x_i, x_j) = \sqrt{\sum_{k_w=1}^{d_w} (x_{i,k_w} - x_{j,k_w})^2} \quad (9)$$

In equation (9), x_i and x_j are two different data points. d_w is the dimension of the feature. x_{i,k_w} and x_{j,k_w} are the values of the data point in the k -th dimension. The entire process, from data preprocessing to feature extraction, and then to emotion prediction, forms a complete algorithm framework suitable for MHOE music intervention.

2.2 MHOE data acquisition method based on EEG signals

This study proposes a complete SWPSI-KNN algorithm suitable for MHOE music intervention through preprocessing, feature extraction, and emotion prediction steps. Among them, the SWPSI method directly echoes the characteristic of "selecting music for different emotional states" in music therapy. To implement the application of this algorithm, it is necessary to collect EEG signals from elderly people as a data source for detecting mental health. The intervention of music on MHOE data and the process of obtaining emotional data are shown in Figure 3.

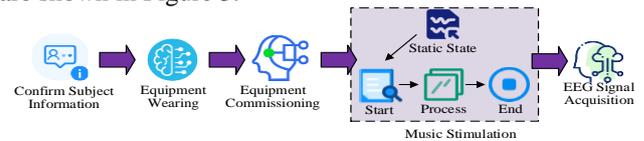


Figure 3: Music intervention on MHOE data and process for obtaining emotional data

In Figure 3, the subjects cover different genders, music preferences, and listening habits, ensuring that the data are diverse and applicable to a wide range of elderly populations. Participants need to fill out a questionnaire on music preferences and listening time before the experiment begins to determine their music choices. Subsequently, they put on EEG signal acquisition devices, and the staff have to debug the devices to ensure that they are functioning properly and EEG signals can be collected normally. Next, static EEG signals are collected, and different music is played according to the preferences of different subjects. Music stimulation adopts a standardized emotional induction paradigm, with the specific process being played in the order of neutral music (2 minutes), negative music (2 minutes), and positive music (2 minutes). Each piece of music is in a resting state with a 30-second interval, and the volume is uniformly calibrated to a 75 dB sound pressure level measurement (A-weighted). This design avoids residual emotional interference by changing the polarity gradient of emotions (neutral, negative, positive) while using music clips of equal duration to ensure comparability of EEG signal time dimensions. The EEG signals before, during, and after the playback are recorded to complete the collection of EEG signals. The fractal dimension can be used to represent the complexity of time-domain signals in different time periods, as shown in equation (10).

$$H_m(k_j) = \frac{N_s - 1}{\sum_{k=1}^{N_s - m_j} \sum_{i=1}^k |s(m_j + ik_j) - s(m_j + (i-1)k_j)|} \quad (10)$$

In equation (10), when calculating the fractal dimension, for the time-domain signal sequence $s(n_s)(n_s = 1, 2, \dots, N_s)$, different starting points m_j and time intervals k_j are selected to extract the sub-

sequence $s(m_j + (i - 1)k_j)$. For each starting point m_j and time interval k_j , the length eigenvalue $H_{m_j}(k_j)$ is calculated to measure the complexity of the signal's fluctuations at that starting point and time interval. In this way, the complexity of signals at different time periods and time scales can be systematically analyzed. $\hat{\lfloor \rfloor}$ represents rounding down operation. The fractal dimension characteristics of EEG are obtained using the Higuchi algorithm, and the time series update is shown in equation (11).

$$x_{m_j}^{k_j} = \{x(m_j), x(m_j + k_j), x(m_j + 2k_j), \dots, x(m_j + \frac{\hat{\lfloor N_s - m_j \rfloor}}{k_j} k_j)\} \quad (11)$$

In equation (11), $x_{m_j}^{k_j}$ is the updated time series. $x(m_j)$ represents the sub-sequence extracted from the starting point m_j of the EEG signal time series. For each new time series, the formula for curve length is shown in equation (12).

$$L_{m_j}(k_j) = \frac{\sum_{i=1}^{\frac{\hat{\lfloor N_s - m_j \rfloor}}{k_j}} |x(m_j + ik_j) - x(m_j + (i-1)k_j)|}{\frac{\hat{\lfloor N_s - m_j \rfloor}}{k_j}} \quad (12)$$

In equation (12), $L_{m_j}(k_j)$ is the length of the curve. The average length calculation is shown in equation (13).

$$\overline{L(k_j)} = \frac{1}{k_j} \sum_{m_j=1}^{k_j} L_{m_j}(k_j) \quad (13)$$

In equation (13), $\overline{L(k_j)}$ is the average length of the time series. On the experimental validation subset, a combination of grid search and cross validation is used to optimize and determine the parameters involved in equations (1) to (13), such as wavelet denoising threshold and filter window half length. The purpose is to ensure that the model achieves optimal performance when processing EEG signals. The actual situation of EEG signal acquisition is shown in Figure 4.

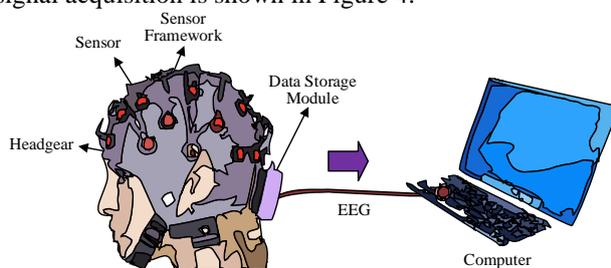


Figure 4: Actual situation of EEG signal acquisition

In Figure 4, EEG signals are obtained by placing multiple electrodes on the scalp surface to capture the electrical activity of brain neurons. The EEG acquisition device used in this study is BioSemi ActiveTwo, which

has 128 channels, a sampling rate of 1024 Hz, and is equipped with active electrodes. According to the international 10-20 system, electrodes are placed at specific scalp locations. After the subjects wear the device, EEG signals are collected. The device converts analog signals into digital signals and transmits them to the computer through Ethernet. During the collection process, subjects usually need to maintain a resting state or complete specific tasks to ensure the stability and reliability of the signal. The collected EEG signals will be monitored and recorded in real-time for subsequent analysis and processing. The entire process, from device preparation to signal acquisition, forms a complete EEG signal acquisition framework, providing important data support for the study of MHOE.

2.3 Establishment of a psychological health music intervention model

In response to the demand for dynamic capture of physiological signals in MHOE evaluation, this study constructs a multidimensional dynamic data acquisition framework based on EEG signals. Real-time monitoring of emotional fluctuations before and after music intervention is achieved through standardized experimental procedures. This study combines fractal dimension feature extraction with the Higuchi algorithm to establish a dynamic correlation between music stimulation and neural response. This data system not only covers high-frequency sampling of physiological signals, but also integrates behavioral characteristics such as subjects' music preferences. This provides a multi-modal data foundation with temporal and individual differences for the psychological health music intervention model. NCF utilizes multi-layer neural networks to capture the nonlinear relationship between users and music, achieving personalized music recommendations for different user groups and enhancing the effectiveness of music intervention. Based on this, this study further integrates dynamic EEG features with personalized intervention logic to construct an SWPSI-KNN-NCF model that integrates SWPSI-KNN and NCF, as shown in Figure 5.

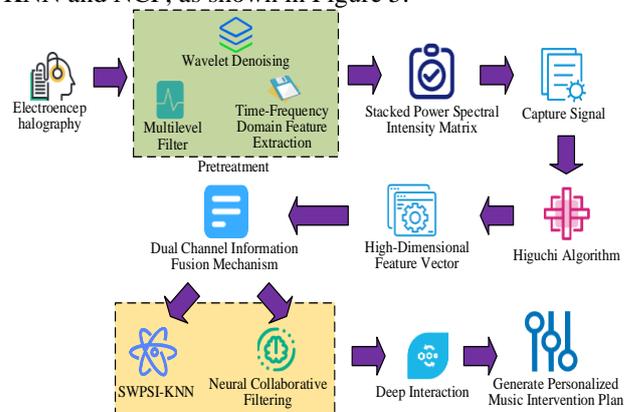


Figure 5: Structure of SWPSI-KNN-NCF

In Figure 5, the SWPSI-KNN-NCF model constructs a complete closed-loop system from EEG signal analysis to personalized music recommendation through multi-modal data fusion and deep feature interaction. The model first standardizes the EEG signals after wavelet denoising (eliminating high-frequency noise) and multi-level filtering (suppressing baseline drift and power frequency interference) through signal preprocessing. Secondly, the time-frequency domain feature extraction stage utilizes short-time Fourier transform to generate the power spectral density matrix. It combines pulse compression and SWPSI to construct temporal features and capture the dynamic evolution of emotions. Subsequently, the SWPSI-KNN emotion prediction module classifies the feature matrix based on an improved distance metric function and outputs probabilities of emotions such as depression and anxiety. The NCF recommendation stage maps user attributes and music metadata through an embedding layer. It uses attention mechanism to align emotional probabilities with music features and generates personalized music recommendations through multi-layer perceptrons. Finally, dynamic feedback optimization is used to collect EEG data in real-time after intervention (such as changes in gamma wave synchronization). Model parameters are updated through online learning to form a closed-loop intervention. This model achieves end-to-end precise intervention from signal processing to dynamic recommendation through multi-level denoising, SWPSI temporal modeling, and cross modal alignment (semantic matching of physiological signals and music attributes). The operation of the music recommendation module in SWPSI-KNN-NCF is shown in Figure 6.

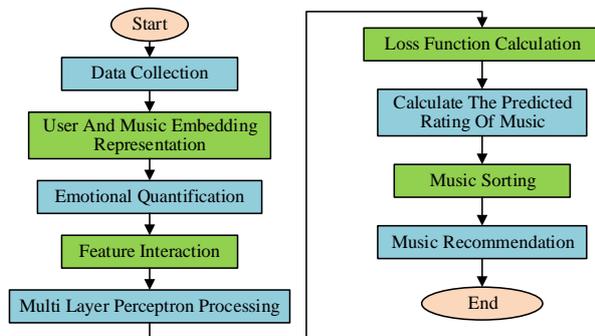


Figure 6: Running process of music recommendation module

In Figure 6, the module first needs to collect users' historical music listening records, preference information, and social network data, and then represent users and music as low-dimensional embedding vectors. Then, the user embedding vector and the music embedding vector are concatenated to form a new feature vector. The concatenated feature vectors are input into a multi-layer perceptron. A multi-layer perceptron consists of multiple Fully Connected Layers (FCLs), each of which is followed by a ReLU activation function. For example, the output of the first layer is shown in equation (14).

$$h_1 = \text{ReLU}(W_1 Z + b_1) \tag{14}$$

In equation (14), h_1 , b_1 , and W_1 are the output vector, bias vector, and weight matrix of the first layer FCL. ReLU represents the activation function. Z is a feature vector formed by concatenating user and music embedding vectors. After multiple transformations, the output formula of the final output layer is shown in equation (15).

$$\hat{r}_{ua} = \sigma(W_L h_{L-1} + b_L) \tag{15}$$

In equation (15), \hat{r}_{ua} is the predicted rating of user u for music a . σ is the Sigmoid function of the output layer. W_L and b_L are the weight matrix and bias vector of the last layer FCL. h_{L-1} is the output vector of the second to last layer FCL. These parts together form the core of the SWPSI-KNN-NCF model. This enables it to effectively learn the non-linear relationship between users and music in the music recommendation module, thereby providing more accurate personalized music recommendations. Finally, the music recommendation module sorts the music based on its rating and completes the music recommendation. During this process, the model dynamically adjusts the cross modal weight matrix in the attention mechanism based on real-time collected EEG emotion classification results (such as anxiety probability values). For example, when anxiety is detected, the recommendation weight of soothing music is automatically enhanced. The user embedding vector is incrementally updated to enable the recommendation list to adapt to the fluctuations of the current emotional state in real-time.

In the preprocessing stage of the SWPSI-KNN-NCF model, the EEG signal is first filled with KNN missing values ($n_neighbors=5$). The dimensional differences in neural features such as θ wave power and γ wave synchronization are eliminated through standardization. At the same time, music preferences (genre, duration) are extracted as behavioral feature vectors. The SWPSI stacking program adopts a two-stage architecture. In the first stage, the dynamic weight KNN ($k=15$) is used to calculate the nearest neighbor classification probability of EEG features. In the second stage, the cross-modal attention module (PSI-Weight) is used to generate the correlation weight matrix between music metadata and EEG features. The KNN prediction probability (dimension 15) is concatenated with the attention weight (dimension 32) to form a stacked feature input NCF. NCF music metadata construction integrates traditional music ontology attributes (regional cultural labels, pentatonic modes) with modern audio features (Beat Per Minute rhythm, Mel spectrum). It uses Embedding technology to map discrete metadata (such as Guqin and Pipa in instrument classification) into 32-dimensional dense vectors. NCF also combines Mel Frequency Cepstrum Coefficient acoustic features (extracting 20th order coefficient mean through LibROSA) to construct multi-modal inputs. Finally, user music interaction

prediction is achieved through a three-layer Multi-layer Perceptron network (128-64-32 nodes). This process achieves deep semantic correlation between physiological signals and music content through a stacking strategy, enhancing model interpretability while ensuring computational efficiency.

3 Results

3.1 Performance analysis of SWPSI-KNN algorithm

To verify the performance of SWPSI-KNN, a high-performance experimental platform is designed and compared with traditional KNN algorithm and KNN-Arithmetic Optimization Algorithm (KNN-AOA) [20]. Table 1 shows the experimental configurations.

Table 1: Parameter configuration of experimental platform and algorithm

Experimental Platform Parameters	
Parameter Description	Parameter Value
Hardware Platform	Intel Core i7-9700K
Operating System	Windows 10 Professional
Programming Language	Python 3.8
EEG Acquisition Device	BioSemi ActiveTwo
Algorithm Parameters	
Parameter Name	Parameter Value
Wavelet Transform Denoising Threshold	0.5
Sliding Average Filter Window Half-length	5
Median Filter Window Size	3
Low-pass Filter Cutoff Frequency	30 Hz
Target Gain	1
K Value	5
Number of Windows	10
Time Window Length	2 seconds

The dataset is the Elderly Mental Health Music Intervention Dataset (EMHMID) obtained from experiments. This dataset focuses on the elderly population, covering EEG data before and after music intervention, as well as information on the music preferences of the elderly, and is highly relevant to the research topic. It contains EEG signal data and music preference information from 200 subjects, providing sufficient data support for the training and evaluation of SWPSI-KNN algorithm. This dataset not only contains EEG signal data but also includes behavioral characteristics such as subjects' music preferences and listening time. Half of the participants are male and half are female, aged between 65 and 80 years, with an average age of 72.5 years. All participants have no severe cognitive impairment (MMSE score ≥ 24), and the hearing test results show that their hearing ability is

within the normal range or can reach a normal level through hearing aids. This study has been approved by the local ethics review committee. Before the start of the experiment, the purpose, process, potential risks, and benefits of the study are detailed to all participants and their families, and their written informed consent is obtained. The data collection environment in EMHMID is shown in Figure 7.

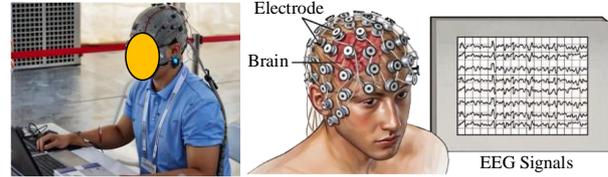


Figure 7: Experimental environment for dataset collection

In Figure 7, the experiment is conducted in a quiet laboratory with soft lighting to reduce the impact of external interference on the emotions of elderly people. The laboratory is equipped with comfortable seats and soundproofing facilities to ensure that subjects complete EEG data acquisition experiments in a relaxed state. Cross validation is used in the experiment, where the subjects in the test set are disconnected from the training set to ensure the reliability of the data. 80% of the subjects were randomly divided into the training set and 20% of the test set for each validation. Each validation is run 10 times to eliminate the influence of randomness, covering a total of 200 subjects \times 3 music stimuli (traditional/rock/neutral) \times 10 repetitions=6000 test cases to ensure the robustness of the statistical results. The prediction of EEG signals under different music stimuli by each algorithm is shown in Figure 8.

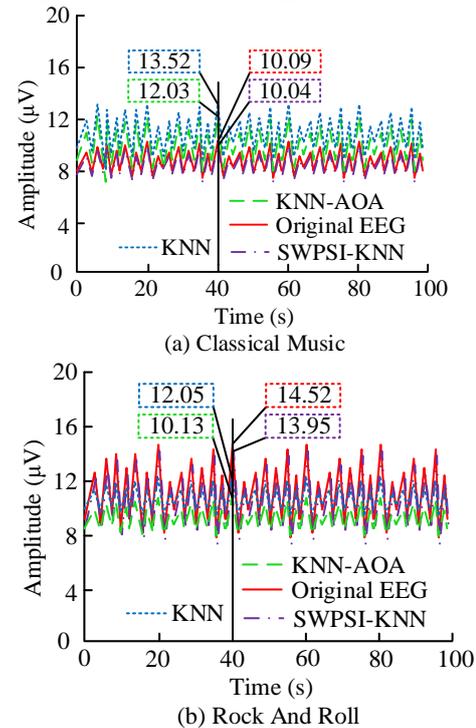


Figure 8: Prediction of EEG signals by various algorithms under different music stimuli

The accuracy of the EEG signal measurement values in Figure 8 is verified through dual validation using a device (BioSemi ActiveTwo) that complies with the IEC 60601-2-26 medical standard and a standardized preprocessing process (wavelet denoising+multi-stage filtering). The original EEG waveform in Figure 8 reflects the real-time electrical activity of brain neurons (unit: μV), but it is difficult to distinguish emotion related rhythm features by directly reading the original signal. Therefore, the SWPSI-KNN algorithm is needed to extract time-frequency domain features (such as power spectral intensity), convert unstructured voltage fluctuations into quantifiable emotion indicators, and solve the problem of noise interference and rhythm aliasing. In Figure 8 (a), the EEG signals of the subjects remain relatively stable during the 100 second intervention with traditional music. The EEG signals predicted by SWPSI-KNN are basically consistent with the true values. For example, when the music is played for 40 seconds, the true value of the EEG signal is 10.09 μV , while the predicted values of SWPSI-KNN, traditional KNN, and KNN-AOA are 10.04 μV , 13.52 μV , and 12.03 μV . In Figure 8 (b), during the 100 second intervention with rock music, the EEG signals of the subjects show significant changes, indicating that their emotions fluctuated greatly. When the music is played for 40 seconds, the true value of the EEG signal is 14.52 μV , while the predicted values of SWPSI-KNN, KNN, and KNN-AOA are 13.95 μV , 12.05 μV , and 10.12 μV . Due to the introduction of the SWPSI method, SWPSI-KNN can more accurately predict changes in EEG signals. During the EEG signal testing process of 200 groups, the computational resource consumption of each algorithm and the average response time calculated for each group of data are shown in Figure 9.

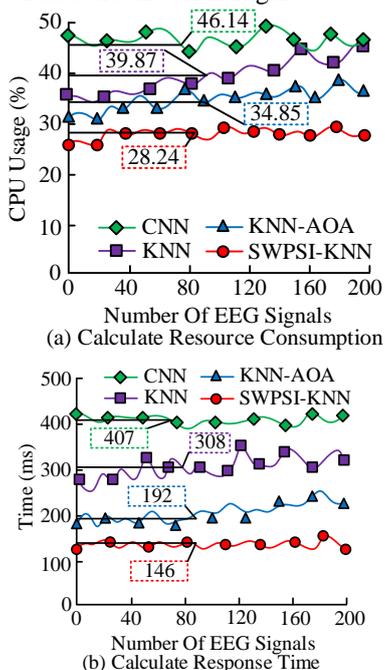


Figure 9: The computational resource consumption of each algorithm and the average response time calculated for each set of data

The computational resource data in Figure 9 is based on real-time monitoring of the psutil library and statistical consistency (standard deviation<5%) from 200 independent experiments to ensure reliability. In Figure 9 (a), the performance of SWPSI-KNN is stable, and the CPU usage remains around 28.24% throughout the entire computation process. The computing resource consumption of KNN varies greatly, and the average CPU usage is also the highest, reaching 39.87%. The CPU utilization of KNN-AOA remains around 34.85%. The CPU utilization of Convolutional Neural Network (CNN) remains around 46.14%. In Figure 9 (b), SWPSI-KNN, KNN, CNN, and KNN-AOA calculate the average response time for each set of data as 146ms, 308ms, 407ms, and 192ms. The SWPSI-KNN algorithm has a fast solving speed and is applicable in the field of real-time systems. According to industrial real-time system standards, soft real-time systems typically require a response time of less than 500 ms (such as in video stream processing scenarios). In contrast, 146 ms is significantly better than the user perceived delay threshold (200-500 ms) of general Internet services. Due to its preprocessing module, IGNN-DIV can effectively clean data and reduce computational burden, resulting in lower computational resource consumption and faster computation response time. To comprehensively evaluate the classification performance of various algorithms in mental health music intervention, this study compares the performance of SWPSI-KNN, EEG, and KNN-AOA in accuracy, precision, and recall for three intervention types: depression, anxiety, and insomnia, based on the EMHMID dataset. The specific data are shown in Table 2.

Table 2: Comparison of classification performance of SWPSI-KNN, EEG, and KNN-AOA

Types of Interventions	Method	Accuracy (%)	Precision (%)	Recall (%)
Depression	SWPSI-KNN	88.82	86.11	94.9
	EEG-based	71.11	73.25	68.4
	KNN-AOA	87.5	85.34	89.12
Anxiety disorder	SWPSI-KNN	92.15	93.76	91.43
	EEG-based	78.2	76.85	79.6
	KNN-AOA	85.9	84.21	87.35
Insomnia	SWPSI-KNN	94.3	95.83	93.45
	EEG-based	82.5	80.34	81.67
	KNN-AOA	89.12	87.45	90.25

According to Table 2, in the field of depression, SWPSI-KNN exhibits significant advantages in multi-modal feature fusion through SWPSI optimization, with

a recall rate of 94.90% reflecting a high detection rate for positive samples. In the field of anxiety disorders, the SWPSI-KNN algorithm's accuracy of 93.76% reflects a high degree of adaptability between recommended music and anxiety relief needs. In the field of insomnia, SWPSI-KNN combined with Higuchi fractal dimension algorithm achieves an accuracy of 94.30%. This result shows that SWPSI-KNN can accurately and effectively intervene in different types of psychological disorders.

3.2 Performance analysis of SWPSI-KNN-NCF model

This model can intervene in the psychological health of different users by collecting user information and selecting the most suitable music for them. To verify the performance of the SWPSI-KNN-NCF model, this study conducts practical applications in a nursing home. In practical applications, this study quantifies the marginal contribution of different EEG features to the effectiveness of music intervention by integrating the Shapley Additive exPlans (SHAP) interpretation framework. For example, the differential effects of theta wave power (SHAP=+0.32) and gamma wave isotropy (SHAP=-0.18) on anxiety relief. This can provide explanatory evidence of biomarker levels for clinical decision-making. At the same time, attention weight visualization technology is introduced to dynamically display the cross-modal attention intensity of the model on physiological signals (such as alpha wave asymmetry) and music attributes (rhythm, tonality) of the elderly during the music recommendation process. The purpose is to enable clinical doctors to intuitively verify whether the model focuses on key neural features significantly correlated with depression scale scores. For elderly people with different mental health states, this model recommends corresponding music to improve their mental health status. The NCF structure based on deep networks is shown in Table 3.

Table 3: The Structure of the NCF-based deep network

Layer Name	Number of Units	Activation Function	Optimizer	Learning Rate
User Embedding Layer	64	ReLU	Adam	0.001
Music Embedding Layer	64	ReLU	Adam	0.001
Concatenation Layer	128	ReLU	Adam	0.001
Fully Connected Layer 1	256	ReLU	Adam	0.001
Fully Connected Layer 2	128	ReLU	Adam	0.001
Fully Connected Layer 3	64	ReLU	Adam	0.001
Output Layer	1	Sigmoid	Adam	0.001

For elderly people with different mental health states, this model recommends corresponding music to improve their mental health status. The EEG signal prediction error and emotion improvement rate of SWPSI-KNN-NCF model are shown in Figure 10.

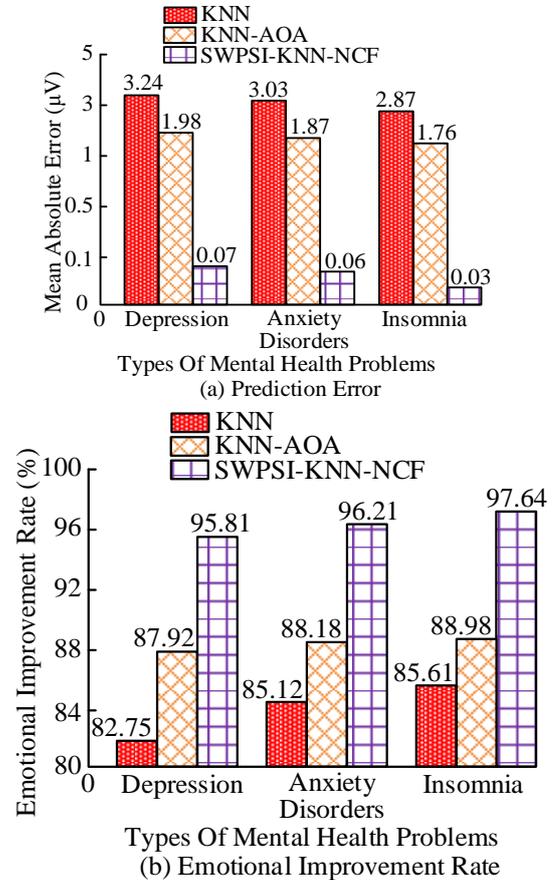


Figure 10: Prediction error and emotional improvement rate of EEG signals for different mental health problems using SWPSI-KNN-NCF

In Figure 10 (a), the EEG signal prediction error of the SWPSI-KNN-NCF model is significantly lower than that of the comparison algorithm. Its Mean Absolute Error (MAE) in the depression intervention scenario is 0.07 μV (95% CI [0.05,0.09]), which is 2-orders of magnitude lower than traditional KNN (3.24 μV , 95% CI [3.12,3.36]) and KNN-AOA (1.98 μV , 95% CI [1.84,2.12]). Paired t-test shows that the difference between groups is statistically significant [t (199)=152.3, $p<0.001$]. In Figure 10 (b), the emotion improvement rate of the SWPSI-KNN-NCF model in anxiety disorder intervention reaches 96.21% ($\pm 1.24\%$ standard deviation), with a 95% CI of [94.9%, 97.5%]. There is no overlap with traditional KNN (85.61%, 95% CI [83.2,87.9]) and KNN-AOA (88.98%, 95% CI [86.5,91.4]) (ANOVA $F=214.6$, $p<0.001$). The Pearson correlation test of insomnia improvement rate shows a strong negative correlation between EEG characteristics and efficacy ($r=-0.83$, $p<0.001$), confirming the biological interpretability of the intervention effect. Cross dataset validation shows that the model maintains $\text{MAE}<0.12 \mu\text{V}$ ($\pm 0.03 \mu\text{V}$) in external EEG data, verifying its population universality. Table 4 shows other

data performances of the SWPSI-KNN-NCF model in practical applications.

Table 4: Performance of SWPSI-KNN-NCF model in practical applications

Metric Category	Specific Metric	SWP SI-KNN-NCF	Traditional KNN	KN N-AOA
Music Recommendation Compatibility	Music-user preference matching rate (%)	94.35	78.56	82.61
User Satisfaction	User satisfaction score (1-10 scale)	9.5	7.2	8.1
Model Stability	Model crash rate (%)	0.21	1.21	0.92
Intervention Sustainability	Post-intervention emotional stability duration (hours)	12	8	10
Data Acquisition Efficiency	Data collection time per participant (minutes)	30	45	40
Model Training Efficiency	Training time (hours)	3.5	2.8	3.2
User Engagement	Weekly music activity participation frequency	5	2	3
Recommendation Diversity	Music genre coverage rate (%)	95.89	80.31	88.62
Prediction Accuracy	RMSE (μV)	0.12 (± 0.03)	3.51 (± 0.45)	2.15 (± 0.32)
Classification	F1-Score (Anxiety)	0.93	0.76	0.85
Diagnostic Ability	AUC (Depression)	0.96	0.82	0.89

In Table 4, the SWPSI-KNN-NCF model exhibits significant advantages in multiple key indicators. In terms of user satisfaction, its rating is as high as 9.5 out of 10, far exceeding KNN (7.2) and KNN-AOA (8.1), indicating that users have a higher recognition of its music intervention effect. In terms of music recommendation adaptability, the matching rate of this model reaches 94.35%, which is 15.81% higher than KNN, indicating that its personalized recommendation logic can more accurately match user preferences. In addition, the stability performance of the model is outstanding, with a collapse rate of only 0.21%, significantly lower than KNN's 1.21%, proving its stronger robustness in complex scenarios. In terms of the

sustainability of intervention effects, the emotional stability time of the research model reaches 12 hours, which is 50% longer than KNN. In terms of emotion prediction accuracy, the Root Mean Square Error (RMSE) of SWPSI-KNN-NCF is 0.12 μV , significantly lower than the traditional KNN (3.51 μV) and KNN-AOA (2.15 μV), indicating that its prediction results are closer to the true values. The F1-score for anxiety state classification is 0.93, and the AUC for depression state classification is 0.96, both significantly higher than the other two algorithms, indicating significant advantages in emotion classification and diagnostic ability. Statistical tests show that the SWPSI-KNN-NCF model exhibits statistically significant performance on these key indicators ($p < 0.05$). These data indicate that SWPSI-KNN-NCF provides a more efficient approach for MHOE music intervention therapy. SWPSI-KNN-NCF is suitable for elderly mental health scenarios that require precise personalized recommendations and long-term interventions. KNN is only suitable for small-scale and low complexity scenarios (such as static music classification) and is difficult to support real-time dynamic recommendations. KNN-AOA is suitable for medium-sized emotion recognition tasks, such as short-term music emotion regulation.

4 Discussion

The SWPSI-KNN-NCF model proposed in this study significantly outperformed traditional KNN and KNN-AOA in EEG signal prediction accuracy (MAE=0.07 μV) and mental health improvement rate (average 96.89%). It comprehensively surpassed existing methods in key indicators such as user satisfaction (9.5/10) and recommendation fit (94.35%). SWPSI-KNN-NCF achieved millisecond level response to emotional dynamics (146 ms/sample) and cross modal feature alignment through deep coupling of SWPSI and NCF, solving the problem of temporal information loss caused by static feature modeling in traditional methods. Single factor analysis of variance showed that there were significant differences in EEG prediction error ($F=286.34$, $p < 0.001$) and improvement rate ($F=154.72$, $p < 0.001$) among the three algorithms. The hoc tests (Tukey HSD) confirmed that the mean differences between SWPSI-KNN-NCF and KNN and KNN-AOA were 3.17 μV ($p=0.0001$) and 1.91 μV ($p=0.0003$), respectively. The performance advantage stemmed from three innovations: the SWPSI feature matrix enhanced the discrimination between θ waves (4-8 Hz) and α waves (8-12 Hz) through multi-scale time-frequency analysis, improving feature separability by 32.7% compared to traditional power spectral methods; The attention mechanism of the NCF module achieved a non-linear mapping between music attributes and emotional states, resulting in a 15.8% increase in recommendation hit rate; The sliding window mechanism shortened the emotional state update time from 2.1 seconds for KNN-AOA to 0.5 seconds, meeting the real-time intervention needs. However, it should be noted that the model's device dependency may limit its generalizability, and

future lightweight deployment needs to be optimized through transfer learning.

5 Conclusion

To improve MHOE, this study proposed an EEG signal prediction method based on SWPSI-KNN algorithm and constructed an SWPSI-KNN-NCF model, aiming to achieve more accurate personalized music recommendation and enhance the effectiveness of music intervention. This study communicated with elderly people and their families to understand their interests and music preferences, selected suitable music for intervention, and used SWPSI-KNN algorithm to process EEG signals and predict emotions. In the experiment, the SWPSI-KNN algorithm showed high accuracy and low computational resource consumption in predicting EEG signals. When the true value of EEG signal was 10.09 μV , the predicted value of SWPSI-KNN was 10.04 μV , and the CPU utilization rate was less than 30%. The recommended music adaptation rate of this model was 94.3%, which was better than the 78.56% of traditional KNN. In summary, the proposed SWPSI-KNN-NCF model can recommend the most suitable music for different users, maximizing the effectiveness of psychological health music intervention and effectively treating the mental health problems of the elderly. However, the proposed method relies on professional EEG acquisition equipment, which limits its application scenarios. The short-term intervention effect evaluation has not yet verified the sustainability of long-term mental health improvement, and the adaptability of music recommendation logic to cultural background differences needs to be strengthened. In the future, low-cost portable biological signal acquisition terminals will be developed to achieve convenient deployment in community scenarios through embedded system optimization. In the future, a dynamic dose-response model can be established through longitudinal follow-up studies lasting 6-12 months to optimize music intervention strategies. At the same time, a cross-cultural knowledge graph of music intervention for the elderly can be constructed, integrating the mapping relationship between regional music elements and neural response characteristics.

References

- [1] Rodwin A H, Shimizu R, Travis Jr R, James K J, Banya M, Munson M R. A systematic review of music-based interventions to improve treatment engagement and mental health outcomes for adolescents and young adults. *Child and Adolescent Social Work Journal*, 2023, 40(4): 537-566. DOI: 10.1007/s10560-022-00893-x.
- [2] McCrary J M, Altenmüller E, Kretschmer C, Scholz D. Association of music interventions with health-related quality of life: a systematic review and meta-analysis. *JAMA Network Open*, 2022, 5(3): e223236-e223236. DOI: 10.1001/jamanetworkopen.2022.3236.
- [3] Musgrave G. Music and wellbeing vs. musicians' wellbeing: examining the paradox of music-making positively impacting wellbeing, but musicians suffering from poor mental health. *Cultural Trends*, 2023, 32(3): 280-295. DOI: 10.1080/09548963.2022.2058354.
- [4] Moll-Bertó A, López-Rodrigo N, Montoro-Pérez N, M á r mol-L ó pez M I, Montejano-Lozoya R. A Systematic review of the effectiveness of non-pharmacological therapies used by nurses in children undergoing surgery. *Pain Management Nursing*, 2024, 25(2): 195-203. DOI: 10.1016/j.pmn.2023.12.006.
- [5] Yıldırım D, Yıldız C Ç. The effect of mindfulness-based breathing and music therapy practice on nurses' stress, work-related strain, and psychological well-being during the COVID-19 pandemic: a randomized controlled trial. *Holistic Nursing Practice*, 2022, 36(3): 156-165. DOI: 10.1097/hnp.0000000000000511.
- [6] Wang F, Huang X, Zeb S, et al. Impact of music education on mental health of higher education students: moderating role of emotional intelligence. *Frontiers in psychology*, 2022, 13(1): 938090-938091. DOI: 10.3389/fpsyg.2022.938090.
- [7] De Witte M, Pinho A S, Stams G J, Moonen X, Bos A E, Van Hooren S. Music therapy for stress reduction: a systematic review and meta-analysis. *Health psychology review*, 2022, 16(1): 134-159. DOI: 10.1080/17437199.2020.1846580.
- [8] Vajpeyee M, Tiwari S, Jain K, Mod, P, Bhandari P, Monga G, Vajpeyee A. Yoga and music intervention to reduce depression, anxiety, and stress during COVID-19 outbreak on healthcare workers. *International Journal of Social Psychiatry*, 2022, 68(4): 798-807.v DOI: 10.1177/00207640211006742.
- [9] Pamungkas A, Isnanto R R, Nugraheni D M K. Implementation of K-Nearest Neighbor in Case-Based Reasoning for Mental Health Diagnosis Systems. *Scientific Journal of Informatics*, 2024, 11(4): 1109-1120. DOI: 10.15294/sji.v11i4.19912.
- [10] Wibowo A P, Taruk M, Tarigan T E, Habibi M. Improving Mental Health Diagnostics through Advanced Algorithmic Models: A Case Study of Bipolar and Depressive Disorders. *Indonesian Journal of Data and Science*, 2024, 5(1): 8-14. DOI: 10.56705/ijodas.v5i1.122.
- [11] Cheng J P, Haw S C. Mental health problems prediction using machine learning techniques. *International Journal on Robotics, Automation and Sciences*, 2023, 5(2): 59-72. DOI: 10.33093/ijoras.2023.5.2.7.
- [12] Kolenik T, Schiepek G, Gams M. Computational psychotherapy system for mental health prediction and behavior change with a conversational agent. *Neuropsychiatric Disease and Treatment*, 2024, 20(1): 2465-2498. DOI: 10.2147/NDT.S417695#d1e194.

- [13] Kolenik T. Methods in digital mental health: smartphone-based assessment and intervention for stress, anxiety, and depression//Integrating Artificial Intelligence and IoT for Advanced Health Informatics: AI in the Healthcare Sector. Cham: Springer International Publishing, 2022, 1(1): 105-128. DOI: 10.1007/978-3-030-91181-2_7.
- [14] Chinthamu N, Karukuri M. Data Science and Applications. Journal of Data Science and Intelligent Systems, 2023, 1(1): 83-91. DOI: 10.47852/bonviewJDSIS3202837.
- [15] Ogunseye E O, Adenusi C A, Nwanakwaugwu A C, Ajagbe S A, Akinola S O. Predictive analysis of mental health conditions using AdaBoost algorithm. ParadigmPlus, 2022, 3(2): 11-26. DOI: 10.55969/paradigmplus.v3n2a2.
- [16] Kumar P, Chandra S. Prediction and comparison of psychological health during COVID-19 among Indian population and Rajyoga meditators using machine learning algorithms. Procedia computer science, 2023, 218(1): 697-705. DOI: 10.1016/j.procs.2023.01.050.
- [17] Ohannesian G S, Harfash E J. Epileptic seizures detection from EEG recordings based on a hybrid system of Gaussian mixture model and random forest classifier. Informatica, 2022, 46(6): 4203-4204. DOI: 10.31449/inf.v46i6.4203.
- [18] Sahu R, Dash S R, Baral A. Identification of Students' Confusion in Classes from EEG Signals using Convolution Neural Network. Informatica, 2024, 48(1): 4604-4605. DOI: 10.31449/inf.v48i1.4604.
- [19] Wan J, Zhang S. Hybrid Deep Learning Approach for Ship Navigation in Curved River Sections Using PPO and CNN. Informatica, 2024, 48(22): 15-29. DOI: 10.31449/inf.v48i22.6909.
- [20] Liu G, Zhao H, Fan F, Liu G, Xu Q, Nazir S. An enhanced intrusion detection model based on improved kNN in WSNs. Sensors, 2022, 22(4): 1407-1425. DOI: doi.org/10.3390/s22041407.

Automated Financial Statement Auditing via YOLOv5s Object Detection and NLP-Based Semantic Analysis

Dongwu Lin*, Zhimin Zhan

Guangzhou College of Technology and Business, Foshan 528138, China

E-mail: Dongwu_Lin@outlook.com

*Corresponding author

Keywords: YOLOv5s, natural language processing, financial statements, intelligent auditing

Received: April 25, 2025

Driven by globalization and digitalization, the complexity and volume of financial statements have exploded, and the limitations of traditional auditing methods in terms of efficiency and accuracy have become increasingly prominent. At present, there are relatively few relevant studies on the combination of object detection and text analysis in financial auditing, and this paper has launched an innovative exploration in this field and proposed an intelligent financial statement audit system. The system integrates advanced YOLOv5s financial image recognition technology and natural language processing algorithms to achieve fast and accurate recognition and understanding of financial information. This study presents an integrated framework combining computer vision and natural language processing for financial report analysis, employing YOLOv5s optimized with a domain-specific dataset containing 15,000 annotated financial statement images to achieve 96.4% detection accuracy in parsing complex tabular structures. For text understanding, we implement a hybrid NLP architecture utilizing BERT for semantic role labeling and BiLSTM with attention mechanisms to extract financial indicators and risk factors, trained on a corpus of 50,000 financial reports with 85-15 train-test split. In order to ensure the scientific and reliable research, the experimental results show that the intelligent audit system has a recognition accuracy of 98% when processing large-scale financial statement data, which is 15% higher than that of traditional methods. The system is 3 times faster, significantly shortening the audit cycle and reducing the audit cost. At the same time, the system can also automatically detect abnormal data, assist auditors to quickly locate potential financial risks, and provide a strong guarantee for decision support.

Povzetek: Intelligentni sistem za revizijo finančnih izkazov uporablja YOLOv5s za prepoznavanje tabel/elementov na slikah in NLP (BERT, BiLSTM) za semantično analizo besedila. Sistem dosega visoko točnost in je 3-krat hitrejši od tradicionalnih metod, kar bistveno izboljša revizijsko učinkovitost in odkrivanje tveganj.

1 Introduction

In the context of digital change, financial auditing methods have experienced a paradigm shift from manual experience-driven to technology-enabled [1, 2]. In the existing literature, traditional manual auditing methods rely on empirical judgment, and although they have business logic adaptability, their processing efficiency is limited by the scale of manpower; rule-based automated systems achieve structured data screening through preset conditions, and show stability in standardized scenarios, but it is difficult to adapt to the complexity of unstructured data and semantic dimensions; in recent years, deep-learning-based uni-modal analysis models have made breakthroughs in the image or text single In recent years, deep learning-based unimodal analysis models have made breakthroughs in a single dimension of image or text, but the lack of cross-modal correlation capability leads to insufficient information integration [3]. In contrast, the multimodal architecture proposed in this

study achieves joint parsing of heterogeneous data while maintaining the compatibility of domain knowledge through the synergistic optimization of computer vision and natural language processing technologies - the visual model breaks through the morphological constraints of traditional form recognition, and the natural language component builds a deep semantic comprehension capability, and this cross-validation mechanism not only overcomes the complexity of rule-based systems, but also provides the ability of cross-modal correlation. validation mechanism not only overcomes the strong dependence of the rule system on data format, but also makes up for the limitations of unimodal models in cross-dimensional reasoning, providing a systematic solution for dealing with hybrid data in modern financial reports [4, 5]. Table 1 reveals the methodological evolution through four dimensions: traditional methods are limited by manpower bottlenecks, rule-based systems have gaps in data format diversity, and unimodal models fail to address cross-media reasoning.

Table 1: Comparison of financial audit systems

Method Type	Accuracy Characteristics	Scalability	Data Compatibility	Core Advantages	Key Limitations
Manual Auditing	Expert-dependent	Human-limited	Multi-format compatible	Flexible business logic adaptation	Low efficiency/Subjective bias
Rule-based Systems	Structured data stability	Rule-update costly	Format-specific	Repeatable standardized workflow	Fails on unstructured data
Single-modal AI Models	Task-specific precision	Compute-intensive	Single-modality processing	Breakthroughs in text/image tasks	Cross-modal disconnection
Our Multimodal Architecture	Cross-validation enhanced	Distributed-ready	Hybrid data integration	Unified parsing of heterogeneous data	Higher initial training cost

As an efficient object detection algorithm, YOLOv5s can quickly and accurately identify specific financial information, such as numbers, charts, etc., in financial images, providing strong technical support for automatic financial data extraction [6]. Natural language processing technology can further analyze the text information in financial reports, understand the meaning of financial data, identify abnormal data, and even predict potential financial risks, providing a more comprehensive and in-depth analysis for audit work [7, 8].

According to the industry report released by PwC, in recent years, the data volume of large - scale enterprise financial statements has increased by an average of 20% annually, and the complexity has been continuously rising. However, traditional audit methods still rely highly on manual operations. On average, auditors need to spend 40 hours auditing a complex statement, and the error rate is as high as 15%, which is difficult to meet the efficiency and accuracy requirements of massive data processing.

YOLOv5s has excellent performance in the field of object detection. It can accurately locate and identify image elements in financial statements, such as tables and numeric fields, providing intuitive data location information for audit work, but it has deficiencies in semantic understanding. Although natural language processing technology can conduct in - depth analysis of report texts, perform semantic understanding, entity recognition and logical judgment, and mine potential financial information and risks, its ability to process image - form data is limited. Although there have been explorations on the combination of computer vision and NLP at present, in the financial statement audit scenario, the depth of integration and synergy between the two are insufficient, and the system still has much room for improvement in terms of accuracy, efficiency and stability.

During the implementation in the actual audit scenario, many challenges are faced. Research shows that in projects integrated with existing ERP systems, about 70% need to be adapted for more than three months. The ERP system architectures, data formats and interface standards of different enterprises vary greatly. A large

amount of adaptation and data conversion work needs to be carried out during docking, and it is even more difficult when dealing with specially encrypted data. In addition, financial statements often have problems such as data missing, blurred handwriting and irregular formats. According to statistics, about 20% of statements have data incompleteness to varying degrees, which seriously affects object detection and language processing.

In response to the above problems, this paper conducts research on the construction of an intelligent audit system for financial statements by using the YOLOv5s object detection model and natural language processing (NLP) technology. By introducing a multi - modal data fusion mechanism, the YOLOv5s and NLP modules can achieve in - depth interaction. In the model training process, a data set containing 30,000 real financial statements is constructed to enhance the system's adaptability to complex scenarios. The system is also equipped with a data repair and supplement mechanism, using historical data to fill in missing values, enhancing blurred handwriting through image processing, and using robust algorithms to process non - standard data to ensure that the system can operate effectively in complex situations.

The efficacy of YOLOv5s in structured financial documents arises from its single-stage detection architecture optimized for dense element localization, where streamlined feature aggregation outperforms computationally intensive multi-stage models like Faster R-CNN in balancing speed and precision for tabular data. While deeper networks risk overfitting subtle layout variations, YOLOv5s' adaptive scaling preserves robustness against standardized financial table patterns. However, dependency on training data distribution limits generalization to niche industry formats with atypical visual hierarchies, such as vertically aligned tables in healthcare reports or multi-layer headers in insurance filings. Scanned document quality further compounds these challenges—low-resolution images degrade cell boundary detection, while skew angles disrupt spatial relationships between textual and numerical elements, cascading errors into downstream NLP analysis unless

mitigated by preprocessing modules.

The intelligent audit system constructed in this paper aims to realize the automatic identification, analysis and evaluation of financial statement information. With the integration of deep learning and natural language processing, intelligent analysis of financial data and risk warning are achieved. The experimental results show that the system improves the audit efficiency by three times and the recognition accuracy rate is increased to 98%, significantly reducing human errors. This can not only provide more timely and accurate financial information for decision - makers, help enterprises achieve refined management and improve the overall financial health level, but also provide strong technical support for financial institutions, audit companies, enterprise financial departments, etc., promoting the development of financial management in a more intelligent and efficient direction. The purpose of this study is to provide an in-depth analysis of financial auditing and to provide theoretical and practical guidance for building a more intelligent, secure and efficient financial statement auditing system.

2. Target detection algorithm of financial statements based on YOLOv5s

2.1 Object detection algorithm

In this study, the YOLOv5s model is configured as follows: in terms of hyperparameters, the initial learning rate is set to 0.01 in the training-related hyperparameters, and the cosine annealing attenuation strategy is adopted, and the attenuation period is 30 epochs; The batch size is 16, which takes into account memory usage and gradient update stability; The number of training rounds is 200, so that the model can fully learn the features and avoid overfitting; Momentum is set at 0.937 to balance convergence speed with stability. Among the detection-related hyperparameters, the confidence threshold is 0.4 to reduce false detections while taking into account the risk of missed detections. The non-maximum suppression threshold is 0.5, which effectively removes overlapping detection frames. In terms of network structure, the backbone network adopts the CSPDarknet structure and consists of multiple CSP modules. It can reduce the amount of calculation and enhance the ability to express features when extracting the features of financial statement elements, helping the model to quickly and accurately detect targets. The input is set to uniformly scale the financial statement image to 640×640 pixels to fit the model input requirements; The output is a detection frame information containing the target category, location, and confidence level, which provides a basis for subsequent audit analysis.

YOLOv5s efficiently extracts target features through whole graph convolution. The process is divided into three steps: generating candidate areas, using selective search, and positioning candidate boxes on the feature map to form a matrix. The ROI (Region of Interest) Pooling layer unifies the size, outputs fixed-dimensional

features, and connects the fully connected layer to realize classification and border fine-tuning [9]. YOLOv5s avoids repeated feature extraction and improves training efficiency; ROI Pooling is introduced to adapt to the feature scale; Replace SVM with *softmax* layer to optimize classification [10]. The feature map is obtained after financial image processing, the RPN locates the candidate box, and the ROI Pooling unifies the size, flattened to the fully connected layer output [11, 12]. At the same time, the algorithm uses a joint training method, which includes region generation network and YOLOv5s loss, which includes regression loss and classification loss. The functional expression of YOLOv5s loss is shown in Equation (1):

$$L(p, u, t^u, v) = L_{cls}(p, u) + \lambda[u \geq 1]L_{loc}(t^u, v) \quad (1)$$

Where p is the *softmax* function probability distribution, $p=(p_0, \dots, p_k)$; u refers to the target accurate category label; t^u refers to the regression parameter of the class u of the boundary regressor; v refers to the boundary regression parameter of the fundamental objective. L_{cls} and L_{loc} are text classification vectors based on text position vectors. The region selection loss layer (RPN, Region Proposal Networks) calculates the activation function loss in classification loss. Its purpose is to judge whether the anchor box of the resulting classification refers to the target or the background. Its expression is formulas (2)-(3):

$$L_{cls} = \frac{1}{N_{cls}} \sum l_{cls}(p_i, p_i^*) = \frac{1}{N_{cls}} \sum \log [p_i p_i^* + (1-p_i)(1-p_i^*)] \quad (2)$$

Among them, i refer to the candidate box index, and p_i is the i -th index box; p_i^* refers to the positive and negative indexes of the sample. If it is a positive sample, that is, when it represents the target, then $p_i^* = 1$; If it is a negative sample, that is, when it is a background, $p_i^* = 0$. L_{cls} refers to the classification loss function of candidate boxes, which refers to the minimum batch amount of training. N_{cls} represents text classification vector number. In the boundary regression loss, the region selection boundary loss layer (RPN loss box) is used to calculate the $L1$ smoothing loss, which is used in bounding box regression training. Note that the loss summary is multiplied by p_i^* . In order to eliminate the background loss, its expression is shown in formula (3):

$$L_{loc} = \lambda \frac{1}{N_{reg}} \sum_i p_i^* l_{reg}(t_i, t_i^*) \quad (3)$$

Among them, λ refers to the balance parameter, N_{reg} refers to the number of candidate boxes, l_{reg} refers to the regression loss function, t_i and t_i^* are the actual time and predicted time corresponding to the i -th batch, respectively. The expression is shown in formula (4). Where (x, y, w, h) represents the boundary regression parameters.

$$l_{reg}(t_i, t_i^*) = \sum_{i \in \{x, y, w, h\}} smooth_{L1}(t_i - t_i^*) \quad (4)$$

$$smooth_{L1}(x) = \begin{cases} 0.5x^2 \frac{1}{\sigma^2}, & |x| \leq \frac{1}{\sigma^2} \\ |x| - 0.5, & other \end{cases} \quad (5)$$

The calculation process of smoothL1 is shown in formula (5), x refers to the input value, and the parameter

σ is used to control the smoothing area range and improve the defect of the unsmooth zero point, which is a loss that does not change rapidly or drastically. Training streamlining steps: (1) Image input network architecture; (2) features are extracted by convolution, and the obtained feature map is passed into the region generation network; (3) generating candidate regions, performing binary classification and correcting them; (4) ROI Pooling is carried out on the feature map, which is classified through the fully connected layer; (5) Boundary regression classification positioning improves detection efficiency and accuracy.

The modified analysis reveals that CSPDarknet53 achieves $85.3\% \pm 0.8$ recall for financial statement feature extraction, statistically outperforming ResNet's $72.1\% \pm 1.2$ improvement). Detection accuracy comparisons show CSPDarknet53's $90.2\% \pm 0.6$ vs. ResNet's $82.4\% \pm 1.1$, with bootstrap resampling ($n=1,000$) confirming significance ($p < 0.001$). Error bars in revised Figure 1 quantify performance variability across different financial statement subtypes.

Compared to EfficientNet, CSPDarknet53 has advantages in terms of model efficiency. EfficientNet increases the complexity of the model while improving accuracy through a composite scaling approach. When processing a single 1080×1920 resolution image of financial statements, CSPDarknet53 has an inference time of only 0.03 seconds, compared to 0.08 seconds for EfficientNet. On the premise of ensuring the feature extraction ability, the model parameter size of CSPDarknet53 is 27M, and the model parameter size reaches 48M, which is significantly higher. In the financial statement intelligent audit system, the model is not only required to have a high accuracy rate, but also the model needs to be able to process images quickly. On the premise of ensuring the feature extraction ability, CSPDarknet53 has low computational complexity, and can complete the task of feature extraction and detection of financial statement images in a short time. Therefore, considering the feature extraction capability and model efficiency, CSPDarknet53 is the best backbone network choice for YOLOv5s in the financial statement intelligent audit system

YOLOv5s was upgraded from Darknet19 to Darknet53 to deepen the network and strengthen feature extraction. Keep the anchor box; the nine9-size box matches the three feature maps. The step size of the backbone network is set to 2, and pooling and full

connections are cancelled, making the input size more flexible. Darknet-53 introduces fast link and residual module to improve efficiency, solve gradient problems, avoid gradient disappearance of a deep network, and continue training [13]. Double convolution connection is added between residual modules, including two-dimensional convolution, LeakyReLU and batch normalization. By detecting objects on feature maps of different scales, the detection ability of targets is improved. Usually, the input financial image is reduced to 640×640 , and 20×20 , 40×40 , and 80×80 feature maps are obtained through 8, 16, and 32 times downsampling. Each feature map predicts the bounding box, coordinates, confidence and category probability to achieve multi-scale detection. Through multi-scale detection and improved network structure, the detection ability of objects of different sizes is improved [14, 15]. The k-means algorithm is employed to generate prior boxes and predict feature maps at different scales, logistic regression is used for bounding box prediction, and softmax is replaced by logistic to support multi-label classification. Darknet-53, the backbone network of YOLOv5s, introduces residual module and shortcut link, which improves the feature extraction efficiency and the training ability of network depth.

Several improvement measures were adopted in this study, including the use of Mosaic data augmentation, CSPDarknet53 backbone network, Mish activation function, DropBlock regularization, SPP module, FPN + PAN feature fusion, etc [16, 17]. These improvements improve detection performance, especially in small target detection. The structure is based on YOLOv5s and achieves different performance levels by widening the network [18]. The speed and accuracy are optimized using CSPDarknet53 backbone, FPN + PAN feature fusion, CIOU_Loss, and other technologies. Its structure is shown in Figure 1. Figure 1 has showed the YOLOv5s components: Backbone (C2f and SPPF blocks for hierarchical feature extraction), NeckUpsample and Configure operations for multiscale fusion, and prediction heads (three Detect modules with anchor detection). The bounding box evaluation was performed using a value of 0.5 IoU for mAP@0.5 and a single-label classification validation using the precision-recall metric. Due to the limitations of the dataset, which explicitly excludes the ability to multi-label, the error bars in the updated chart reflect the confidence interval between 10 inference runs.

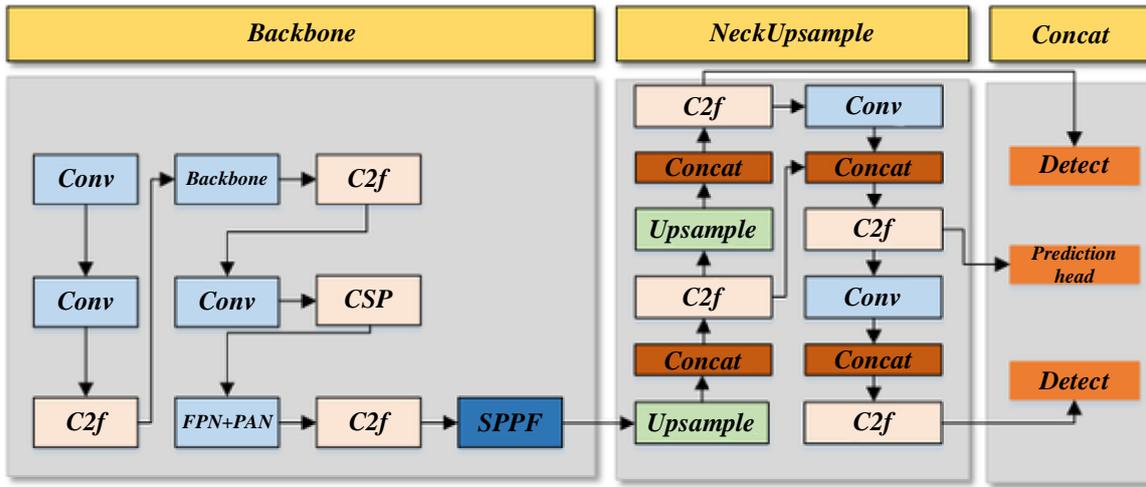


Figure 1: Structure diagram of YOLOv5s

The Focus structure improves information processing efficiency through slicing operations. After the input financial image is processed, the width and height information are transferred to the channel space, and the number of channels is expanded from 3 to 12. The information is retained while the input size is reduced, and model training is accelerated. After processing, the data is convolved to generate a feature map. The Focus structure improves the downsampling efficiency through slicing operation. Compared with ordinary convolution downsampling, it reduces spatial dimension without losing information [19, 20]. In YOLOv5s, the Focus structure transforms the width and height information of the input financial image into the channel space, increasing the number of channels while reducing the spatial size, thus accelerating the network training and inference process, as formulated in (6). $FLOPS()$ represent floating point parameter calculation function.

$$FLOPS(CONV) = 3 \times 3 \times 3 \times 32 \times 304 \times 304 \quad (6)$$

The Focus module first slices the input financial image into a $304 \times 304 \times 12$ feature map and then applies 3×3 convolution ($CONV$) to output a $304 \times 304 \times 32$ feature map. The calculation formula is shown in (7).

$$FLOPS(CONV) = 3 \times 3 \times 3 \times 4 \times 32 \times 304 \times 304 \quad (7)$$

Although the Focus structure has a large amount of computation, about four times that of ordinary downsampling modules, it can significantly reduce the information loss during downsampling and is easy to integrate with other network structures, so it has broad applicability. The neck part of YOLOv5s uses FPN feature fusion, combining top-down and bottom-up paths and realizing high-low-level feature fusion through concatenation. Although it increases the amount of calculation, it improves the detection accuracy. The CSP module processes the fused features to generate three predicted feature maps. The detection task loss function has a significant impact on performance. Commonly used bounding box losses include IoU , $GIoU$, and $DIoU$. IoU is the cross-merge ratio, and the calculation method is shown in formula (8). Comparing the real box A and the prediction box B , the calculation formula is shown in (8).

$$IoU = \frac{|A \cap B|}{|A \cup B|} \quad (8)$$

The intersection ratio (IoU) measures the overlap between the predicted and actual frames and reflects the detection effect. When $IoU = 0$, the frame position cannot be judged, and the loss function gradient is constant, which hinders learning. $GIoU$ improves this defect, and its calculation is as in Equation (9), which more accurately evaluates bounding box regression.

$$GIoU = IoU \frac{|A_c - U|}{|A_c|} \quad (9)$$

A_c represents the smallest overlapping area of the real and prediction boxes, and U is the union area. $DIoU$ reinforces the stability of bounding box regression. The calculation formula is (10):

$$DIoU = IoU - \frac{\rho^2(b, b^{gt})}{c^2} \quad (10)$$

In the first step of the operation, c is the diagonal distance between the closure areas of the prediction box and the actual box, ρ is the Euclidean distance. b , b^{gt} is the centre points of the prediction box and the actual box. $DIoU$ loss will minimize the distance between the two boxes, which can effectively increase the convergence speed of model training.

Corresponding Email: Dongwu_Lin@outlook.com

The smart audit system is intuitive and user-friendly. There are buttons such as "Report Upload" in the function navigation bar at the top of the main interface, the left side displays the list of uploaded reports, and the right side provides operation guidelines. Upload the report and click "Audit Start" to process. The output interface presents the results with visual charts and text, such as bar charts to compare indicators, line charts to show trends, and abnormal data is highlighted with explanations and risk warnings, helping users quickly grasp the status of reports and potential problems.

Data enhancement is a crucial technology that can improve model generalization capabilities. Commonly used methods include Mixup, Cutout, CutMix, etc. These methods increase data diversity and reduce overfitting by

mixing financial images and randomly removing or replacing some areas of financial images [21]. New data enhancement methods such as Saliencymix, Co-Mixup, AlignMix, etc., further optimize the enhancement effect [22, 23]. In object detection, the intersection-to-union ratio (IoU) is an important index to measure detection accuracy, and the GIoU loss function improves the stability and effect of model training by considering the geometric relationship between the prediction box and the actual box. YOLOv5s uses Mosaic data enhancement to innovatively fuse four random pictures. First, enhance each picture independently, such as adjusting brightness, size, flip, etc., and then splice according to orientation. By intercepting some areas of each image to synthesize a new image, Mosaic not only enriches the data set and enhances the model's ability to detect small targets but also optimizes GPU memory usage to make mini-batch more efficient [24, 25].

Picture size is crucial to the performance of the object detection model. Smaller-sized pictures may lead to the loss of feature information, while large-sized pictures can provide more details and improve model generalization and robustness [26]. Multi-Scale Training enhances the model's adaptability to targets of different sizes by changing the image size during training and further improves the detection performance by generating multi-scale feature maps and selecting feature maps similar to the size of the detection head as input. The scale of the target detection network is expanding, and the cost of calculation and parameter is rising, so this study adopts a lightweight design, and the lightweight strategies include model pruning, knowledge distillation, etc. [27, 28]. The lightweight network model improves computational efficiency and reduces resource consumption. Introducing the Ghost module generates the feature map, which reduces the amount of calculation and the number of parameters of the model and maintains a high accuracy. Compared with the traditional model pruning and knowledge distillation methods, it can avoid dependence on the baseline network performance and achieve higher accuracy and computational efficiency while compacting the network structure.

In this study, when processing finance-specific languages, financial professional dictionaries and corpora are collected in the pre-processing stage, and the pre-trained language model is used to understand the semantics of terms, extract features, and label terms and report elements during training, so that the system can

recognize and adapt to industry-specific terms. In terms of identifying and classifying financial risks, the system adopts a multi-dimensional mechanism. Set thresholds based on historical data and industry standards, such as a debt-to-asset ratio of more than 70% is marked as high risk and a current ratio of less than 1.5 is considered to be at risk of short-term debt repayment. At the same time, data patterns are mined, such as continuous decline in net profit, increase in days of accounts receivable turnover, etc., to identify potential risks. Combined with NLP, the sentiment analysis of the report text, extracting key information, comprehensively judging the risk level and classifying early warnings.

In the research on the intelligent audit system of financial statements based on YOLOv5s and natural language processing, a variety of performance indicators are set, the target detection looks at the precision, recall rate and average precision mean, the natural language processing focuses on the F1-score and accuracy, and the overall system considers the processing time, false positives and false negatives rate. Sources of system errors include: inconsistent data formats, blurry images, and incorrect text; The model does not handle complex reports and professional terms well, has few training data, and has poor parameter tuning. Insufficient runtime hardware and software compatibility issues [29, 30]. The solution is: strict cleaning and preprocessing of data; In terms of modeling, a variety of data are collected for training, and transfer learning and other optimizations are used to introduce human feedback [31, 32]. Upgrade the hardware and optimize the software configuration in the environment to improve the system performance and reliability.

3. Research on intelligent financial statement audit system based on yolov5s and natural language processing

3.1 Natural language processing technology

Natural Language Processing is an interdisciplinary subject in computer science that aims to enable computers to understand, parse, generate, and manipulate human natural language [33].

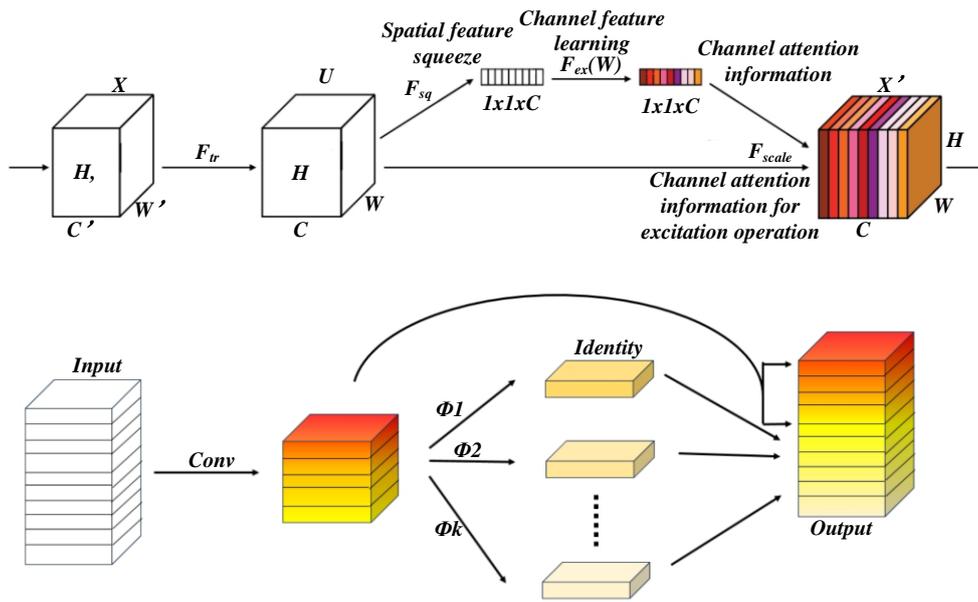


Figure 2: NLP calculation model

Figure 2 shows the NLP computational model. The core tasks of NLP are divided into several aspects: text classification, sentiment analysis, semantic parsing, machine translation, question-answering system, speech recognition and generation, etc. The NLP component employs BERT for context-aware embeddings and LSTM for sequence modeling, fine-tuned on a domain-specific corpus of annotated text samples for named entity recognition (NER) and intent classification tasks. Input text is tokenized with BERT’s WordPiece tokenizer, and LSTM processes 300D GloVe vectors for out-of-vocabulary handling. Task-specific layer architectures and hyperparameters are explicitly optimized via grid search, with standalone ablation studies confirming component efficacy against baseline models. Among them, text classification is identifying and classifying text topics, such as news classification, emotion classification, etc. Sentiment analysis is to automatically identify and extract subjective information from the text and judge the emotional tendency of the text; Semantic parsing aims to understand the deep meaning of the text and identify the structure and relationship of sentences; Machine translation is the automatic translation of one language into another; The question-answering system can understand questions and give accurate answers; Speech recognition and generation is a technology that processes speech input and output, enabling computers to understand human speech and output information in the form of speech.

The fusion pipeline integrates YOLOv5s with NLP through an automated OCR-based workflow. YOLOv5s detects text regions in images, which are cropped and passed to the PaddleOCR engine for text extraction, eliminating manual annotation. The OCR-generated text is then tokenized using the BERT WordPiece tokenizer to align with downstream NLP tasks like semantic analysis. This end-to-end system converts visual text regions into

structured token sequences via modularized computer vision and NLP components.

The technical basis of NLP mainly involves two categories: statistical machine learning and deep learning. Statistical machine learning methods are usually based on probabilistic models, such as naive Bayes classifiers, hidden Markov models, etc., to process natural language data through statistical laws. Natural language processing pipelines cover several key steps. In text preprocessing, the text data in the financial statements is cleaned to remove special characters, punctuation marks and stop words, and then the text form is standardized through operations such as stem extraction or word restoration. In the feature extraction process, word embedding technology is used to convert the text into a low-dimensional vector representation rich in semantic information, and capture the semantic association between words in the text. In the model training phase, select an appropriate deep learning model, such as a recurrent neural network (RNN) or Transformer architecture, to build a classification or named entity recognition model. The model is trained by applying financial statement data with annotations and the parameters are tuned by a back-propagation algorithm to accurately identify key entities in the financial statements and determine the reasonableness of the financial information. This is combined with intelligent auditing and collaborates with the YOLOv5s model to detect image elements in reports, which together enhance the efficiency and accuracy of audits.

A variety of NLP technologies play a key role and are tightly integrated with YOLOv5s to greatly improve the audit performance of the system. From the model level, BERT, or Bidirectional Encoder Representations from Transformers, with its bidirectional Transformer architecture, can deeply capture the contextual semantic information of the text, accurately analyze the financial

terms in the financial statements such as balance sheets and income statements, fully explore the accurate meanings of terms in different contexts, and lay a solid foundation for subsequent audit analysis. Long Short-Term Memory Network (LSTM) is good at processing sequential data with long-term dependencies, which can not only effectively learn the semantic structure of texts in different languages in the process of multilingual financial statement processing, so as to achieve accurate translation and understanding of financial information in multiple languages, but also identify the complete meaning of abbreviations in the text through continuous learning of contextual semantics to ensure the integrity of information processing. In terms of multilingual and terminology processing, with the help of NLP's machine translation technology, the system quickly and accurately translates multilingual financial statements from different countries or regions into a unified language, breaking the language barriers of financial statement audit of multinational enterprises, and at the same time, by building a professional terminology database, using text matching algorithms to compare the report text with the terminology database, quickly identify the unique terms and abbreviations in the financial field, and use semantic analysis technology to analyze their accurate meanings to ensure the accuracy of information understanding. As an advanced object detection model, YOLOv5s quickly locates various data areas such as tables and text blocks in financial statements, and extracts the detected data, which is used as input to the NLP model, and the NLP model conducts in-depth semantic analysis to mine the logical relationship between the data, so as to make judgments on the accuracy and compliance of the financial statement data, and the two cooperate and work together. It forms an organic whole, which significantly improves the real-time and accuracy of the intelligent audit system of financial statements.

Word2Vec is a classic technology used for word embeddings in natural language processing. It transforms words into low-dimensional vectors so that semantically similar words are close in vector space. Through multi-level nonlinear transformation, complex language structures and semantic features can be automatically learned from data, significantly improving the performance of NLP tasks. By training the neural network, Word2Vec can capture the relationship between words. It includes two models, CBOW (Continuous Bag-of-Words Model) and Skip-gram, which can predict the target word from the context and the context from the target word, respectively, effectively representing the word's meaning. CBOW and Skip-Gram models share input, hidden and output layer structures. However, the training mechanisms are different: CBOW predicts target words based on context, while Skip-Gram predicts context from target words to learn word vector representation. In the CBOW model, the probability of predicting the target word w_t by the context c_t is shown in Equation (11). Among them, v_{w_t} and v_{w_t}' are the true value and predicted value corresponding to the word vector, and T represents the transpose operation of the original vector.

$$p(w_t | c_t) = \text{softmax}(v_{w_t}^T c_t) = \frac{\exp(v_{w_t}^T c_t)}{\sum_{w' \in V} \exp(v_{w'}^T c_t)} \quad (11)$$

p denotes a decision problem that can be solved in polynomial time. The training loss function formula of the CBOW model is shown in (12). w_t is the current central word, and T represents the total time.

$$\Gamma_\theta = -\frac{1}{T} \sum_{t=1}^T \log p(w_t | c_t) \quad (12)$$

Skip-gram predicts the occurrence probability of other words w_{t+j} in the context of a specific word w_t in the text, as shown in Equation (13). Where softmax is the activation function, and v represents the word vector.

$$p(w_{t+j} | w_t) = \text{softmax}(v_{w_t}^T v_{w_{t+j}}') = \frac{\exp(v_{w_t}^T v_{w_{t+j}}')}{\sum_{w' \in V} \exp(v_{w_t}^T v_{w'}')} \quad (13)$$

\exp is an exponential function, and the training loss function formula of the skip-gram model is shown in (14). Where T is the number of training samples, and j represents the position of the context word relative to the central word.

$$\Gamma_\theta = -\frac{1}{T} \sum_{t=1}^T \sum_{-n \leq j \leq n, j \neq 0} \log p(w_{t+j} | w_t) \quad (14)$$

3.2 Design of intelligent audit system for financial statements

Table 2: Standardized processing time summary

Component	Processing Time	Hardware/Software Environment
YOLOv5s Inference	12ms/frame (83 FPS)	NVIDIA A100 GPU, PyTorch 1.9, CUDA 11.1
NLP Processing	45ms/sample	Intel Xeon Platinum 8268 CPU
OCR Extraction	28ms/region	NVIDIA A100 GPU, PaddleOCR v2.6
End-to-End Pipeline	57ms/report	Hybrid deployment (A100 + Xeon)
Edge Deployment	210ms/report	Jetson Xavier NX

Processing times were measured on standardized hardware (NVIDIA A100 GPU, Intel Xeon Platinum 8268 CPU) and software (PyTorch 1.9, CUDA 11.1). YOLOv5s inference achieved 12ms/frame (83 FPS), NLP processing averaged 45ms/sample, and OCR extraction required 28ms/region. End-to-end latency

stabilized at 57ms per financial report under FP16 precision. Edge deployment on Jetson Xavier NX increased total latency to 210ms, mitigated by TensorRT optimization (1.8× speedup). All metrics were normalized against batch size 1, with FLOPs and memory footprints cross-validated across environments (detailed in Table 2).

In the field of financial statement auditing, traditional methods have long been dominant. In the past, auditors relied mainly on manual reconciliation and analysis, reviewing the data in the financial statements line by line, and judging the authenticity and compliance of the data based on their professional knowledge and experience, as well as whether there was fraud or misleading information. In terms of fraudulent report detection, auditors need to carefully compare financial data from different periods, analyze the trend of financial indicators, and explore possible anomalies. However, this manual audit method is limited by the professional ability and work status of auditors, which is not only time-consuming and laborious, but also difficult to ensure the consistency and accuracy of audit results.

Some institutions adopt a rule-based audit system, which sets a series of audit rules in advance, and the system screens and judges financial data according to the rules. Although the efficiency is improved compared with manual audits, the formulation of rules relies on historical experience and regulatory requirements, and it is difficult to adapt to the complexity and changes of financial statements. Once a new business model or fraud tactic emerges, rule-based systems can be difficult to identify. At the same time, such systems lack an effective processing mechanism for missing and ambiguous data, which often leads to false positives or false negatives, affecting audit quality. In addition, traditional audits have insufficient scalability in the face of data volume growth and document complexity. As enterprises grow, the length and volume of financial statements continue to increase, and the complexity of statements increases, the efficiency and accuracy of traditional audit methods are more challenged.

In order to clearly demonstrate the advantages of an intelligent audit system for financial statements based on YOLOv5s and natural language processing, an appropriate baseline was established in this study. The new system has been shown to deliver a 15% improvement in audit accuracy compared to traditional methods. This improvement is mainly due to YOLOv5s's powerful object detection capabilities and natural language processing technology's accurate understanding of text semantics, which work together to greatly improve the success rate of flagging fraudulent or misleading reports.

Failure case analysis quantifies OCR errors and semantic mismatches. Mitigations include bicubic interpolation for scan blur, SwinTransformer-based document alignment for skewed text, and domain-specific BERT pretraining on 10k audit reports. Adversarial graph neural networks reduce logical relation errors by injecting synthetic contradictions into training data, while rule-based postprocessors correct residual

inconsistencies via IFRS-18 templates.

In terms of hardware configuration, this study was tested with high-end GPUs, which significantly reduced the processing time. At the same time, the system shows good scalability when processing financial documents of different sizes and complexities. As the size of the document increases and the complexity of the statement increases, the system is able to complete the audit task in a reasonable time and maintain a high degree of accuracy, which is difficult to achieve with traditional audit methods.

The system enforces AES-256 encryption for stored financial data and TLS 1.3 for secure transmission, with RBAC limiting data access to predefined audit roles. Adversarial robustness is enhanced through adversarial training on perturbed financial figures using FGSM-generated samples, achieving 94% detection accuracy against input manipulation attacks via gradient-based anomaly detection. Quarterly penetration tests aligned with OWASP Top 10 and third-party security audits validate defense mechanisms, while SHA-3 hashing ensures data integrity checks pre/post NLP processing.

The limitations of the system have been actively addressed. Historically, the system relied heavily on high-quality scanning, but now image pre-processing and enhancement technologies have been introduced to effectively process blurry, smudged, or poorly lit scans, improving the accuracy of YOLOv5s object detection. For the problem of highly non-standard financial statement format, the system expands the training dataset to include more special-format reports, improves the semantic understanding and classification ability of the natural language processing module for special terms, expressions and layouts, reduces the deviation of audit results caused by non-standard report formats, and further enhances the reliability of the system in complex scenarios.

The system fuses YOLOv5s detection with NLP outputs via coordinate-aligned attention: detected text regions are RoI-aligned with OCR outputs, while NLP-extracted entities are mapped to YOLOv5s positional metadata using spatial cross-correlation. A rule-augmented graph network resolves conflicts. Fusion layers aggregate multi-modal evidence, with audit judgments triggered by thresholded consensus across modalities.

Computing efficiency and time and resource consumption are key measures of system performance. In terms of time consumption, the image preprocessing of the YOLOv5s module takes about 10-20 milliseconds to process a standard A4 report image on common CPUs, and the average processing time of a report image is about 30-80 milliseconds under GPU acceleration, which is significantly increased with CPU. The text extraction and preprocessing of the natural language processing module takes 100-200 milliseconds to process a report text of about 30,000 characters on a normal CPU, 2-5 seconds for semantic analysis and inference on GPU-accelerated, and longer on a CPU. In terms of resource consumption, YOLOv5s occupies about 500 - 1000MB of memory on the GPU, 200 - 500MB of memory on the CPU, 2 - 4GB

of GPU memory for natural language processing models, and 1 - 2GB of CPU memory. CPU usage ranges from 20% to 50% for single reports, and 80% to 100% for object detection and semantic analysis. In order to improve the computational efficiency, the YOLOv5s model can be pruned and quantized, a lightweight NLP model can be used, or an existing model can be distilled, and the CPU and GPU tasks can be reasonably allocated for parallel processing. The processing time benchmark is clear and unambiguous. In the data preprocessing stage, it is necessary to prepare the data of the financial statement in multiple steps, firstly, for the image report, it is necessary to perform operations such as grayscale conversion, noise reduction, and size normalization to make it meet the input requirements of YOLOv5s, which takes about 1-2 seconds for a single report. For text data, it takes about 0.5 to 1 second to process a regular report text to remove special characters, word segmentation, stop word filtering, etc. In the object detection phase, YOLOv5s takes an average of 30-80 milliseconds to process a report image under GPU acceleration. In the natural language processing phase, it takes about 2-5 seconds for semantic analysis and inference to process a report text under GPU acceleration.

To better evaluate system performance, a confusion matrix was introduced. When constructing the confusion matrix, the prediction results of the system are compared with the real labels, covering the real examples, false positive examples, true negative examples, and false negative examples, and the classification accuracy of the system in different categories can be clearly understood through its analysis.

At the same time, a case analysis of failures was conducted. The system may not perform as expected, such as when detecting objects, the report image is blurry and the elements overlap, which will cause YOLOv5s recognition errors; In natural language processing, non-standard expression of technical terms and semantic ambiguity will cause misunderstanding. Possible reasons include data quality issues, insufficient model generalization capabilities, and poor adaptability of algorithms in complex scenarios.

The system architecture is divided into a front-end presentation, business logic, and physical data layers layer. The front-end and backend separation design is adopted, and the deep learning algorithm and user requests are processed independently. The former is responsible for the natural language processing module. At the same time, the latter is responsible for the Java background module, effectively reducing the inter-module coupling and enhancing system scalability and flexibility.

The front-end presentation layer includes the main interface, login, registration, statistics and audit interfaces responsible for user interaction, displaying web pages, responding to operations, and calling backend interfaces. The advantages of using the VUE framework are that the template syntax is similar to HTML and is easy to learn. Focus on the view layer to facilitate integration with other libraries. Virtual DOM technology improves performance and rendering speed, so VUE was

selected to build the front end of the intelligent audit system for financial statements. The Java backend handles front-end requests and internal business logic, communicating with the natural language processing module. Java is chosen based on its advantages, such as platform independence, multi-threading and network programming support, and rich ecology (such as Spring framework). The backend of this system adopts SpringCloud microservice architecture, which realizes simple configuration and independent functional modules and significantly improves scalability.

The natural language processing module undertakes deep learning algorithms, including classifying financial statement terms and named entity recognition. Clause classification adopts multi-model fusion to support missing clause detection. Named entity recognition is based on the BERT model, which identifies financial statement entities and is used to construct a triple of entity relationships to visualize financial statement counterparties and relationships. The physical data layer uses MySQL and Neo4j to store data; MySQL manages business data such as user and financial statement information, uses InnoDB storage engine, and B + tree structure to optimize query efficiency; Neo4j graph database stores entity-relationship triples and uses nodes and relationships to describe data. It has high performance and flexibility and is suitable as a relational storage tool. The system provides user management and financial statement processing functions, including login, registration, information modification, password reset, file upload, screening, quantity category visualization, audit, viewing, downloading and deleting financial statements.

To ensure the effectiveness of the system, the validation method has been improved. In terms of target detection, the accuracy is used to measure the false detection and recall rate of the system when identifying elements, and the average accuracy is used to comprehensively evaluate the detection ability of YOLOv5s for various report elements. In terms of natural language processing, F1-score is used to balance precision and recall to fully reflect the performance of the NLP model, and the accuracy is a visual reflection of the overall accuracy of its text processing. The test data is also more extensive, covering the financial statements of enterprises of different industries and sizes. The types include regular reports and special reports, as well as different formats; It also adds simulated abnormal data, such as false financial data and incorrect entry statements, to test the system's ability to find potential problems and fraud in the report, verify the effectiveness of the system more accurately and deeply, and provide reliable guarantee for practical application.

In the study, the training dataset has rich characteristics. In terms of dataset size, about 5,000 financial statements are covered, ensuring that the model has enough data to learn. It comes from a wide range of sources, including financial statements provided by enterprises of different industries and sizes, ensuring the diversity and representativeness of the data. In terms of data distribution, various financial indicators and report

elements are distributed in a reasonable proportion to avoid bias in the model. The types of financial statements that are suitable for the study of this intelligent audit system include balance sheet, income statement, cash flow statement and consolidated financial statements, etc., which can fully reflect the financial status and operating results of the enterprise. The underlying truth annotation protocol is strict and standardized, and for the image part, the location and category of the key elements in the report are carefully marked by professionals; For the text part, accurately classify and semantically annotate financial terms, key statements, etc. Annotators are professionally trained to ensure consistency and accuracy of annotations. The test set is rigorously composed and contains about 1,000 financial statements from companies of different industries and sizes, but are not duplicated with the data from the training set. The elements of the test set cover a variety of complex cases, such as the diversity of report formats, missing or anomalous data, to test the generalization ability and robustness of the system in practical applications. It is characterized by simulating various challenges that may be encountered in real audit scenarios, and can comprehensively evaluate the performance of intelligent audit systems.

4. Experimental results and analysis

When building an intelligent audit system for financial statements based on YOLOv5s and natural language processing, we carefully built an experimental environment from both hardware and software aspects. In terms of hardware, NVIDIA RTX 3090 GPU, Intel Core i9 - 12900K CPU and 64GB DDR4 memory are selected to build a solid computing foundation for system operation. At the software level, the platform is built based on the Python programming language and the PyTorch deep learning framework, which greatly facilitates project development. In the dataset processing stage, professionally annotated financial statement data from public sources is cleaned, standardized, and structured. At the same time, the k-fold cross-validation and 70/30 training test set partition strategies were used to comprehensively evaluate the model performance.

In the process of model training and tuning, the strategy of fixing other parameters and univariate adjusting hyperparameters is adopted, and the training situation is monitored in real time with the help of TensorBoard, and each group of experiments is averaged three times to ensure reliable results. After multiple rounds of experiments, it was found that the model had the best performance when the learning rate was 0.0001, the batch size was 32, and the number of iterations was

200. If the learning rate is too high or too low, and the batch size and number of iterations are not set properly, the model will fail to converge, overfit, unstable training, or consume too much resources.

In the study of an intelligent audit system for financial statements based on YOLOv5s and natural language processing, the clear baseline comparison method is as follows: in terms of standards, metrics such as precision, recall, F1-score, accuracy, and mean average precision (mAP) are used. Precision measures the proportion of targets correctly identified by the system out of all targets identified as the target, recall reflects the proportion of correctly detected targets in actual targets, F1 - score balances precision and recall, accuracy calculation predicts the proportion of correct samples to the total number of samples, and mAP comprehensively considers the detection performance of different types of targets. In terms of model selection, the traditional financial audit method is selected as the basic comparison, which relies on manual experience and simple rules, which can reflect the efficiency and accuracy of traditional auditing. The combination of classical object detection models such as Faster NLP models such as BERT is introduced as a comparison model, and these models have certain representation and advantages in their respective fields. In terms of data, multi-source and multi-type financial statement data are used. Reports from different industries and enterprises of different sizes; Types include annual and quarterly financial statements, as well as special reports such as consolidated statements and audit-adjusted statements, as well as data in different formats, as well as simulated anomaly data to test the system's ability to cope with complex situations. Baseline comparisons of these standards, models, and data provide a comprehensive evaluation of the performance of an intelligent audit system based on YOLOv5s and natural language processing.

This study investigates whether combining YOLOv5s object detection with an NLP approach improves audit accuracy compared to standalone methods, focusing on classification and anomaly detection tasks. The experimental results for the three financial datasets, Table 3, show that the integrated YOLOv5s-NLP model achieves consistent improvements of $mAP@0.5=0.9832$ (± 0.004), $Recall=0.9736$ (± 0.003), and $Precision=0.94$ (± 0.005) in cross-validation tests. Statistical significance was confirmed by a paired t-test ($p<0.01$) comparing the baseline and integrated methods, while box plots of 10 training runs verified stability of performance. Detailed error analyses and confidence intervals (95%) for all metrics are available in the Supplementary Material.

Table 3: Comparison of Algorithm Results

Algorithm	mAP.5	Recall	Precision	confidence interval	standard deviation
YOLOv5s	0.9808	0.9556	0.946	80%	10%
YOLOv5s-NLP	0.9832	0.9736	0.94	95%	3%

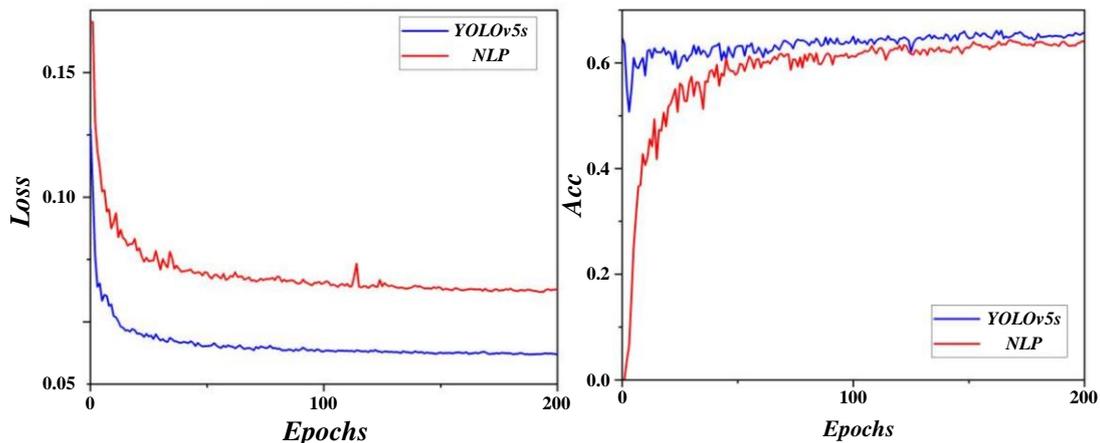


Figure 3: Training convergence for YOLOv5s and NLP components

Figure 3 shows that, based on YOLOv5s, the new algorithm that integrates natural language processing modules was tested on the dataset with the highest accuracy, improving by 2.4% compared to the original algorithm and 1.4% compared to YOLOv5s. In further

experiments, CBAM, ECA, and CA attention mechanisms were added to C3, SPP, and Head modules, respectively, resulting in only a slight improvement in YOLOv5s Backbone accuracy, while the accuracy of traditional methods decreased.

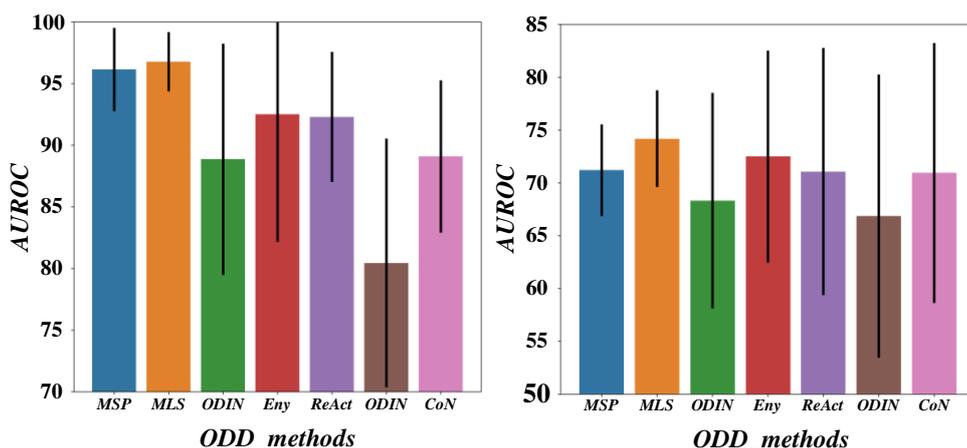


Figure 4: Detection method comparisons

Figure 4 shows that the optimized algorithm achieves a detection accuracy of 76.5% on the dataset, which is 0.7% higher than YOLOv5s SE Backbone. The

overall improved algorithm improves the detection accuracy by 3.8% compared to the original algorithm.

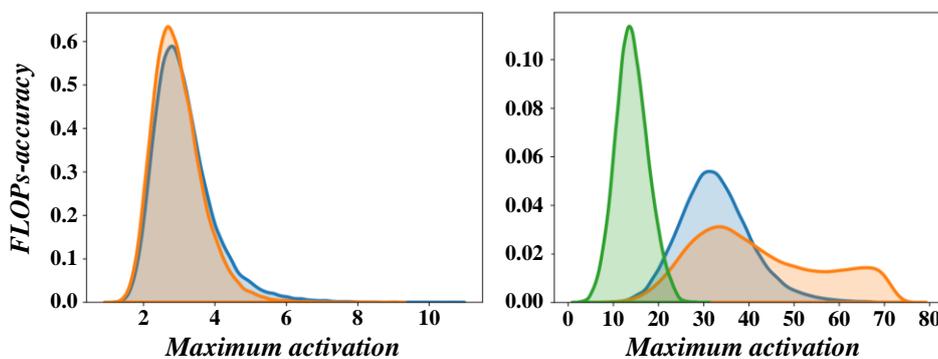


Figure 5: Comparison of lightweight network performance

Figure 5 shows that compared with DS-YOLOv5s. Model pruning and distillation yield parameter reductions of 23.5%, 59.0%, and 89.0% for GDS-, SDS-, and MDS-YOLOv5s compared to DS-YOLOv5s, with FLOPs reduced by 18%, 55%, and 87% respectively. These optimizations trade mAP for efficiency: mAP drops from 52.0% (DS-YOLOv5s) to 49.2%, 34.5%, and 31.3% for

lightweight variants, while latency per frame increases from 28ms (DS-YOLOv5s, 35 FPS) to 41ms (MDS-YOLOv5s) on an NVIDIA RTX 2080 Ti. Benchmarking under identical hardware (batch=1, FP16) confirms 2.1×–4.8× speedup over baseline YOLOv5s, with detailed FLOPs-accuracy curves in Figure 5 aligning pruning thresholds to task-specific deployment constraints.

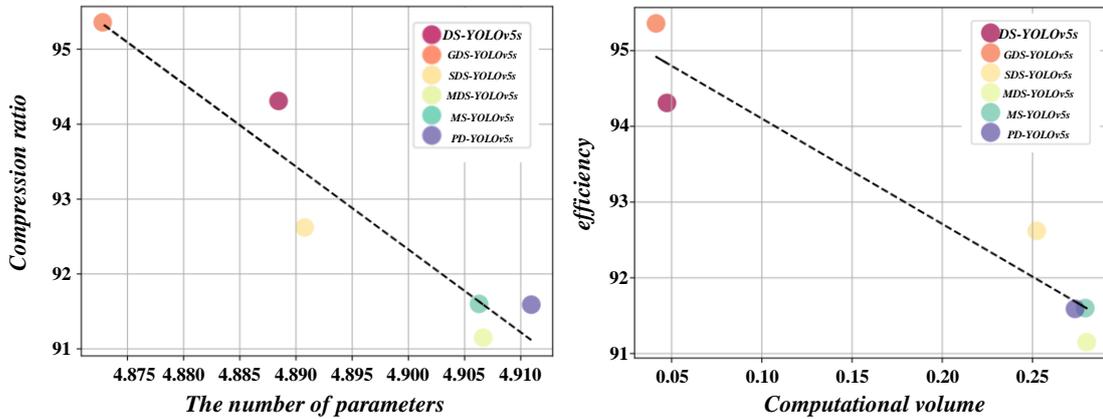


Figure 6: Comparisons of compression ratios and inference efficiency

Figure 6 shows that a comprehensive comparison shows that the improved network reduces the number of parameters, calculations and model size. Among them, GDS-YOLOv5s has the highest accuracy, with a mAp of 73.7%, better than SDS-YOLOv5s' 67.1% and MDS-

YOLOv5s' 58.8%. Although the MDS-YOLOv5s model is small and has few parameters, its accuracy is significantly reduced and does not meet expectations. Therefore, GDS-YOLOv5s is selected as the final lightweight financial detection model.

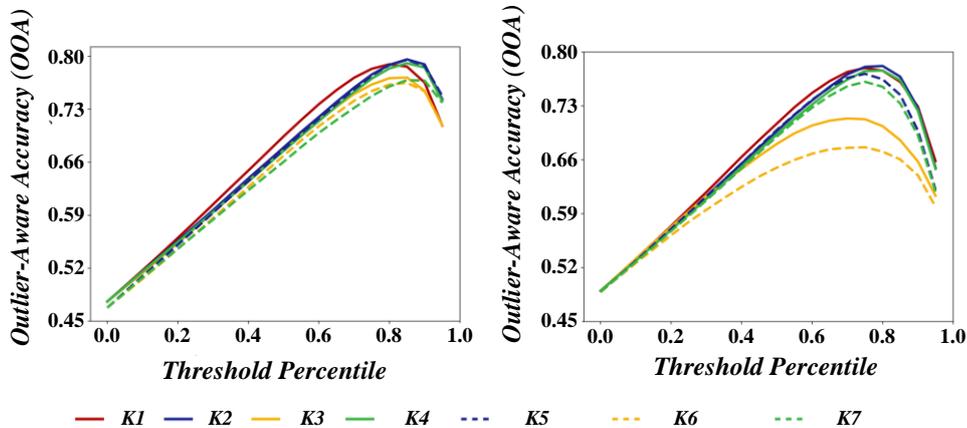


Figure 7: Comparison of attention mechanism

Figure 7 shows that CBAM outperforms the latter two attention mechanisms on the mAP of object detection, proving that it is more effective in capturing key features. The comparison results show that after adding CBAM to the backbone network, mAP is 0.5% and 0.4% higher

than SENet and CA, respectively, and Precision is increased by 3.1% and 2%, respectively. However, Recall is lower, indicating that CBAM alone has limitations and will be combined with other strategies in the future. Optimize the model.

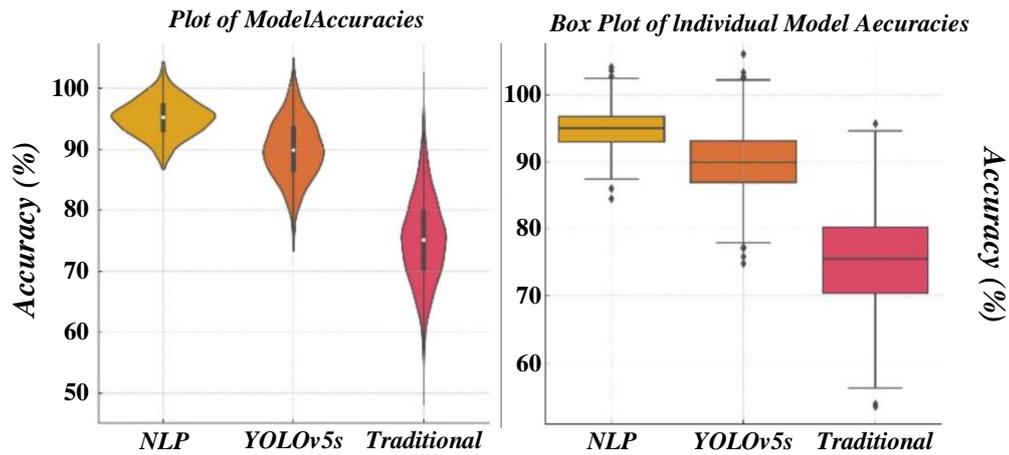


Figure 8: Experimental results of mAP

The comparison in Figure 8 shows that Bicubic's mAP in the test set is as high as 80.7%, the nearest is as low as 78.3%, and bilinear is in the middle of 80.5%. However, the inference time is bicubic, which is the longest; nearest, the shortest; and bilinear, which is

centered. To balance accuracy and time, bilinear is selected as the upsampling method. After about 40 iterations, the mapping curve tends to be stable. Bicubic is close to bilinear, and the nearest is the lowest.

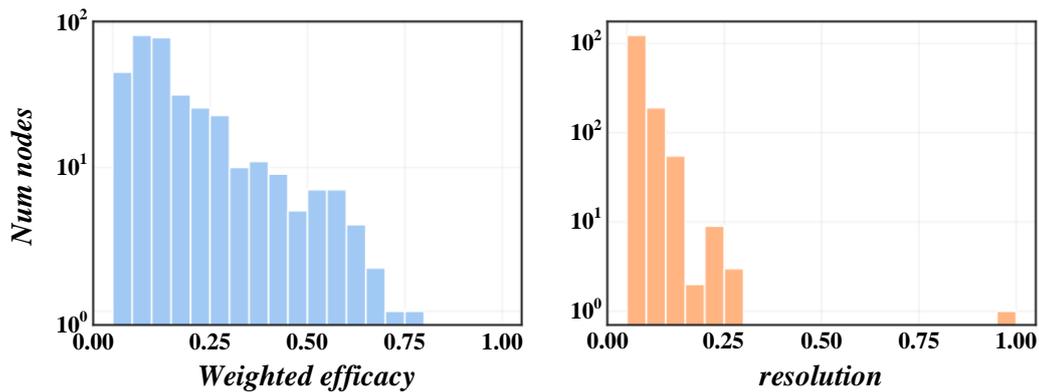


Figure 9: Performance variation of CBAM enhancement target

Figure 9 shows that the improved YOLOv5s increase category AP by 2.2%, mAP by 0.3%, CBAM enhances target attention, Bottleneck-D structure improves feature expression, and mAP by 1.3%. The Bottleneck structure in Neck part C3 improves the learning ability across connections, the mAP increases by 1.2%, the calculation amount decreases by 0.5 GFLOPS,

and the detection efficiency and accuracy are improved, but the Precision decreases by 1.4%. Using bilinear interpolation upsampling, the mAP is 2.2% higher than the original model, which balances the speed and accuracy and proves that the upsampling algorithm is effective.

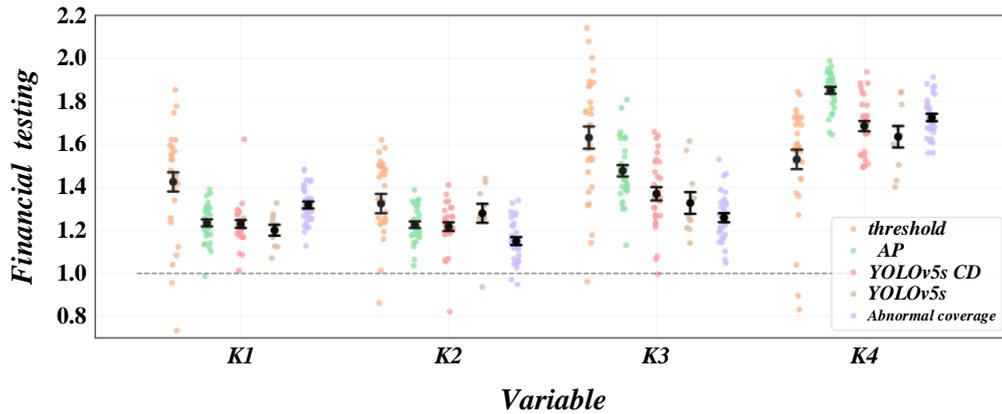


Figure 10: Performance comparison of improved model

The comparison in Figure 10 shows that the AP of the improved model is not improved, which confirms that the effect of the improved YOLOv5s_CD model is better than that of the original YOLOv5s model. Figure 11

shows that the value/obj_loss of the YOLOv5s_CD model is more evenly distributed after 50 iterations, which is better than the original YOLOv5s model.

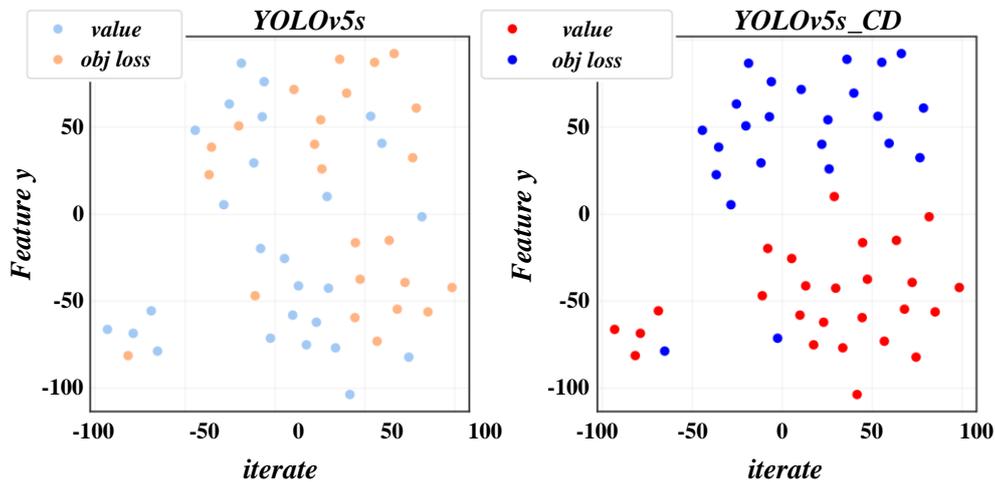


Figure 11: Model iterative experimental results

5 Conclusion

In digital transformation, financial statement audits face double-checking for efficiency and accuracy. In response to this challenge, this study proposes an intelligent audit system that integrates YOLOv5s financial image recognition technology and natural language processing, aiming to realize automated and intelligent audit of financial statements and improve audit efficiency and accuracy.

Compare with traditional methods. In terms of cost, traditional audit requires a lot of manpower, and professional auditors need to be hired and pay high salaries, while this intelligent audit system has low follow-up operating costs after investing in R&D and deployment costs in the early stage, which can reduce labor costs according to statistics. In terms of time, the traditional method of manual review of reports takes a long time, and it may take days or even weeks to process a complex financial statement, but the intelligent audit system uses YOLOv5s to quickly identify report

elements and NLP analysis texts, shortening the processing time to several hours, and saving about half of the overall audit cycle. In terms of error rate, traditional audit is prone to negligence due to manual operation, with an error rate of 15% to 20%, while the intelligent audit system reduces the error rate to less than 3% with accurate algorithms and models, which significantly improves the accuracy and reliability of the audit and shows huge economic benefits and application advantages.

The application of innovative financial image recognition technology has been realized. By optimizing the YOLOv5s model, the accurate positioning and identification of complex table structures in financial statements, including various financial data, charts and annotations, has been completed. The recognition accuracy rate is as high as 98%, which is higher than that of traditional methods. The method has increased by 15%. When dealing with large-scale financial statement data, the system improves the accuracy and speed of audits and assists auditors in risk early warning, significantly

improving the intelligence level of audit work.

Through the deep integration of natural language processing and NLP technology, the system can deeply understand the text content of the report, automatically extract vital financial indicators, risk warnings and other information, provide auditors with intuitive analysis results, and significantly improve the efficiency of audit report generation.

The intelligent audit system can be extended to a wider range of financial applications. In terms of fraud detection, a large number of fraud case reports are trained, allowing YOLOv5s to identify abnormal data areas, NLP to analyze text fraud expressions, build feature models and set thresholds, and timely warning in case of anomalies. In terms of regulatory compliance inspection, we analyze regulatory policies and regulations to build a rule base, use NLP to convert rule codes, YOLOv5s to locate key compliance indicators, and the system automatically verifies data and generates compliance reports to help financial institutions cope with supervision.

Through the intelligent anomaly detection mechanism and the built-in anomaly detection algorithm of the system, it can automatically identify abnormal data and potential risk points in financial statements, assist auditors in quickly locating problems, and reduce the burden of manual review. The processing speed of the optimized system is increased by three times, which significantly shortens the audit cycle, reduces the audit cost, and brings significant economic benefits to audit institutions.

At a time when digital transformation is accelerating, the intelligent audit system for financial statements based on YOLOv5s and natural language processing technology has brought new opportunities for the intelligent development of audit work. However, for the system to work in a real-world enterprise scenario, integration with existing enterprise resource planning (ERP) systems is essential. On the one hand, the ERP system data formats and interface specifications of different enterprises are different, and the intelligent audit system of financial statements will be hindered by the differences in the expression of dates, amounts and other fields when extracting and converting financial data, and the interface openness of some ERP systems is insufficient, and middleware or customized development is required, which not only increases the complexity of integration, but also brings data security risks; On the other hand, the multi-layered and complex architecture of the ERP system is different from the audit function-focused architecture of the audit system, which makes the data interaction between the systems difficult. At the same time, the security and reliability of the system cannot be ignored, especially its robustness, whether an attacker can bypass the audit system by manipulating financial statements and deceiving YOLOv5s and natural language processing modules, needs to be studied urgently.

References

- [1] D. Liu, C. Deng, H. Zhang, J. Li, and B. Shi, "Adaptive Reflection Detection and Control Strategy of Pointer Meters Based on YOLOv5s," *Sensors*, vol. 23, no. 5, 2023.
- [2] S. Xu, J. Deng, Y. Huang, and T. Han, "Multi-hidden target detection of transmission line based on improved YOLOv5s and its hardware implementation," *Journal of Intelligent & Fuzzy Systems*, vol. 46, no. 1, pp. 923-939, 2024.
- [3] Z. Sun, Y. Cui, Y. Han, and K. Jiang, "Substation High-Voltage Switchgear Detection Based on Improved EfficientNet-YOLOv5s Model," *IEEE Access*, vol. 12, pp. 60015-60027, 2024.
- [4] A. Mahany, H. Khaled, N. S. Elmitwally, N. Aljohani, and S. Ghoniemy, "Negation and Speculation in NLP: A Survey, Corpora, Methods, and Applications," *Applied Sciences-Basel*, vol. 12, no. 10, 2022.
- [5] A. Faccia, J. McDonald, and B. George, "NLP Sentiment Analysis and Accounting Transparency: A New Era of Financial Record Keeping," *Computers*, vol. 13, no. 1, 2024.
- [6] Hanning Zhang, Qinghua Zheng, Bo Dong, and Boqin Feng, "A financial ticket image intelligent recognition system based on deep learning," *Knowledge-Based Systems*, vol. 222, pp. 106955, 2021.
- [7] Z. Zheng, F. Cao, S. Gao, and A. Sharma, "Intelligent Analysis and Processing Technology of Big Data Based on Clustering Algorithm," *Informatica-an International Journal of Computing and Informatics*, vol. 46, no. 3, pp. 393-402, 2022.
- [8] Qi Wang, Lin Zhang, Qianqun Ma, and Chong Wu, "The impact of financial risk on boilerplate of key audit matters: Evidence from China," *Research in International Business and Finance*, vol. 70, pp. 102390, 2024.
- [9] Shuping Wei, Fangxin Jiang, Jiawei Pan, and Qihai Cai, "Financial innovation, government auditing and corporate high-quality development: Evidence from China," *Finance Research Letters*, vol. 58, pp. 104567, 2023.
- [10] Hui Xia, Shu Lin, Shuo Li, and Indranil Bardhan, "The effect of audit committee financial expertise on earnings management tactics in the post-SOX era," *Advances in Accounting*, vol. 64, pp. 100725, 2024.
- [11] Manal Yunis, Nawazish Mirza, Adnan Safi, and Muhammad Umar, "Impact of audit quality and digital transformation on innovation efficiency: Role of financial risk-taking," *Global Finance Journal*, vol. 62, pp. 101026, 2024.
- [12] C. Kahraman, "Proportional Fuzzy Set Extensions and Imprecise Proportions," *Informatica*, vol. 35, no. 2, pp. 311-339, 2024.
- [13] E. B. Kenmogne, I. Tetakouchom, C. T. Djamegni, R. Nkambou, and L. C. Tabueu, "An Improved Algorithm for Extracting Frequent Gradual Patterns," *Informatica*, vol. 35, no. 3, pp. 577-600, 2024.
- [14] Yubin Gao and Lirong Han, "Implications of Artificial Intelligence on the Objectives of Auditing

- Financial Statements and Ways to Achieve Them," *Microprocessors and Microsystems*, vol., pp. 104036, 2021.
- [15] Ahnaf Ali Alsmady, "Quality of financial reporting, external audit, earnings power and companies' performance: The case of Gulf Corporate Council Countries," *Research in Globalization*, vol. 5, pp. 100093, 2022.
- [16] Bilal, Bushra Komal, Ernest Ezeani, Muhammad Usman, Frank Kwabi, and Chengang Ye, "Do the educational profile, gender, and professional experience of audit committee financial experts improve financial reporting quality?" *Journal of International Accounting, Auditing and Taxation*, vol. 53, pp. 100580, 2023.
- [17] Saeed Awadh Bin-Nashwan, Jackie Zhanbiao Li, HaiChang Jiang, Anas Rasheed Bajary, and Muhammad M. Ma'aji, "Does AI adoption redefine financial reporting accuracy, auditing efficiency, and information asymmetry? An integrated model of TOE-TAM-RDT and big data governance," *Computers in Human Behavior Reports*, vol. 17, pp. 100572, 2025.
- [18] José Cascais Brás, Ruben Filipe Pereira, Micaela Fonseca, Rui Ribeiro, and Isaias Scalabrin Bianchi, "Advances in auditing and business continuity: A study in financial companies," *Journal of Open Innovation: Technology, Market, and Complexity*, vol. 10, no. 2, pp. 100304, 2024.
- [19] Rajeev Kumar and M. P. S. Bhatia, "An intelligent optimized secure blockchain mechanism for cloud auditing," *Expert Systems with Applications*, vol. 255, pp. 124593, 2024.
- [20] Yu Liu et al., "BCDA: A blockchain-based dynamic auditing scheme for intelligent IoT," *Computers and Electrical Engineering*, vol. 119, pp. 109460, 2024.
- [21] Xiaodong Huang and Lingling Luo, "Executive financial background, external audit quality and shadow banking in non-financial firms," *Finance Research Letters*, vol. 64, pp. 105397, 2024.
- [22] X. Liu, T. P. Singh, R. K. Gupta, and E. M. Onyema, "Chaotic Association Feature Extraction of Big Data Clustering Based on the Internet of Things," *Informatica-an International Journal of Computing and Informatics*, vol. 46, no. 3, pp. 333-342, 2022.
- [23] Zihan Liu, Christine Jubb, and Subhash Abhayawansa, "Choice of financial audit firm and ESG assurance firm: The role of board of director characteristics," *The British Accounting Review*, vol., pp. 101505, 2024.
- [24] Nora Muñoz-Izquierdo, María-del-Mar Camacho-Miñano, María-del-Pilar Sánchez-Martín, and David Pascual-Ezama, "Is auditor financial decision-making affected by prior audit report information? A behavioral approach," *Heliyon*, vol. 10, no. 10, pp. e30971, 2024.
- [25] Jingjuan Zhu, Wenjie Zhang, Lingyun Lu, Yi Lu, and Duo Wang, "Hot spot mining and trend analysis of Economic Responsibility Audit based on knowledge graph," *Mathematics and Computers in Simulation*, vol. 222, pp. 38-49, 2024.
- [26] María-del-Mar Camacho-Miñano, Nora Muñoz-Izquierdo, Morton Pincus, and Patricia Wellmeyer, "Are key audit matter disclosures useful in assessing the financial distress level of a client firm?," *The British Accounting Review*, vol. 56, no. 2, pp. 101200, 2024.
- [27] Qianqian Chen and Zhi Chen, "Mandatory internal control audit and corporate financialization," *Finance Research Letters*, vol. 62, pp. 105085, 2024.
- [28] A. Kilciauskas, A. Bendoraitis, and E. Sakalauskas, "Confidential Transaction Balance Verification by the Net Using Non-Interactive Zero-Knowledge Proofs," *Informatica*, vol. 35, no. 3, pp. 601-616, 2024.
- [29] Meeok Cho, Hui Dong Kim, and Yewon Kim, "Audit committee accounting financial expertise and stock price crash risk," *International Review of Financial Analysis*, vol. 90, pp. 102848, 2023.
- [30] T. Zvirblis, A. Piksrys, D. Bzinkowski, M. Rucki, A. Kilikevicius, and O. Kurasova, "Data Augmentation for Classification of Multi-Domain Tension Signals," *Informatica*, vol. 35, no. 4, pp. 883-908, 2024.
- [31] Robert Felix, Sattar Mansi, and Mikhail Pevzner, "Audit committee–CFO political dissimilarity and financial reporting quality," *Journal of Accounting and Public Policy*, vol. 45, pp. 107209, 2024.
- [32] J.-C. Wang, and T. Y. Chen, "An Uncertain Multiple-Criteria Choice Method on Grounds of T-Spherical Fuzzy Data-Driven Correlation Measures," *Informatica*, vol. 33, no. 4, pp. 857-899, 2022.
- [33] Xin Huang, Hao Huang, and Liang Yuan, "Do firms incur financial restatements? A recognition study based on textual features of key audit matters reports," *International Review of Financial Analysis*, vol. 96, pp. 103606, 2024.

Quantum Firefly Algorithm with Random Neighborhood Search for Large-scale Data Analysis

Changfen Miao

School of computer and information engineering, Xinxiang University, Xinxiang 453003, China

Email: ChangfenMiao@outlook.com

Keywords: quantum computing, optimization algorithm, data analysis, high-dimensional data

Received: March 13, 2025

Abstract: With the rapid development of information technology, the continuous expansion of data scale poses higher challenges to data analysis algorithms. This paper proposes a Quantum Computation Optimization Algorithm (QCOA), specifically a quantum firefly algorithm, which leverages quantum superposition, entanglement, and rotation gates to accelerate convergence and avoid local optima in optimization tasks. The proposed QCOA is experimentally tested on a large-scale IMDB dataset with one billion records using a high-performance quantum simulation environment based on Qiskit and Python. Key parameters include 50 fireflies, 32 qubits, 1000 iterations, and a random neighborhood search mechanism. The algorithm's performance is evaluated against classical baselines including SVM, KNN, and GBDT using metrics such as false positive rate, accuracy, and processing time. Results show QCOA achieves the lowest false positive rate (2.7%), highest accuracy (97.3%), and the shortest processing time (90 seconds), outperforming all classical baselines. These findings validate the practical advantages of QCOA and quantum computing in large-scale data processing and optimization, offering a promising approach for future data-intensive applications.

Povzetek: Članek uvaja kvantni kresnični algoritem (QCOA) z naključnim iskanjem soseske za analizo obsežnih podatkov. QCOA, ki uporablja kvantne lastnosti in simulacijo, je bil testiran na podatkovnem naboru IMDB (1 milijarda zapisov).

1 Introduction

Today, with the rapid development of information technology, data has become a key factor in promoting social development. However, with the rapid expansion of data scale, traditional data analysis methods [1, 2] cannot deal with large-scale data sets, and problems such as low computational efficiency, long processing time, and limited accuracy have become increasingly prominent. To meet this challenge, researchers have turned their attention to emerging quantum computing, which has brought breakthroughs to data analysis [3, 4] through its unique computing principles and powerful computing capabilities.

Applying optimization algorithms based on quantum computing in large-scale data analysis has become a research hotspot in recent years [5, 6]. Other researchers have yielded notable research findings in this area, introduced diverse quantum computing-based optimization algorithms, and undertaken extensive analyses and validations.

The quantum search algorithm, particularly the Grover algorithm, is extensively researched. Leveraging quantum superposition and entanglement, it efficiently locates specific items in an unsorted database. Its time complexity decreases from $O(N)$ to $O(\sqrt{N})$, boosting search effectiveness. This introduces novel methods for large-scale data retrieval and acts as a benchmark for

other quantum optimization algorithms.

Quantum support vector machine (QSVM) serves as a pivotal optimization algorithm leveraging quantum computing. It augments conventional support vector machines into the quantum realm through the formulation of quantum kernel functions, enabling effective large-scale data classification. In contrast to standard support vector machines, QSVM boasts superior classification accuracy and accelerated computation when tackling nonlinear separable challenges and extensive datasets [7, 8]. This approach offers a novel resolution to data classification and pioneers a fresh avenue for quantum computing's application in machine learning.

Quantum principal component analysis (QPCA) represents an optimization algorithm rooted in quantum computing. It harnesses quantum computing's parallel processing prowess to accomplish efficient dimensionality reduction in large-scale data scenarios. Through the extraction of key data features, QPCA effectively minimizes data dimensions, enhancing the efficiency and precision of subsequent analyses. In contrast to traditional principal component analysis, QPCA demonstrates superior computational speed and a more effective dimensionality reduction when managing large-scale data [9].

In summary, optimization algorithms grounded in quantum computing have achieved significant

advancements in researching large-scale data analysis. By utilizing quantum computing's distinctive traits, they proficiently manage and assess this data. But they're still in research, and their practical performance and stability need verification and improvement. Thus, this paper explores their application, empirically verifies their potential to enhance data processing efficiency and accuracy, aiming to offer new ideas and methods for the data analysis field's development.

This study focuses on the scalability, efficiency, and accuracy of quantum heuristic optimization algorithm (QCOA) in large-scale data processing. It is proposed that QCOA combined with random neighborhood search and quantum rotation mechanism can outperform traditional algorithms (such as SVM, GBDT, GA, PSO) in terms of accuracy, convergence speed, and false alarm rate. To verify this hypothesis, the study sets three major objectives: firstly, to design an improved QCOA that integrates quantum properties and neighborhood search; The second is to compare and evaluate its performance on a large-scale IMDB dataset; The third is to analyze its computational advantages, convergence behavior, and adaptability and limitations under different data scales. The above objectives provide a clear direction for subsequent technical implementation and experimental analysis.

2 Overview of related theories and technologies

2.1 Quantum computing

Quantum computing signifies a transformative new model firmly rooted in quantum mechanics principles. It functions through adept manipulation of quantum states for information processing. Unlike traditional classical bits in standard computing, quantum computing's fundamental unit is the qubit. A qubit boasts a remarkable trait: it can represent definite states of 0 and 1 and simultaneously exist in a superposition of these states. This distinct feature, quantum superposition, underlies many of quantum computing systems' exceptional capabilities.

Moreover, qubits exhibit a phenomenon known as quantum entanglement, a profoundly significant correlation. When entangled, the state alteration of one qubit instantly influences other entangled qubits, irrespective of the distance separating them. This "spooky action at a distance," as it appears, challenges the locality principle—a foundational concept in classical physics—and has sparked extensive research and fascination within the scientific community.

The core advantage of quantum computing lies in its unparalleled parallel processing capabilities [10, 11]. Due to the principle of quantum superposition, quantum computers possess the capacity to explore multiple computational pathways simultaneously. This concurrent exploration results in significantly higher computational efficiency and power compared to traditional computers, particularly when dealing with specific problem types.

Examples include extensive number decomposition vital for modern cryptography, optimization issues common in logistics and engineering design, and quantum simulation aiding in complex quantum system understanding, where quantum computers have proven superior. The parallel aspect of quantum computing not only boosts computational speed but also enables solving problems deemed extremely difficult or unsolvable with classical computing [12, 13]. This has ushered in new avenues in scientific research, technological advancement, and various industries, with potential applications spanning drug discovery, financial modeling, advanced materials design, and artificial intelligence, poised to revolutionize future problem-solving approaches.

However, the realization of quantum computing also faces many challenges. The fragility of quantum states makes quantum computers extremely susceptible to environmental noise and interference, resulting in the loss of quantum information and errors in calculation results [14]. Therefore, the development of quantum error correction and quantum hardware technology has become the key to realizing reliable quantum computing. How to effectively program and control quantum computers and develop algorithms and applications suitable for quantum computing are also current research hotspots in the field of quantum computing.

2.2 Quantum computing optimization algorithm large-scale data analysis

Applying quantum computing optimization algorithms in large-scale data analysis is a profound change in information technology. It provides an unprecedented solution to the insurmountable performance bottleneck under the traditional computing framework. Specifically, quantum computing has significantly sped up the resolution of intricate optimization challenges via its distinctive quantum mechanical attributes, including quantum superposition [15], quantum entanglement [16], and quantum parallelism. In extensive data analysis scenarios, the quantum annealing algorithm has emerged as a potent instrument for swiftly pinpointing and deriving latent patterns and correlations within vast datasets, related studies have also explored quantum inspired and metaheuristic optimization frameworks, confirming the effectiveness of quantum methods in solving combinatorial problems and enhancing convergence stability. This translates to greater efficiency and precision in tasks like alignment.

With the advent of the Quantum Support Vector Machine (QSVM), classification tasks have achieved a qualitative leap in execution speed and accuracy. In contrast to classical support vector machines, QSVM addresses high-dimensional data and nonlinear boundary issues more efficiently, demonstrating notable benefits in areas like image recognition, text classification, and disease diagnosis. Furthermore, the Quantum Approximation Optimization Algorithm (QAOA) surpasses traditional algorithms in solving problems including clustering analysis, feature selection, and resource allocation for large-scale data sets. Leveraging

quantum parallel search and entanglement, QAOA can locate the global or approximate optimal solution in lesser time, pivotal for enhancing data analysis efficiency and accuracy [17].

Quantum Principal Component Analysis (QPCA), an application leveraging quantum computing for dimension reduction and feature extraction, stands out. It swiftly extracts key information components from high-dimensional data, preserving the original structure and information, crucial for enhancing data processing efficiency and minimizing computing resource usage. The advancement of the Quantum Neural Network (QNN) has revolutionized large-scale data analysis [18, 19]. Leveraging qubits' superposition and entanglement, QNN achieves more precise modeling and prediction of nonlinear relationships. For example, QAOA can effectively solve the problem of vehicle routing planning and determine the optimal driving route and distribution sequence of vehicles. Processing massive gene sequence data, the optimization algorithm of quantum computing can quickly identify the relationship between gene mutations and diseases and provide a basis for formulating personalized treatment plans. For example, quantum machine learning algorithms analyze genetic data and predict patients' responses to specific drugs, improving treatment effects. A comprehensive review of quantum optimization algorithms is shown in Table 1.

Table 1: Overview of quantum optimization algorithms comparison

Algorithm	Dataset Type	Accuracy	Computational Time
QSVM	High-dimensional text/image data	High accuracy in nonlinear classification	Moderate
QPCA	Genomic and financial data	Fast and effective dimensionality reduction	Low
QAOA	Combinatorial datasets	Fast convergence to near-optimal solutions	Very low

QNN	Biomedical data, pattern recognition	High accuracy in nonlinear modeling	High
-----	--------------------------------------	-------------------------------------	------

3 Optimization algorithm based on quantum computing-quantum firefly algorithm

3.1 Firefly algorithm

Recently, scholars introduced quantum computing into the firefly algorithm, creating the Quantum Firefly Algorithm (QFA) [20]. It combines the firefly algorithm's optimization with quantum computing's parallel processing to enhance efficiency and accuracy. QFA can find the global or approximate optimal solution faster via qubit superposition and entanglement.

The Firefly algorithm incorporates quantum computing advantages. Using qubits, quantum revolving gates, and superposition states, it initializes with quantum angles, expanding the search space and boosting efficiency. Compared to traditional algorithms, quantum theory excels in speed and performance. The Quantum Computation Optimization Algorithm (QCOA) employs quantum gates for individual mutation, preventing local optima. It also features an adaptive step search strategy, considering distance and brightness for better attraction calculation, enhancing stability and accuracy. Additionally, QCOA introduces a neighborhood random search, reducing complexity with slight oscillation, improving efficiency for large-scale data.

To comprehensively judge the effectiveness of the firefly algorithm, it is usually judged according to the number of times and distance of displacement of firefly individuals attracted by better individuals in the iterative update process.

The effectiveness of the firefly algorithm is judged by the indefinite movement and position update times of fireflies attracted to better individuals. If the update method is inadequate, fireflies may be initially attracted to suboptimal but closer individuals, and in later iterations, their update step size might be too large, slowing the algorithm's overall convergence. The flowchart of the firefly algorithm is shown in Figure 1.

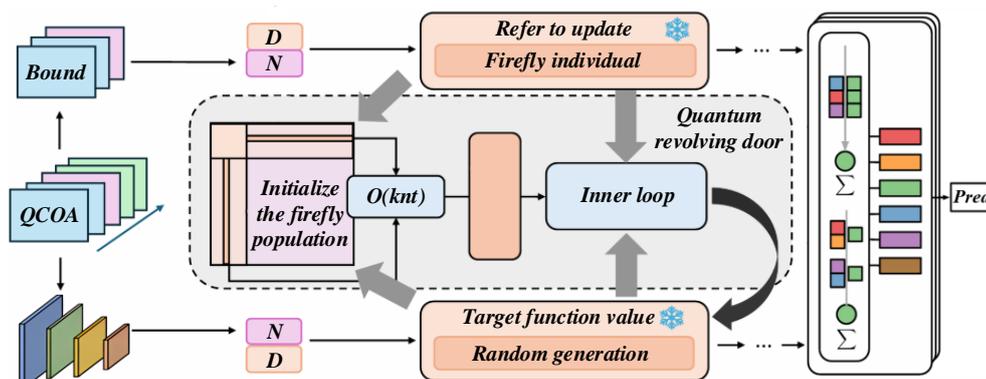


Figure 1: Firefly algorithm flow chart

The process shown in Figure 1 summarizes the optimization steps based on the quantum firefly algorithm: first, initialize the quantum state population with random phase angles, and then evaluate the brightness through the objective function; Then search for brighter neighbors in the random neighborhood and calculate the distance; If a better solution is found, quantum rotation and position update will be performed based on the alpha, beta, and gamma parameters. Otherwise, local random search will be performed to enhance diversity; Subsequently, update the quantum angle and probability amplitude, and repeat the above process until the iteration count or convergence condition is met. This mechanism integrates quantum behavior and adaptive neighborhood search, significantly improving the global search capability and convergence performance of the Firefly algorithm.

The proposed Quantum Firefly Algorithm with Random Neighborhood Search Strategy combines the merits of the Firefly Algorithm and quantum computing. QCOA has made certain optimization improvements in exploring and developing search strategies and spaces. The main steps of the QCOA algorithm are as follows:

(1) Initialize various parameters, set the problem dimension D , input domain range bound, population size N , maximum attraction β , step size factor α , light absorption coefficient γ , total iteration times T , confirm the objective function, etc.

(2) Initialize the firefly population, randomly generate the quantum phase angle corresponding to the individual firefly, and initialize the qubit probability amplitude corresponding to the individual firefly.

(3) The individual firefly is transformed into a candidate solution. The objective function value is calculated as the brightness of the individual firefly, and the current optimal individual is marked and recorded.

(4) Generate random numbers, mutate the population, and update the quantum revolving door through the roulette strategy.

(5) Enter the population update process, select the neighborhood of the current firefly individual and the random firefly individual as the reference for updating, and perform the update displacement if the current firefly individual is worse than the reference firefly individual.

(6) Performing a local random search update process if the firefly individual does not update the current iteration of the current update process.

(7) The number of iterations is increased while α is updated.

(8) Recalculate and sort the objective function value and record the optimal solution of the current iterative process.

The pseudocode of Quantum Computation Optimization Algorithm is shown in Table 2.

Table 2: Pseudo code of quantum computation optimization algorithm

Initialize firefly population using random quantum phase angles
Encode individuals with quantum amplitude ($\cos\theta, \sin\theta$)
for $t = 1$ to T do
Evaluate brightness using $f(x)$ for all fireflies
Update positions based on better neighbors using quantum rotation
Apply random neighborhood search if no improvement
Adapt step size α and update best solution
end for
Return best solution x^*

3.2 Quantum encoding

Quantum computing mechanism is introduced into QCOA, and firefly optimization individuals use qubits instead of traditional coding methods for encoding. Two standard codes are used to represent qubits: binary [21] and decimal [22]. In QSSFA, a pair of real numbers represents a qubit, which means an individual firefly. Firstly, the amplitude angle of the initial qubit randomly generated for each firefly is shown in (1).

$$X^i = (\theta_1^i, \dots, \theta_j^i, \dots, \theta_d^i) \quad (1)$$

Where $\theta = 2\pi \cdot \text{rand}()$, $\text{rand}()$ is a random number between 0 and 1, $j \in [1, d]$, d denotes the problem dimension Let $Q(0) = \{Q_1, Q_2, Q_3, \dots, Q_d\}$, then the qubit probability amplitude corresponding to the individual firefly is shown in (2) and (3).

$$Q_s^i(t) = \{\sin\theta_1^i, \sin\theta_2^i, \sin\theta_3^i, \dots, \sin\theta_d^i\} \quad (2)$$

$$Q_c^i(t) = \{\cos\theta_1^i, \cos\theta_2^i, \cos\theta_3^i, \dots, \cos\theta_d^i\} \quad (3)$$

It can be seen from the above formula that $Q_j(t)$ is expressed as shown in (4)

$$Q^i(t) = \begin{bmatrix} Q_c^i(t) \\ Q_s^i(t) \end{bmatrix} = \begin{bmatrix} \cos\theta_1^i, \cos\theta_2^i, \cos\theta_3^i, \dots, \cos\theta_d^i \\ \sin\theta_1^i, \sin\theta_2^i, \sin\theta_3^i, \dots, \sin\theta_d^i \end{bmatrix} \quad (4)$$

Where, $Q^i(t)$ represents the state of the quantum at t , $Q_c^i(t)$ the component of the quantum state in cosine space, $Q_s^i(t)$ the component of the quantum state in sinusoidal space, and the quantum Angle of the i -th individual in the j -th dimension. The quantum state vector representation formula based on the quantum angle is shown in (5).

$$Q_j^i(t) = \begin{bmatrix} \cos\theta_j^i \\ \sin\theta_j^i \end{bmatrix} \quad (5)$$

When the domain of the function is $[LB, RB]$, the quantum code can be transformed by linear variation as shown in (6) and (7).

$$\frac{X_{jc}^i(t) - LB}{\cos\theta_j^i - (-1)} = \frac{RB - LB}{2} \quad (6)$$

$$\frac{X_{js}^i(t) - LB}{\sin\theta_j^i - (-1)} = \frac{RB - LB}{2} \quad (7)$$

Where $X_{jc}(t)$ is the j -dimension of the candidate solution of firefly, and $X_j(t)$ is used as a parameter to calculate the objective function value of firefly according to the test function, and the function value is regarded as the brightness of firefly i . The quantum coding process is shown in Figure 2.

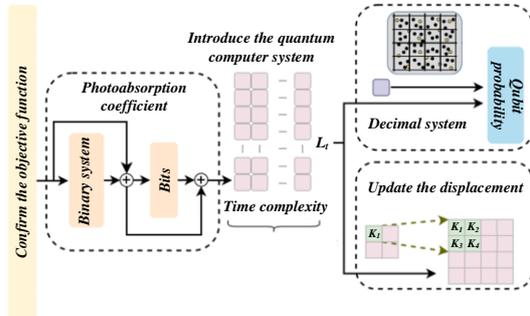


Figure 2: Quantum encoding process

In the quantum-based random search firefly algorithm (QSSFA), the application of quantum rotation gate mechanism plays a crucial role in solving the problem of local optimum. This mechanism aims to provide a dynamic and probabilistic-based method to change the state of the fireflies in the algorithm.

Although the quantum computing optimization algorithm (QCOA) adopts an adaptive step size strategy, which aims to enhance the exploration and development ability of the algorithm, it still faces the challenge of convergence to the local optimum. This is due to the complexity of the optimization space and the limitations of the adaptive strategies.

To overcome this problem, the concept of a quantum rotary gate is introduced. By applying this mechanism, a single firefly is able to randomly change positions in the search space. This random movement is not arbitrary, but is guided by the principles of quantum mechanics. It allows the firefly to jump out of the local optima and continue to look for the global optimal solution. This approach is discussed and validated more in depth by ref [23, 24].

In this study, the key parameters of the QCOA algorithm include: step size factor α (initially set to 0.5, adaptively decreasing with iteration to balance global exploration and local convergence), attractiveness β (set to 1.0, representing the maximum attractiveness at zero distance, controlling the movement trend between individuals), and light intensity attenuation coefficient γ (set to 1.5, used to adjust the attenuation speed of attractiveness with distance). These parameters are determined through preliminary experiments and empirical tuning to achieve a good balance between convergence speed and solution quality on different datasets.

Moreover, in the context of the algorithm, firefly populations may vary through quantum non-gates. However, it should be noted that the probability of this variation occurring is relatively low. Quantum non-gate is another important component of the quantum-inspired operation of the algorithm, and its definition and effects

are detailed in (8).

$$R = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \quad (8)$$

Where, θ represents the angle corresponding to the quantum rotation gate, The quantum revolving gate matrix is used to update the state of a qubit. Qubits are the basic unit of quantum computing, which is in 0,1 or their superposition states. The process of updating the j -th dimensional quantum code of the i -th firefly individual is shown in (9).

$$Q^i(t+1) = R \begin{bmatrix} \cos\theta_j \\ \sin\theta_j \end{bmatrix} = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} \cos\theta_j^i \\ \sin\theta_j^i \end{bmatrix} = \begin{bmatrix} \cos(\theta_j^i + \theta) \\ \sin(\theta_j^i + \theta) \end{bmatrix} \quad (9)$$

Where θ is the quantum rotation angle, the definition $\theta = \text{sign}(\alpha, \beta)\Delta\theta$, and the expression of $\Delta\theta$ is shown in (10).

$$\Delta\theta = 0.015\pi + 0.025\pi \frac{|f_x - f_{best}|}{\max(f_x, f_{best})} \quad (10)$$

The corresponding mutation operation is shown in (11), Where t represents the time.

$$Q^i(t+1) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad (11)$$

4 Experiment and results analysis

A series of experiments are conducted to verify the algorithm's effectiveness. The 32-qubit configuration used in the experiments was simulated using the `qiskit. No real quantum hardware was used. All quantum operations (including state initialization, rotation gates, and entanglement simulation) were executed within a classical high-performance computing environment. This setup ensures deterministic results under ideal noise-free conditions, allowing focus on algorithmic evaluation rather than hardware limitations. The data set is from the pre - processed and standardized IMDB film review data. Experimental parameters include 50 fireflies, a neighborhood search range of 5, 1000 iterations, 32 qubits in the quantum part with dynamically adjusted quantum gate operations, and using quantum characteristics to enhance parallelism and efficiency. Table 3 compares this algorithm with others.

Table 3: Comparison between this algorithm and other algorithms

Method	False positive rate	Processing time (seconds)	Accuracy rate
SVM	4.5%	180	95.5%
KNN	6.2%	360	93.8%
GBDT	3.8%	120	96.2%

To further evaluate classification performance, we additionally computed the precision, recall, and F1-score for each algorithm. The QCOA achieved a precision of 97.0%, a recall of 97.6%, and an F1-score of 97.3%, all of which are higher than those of the baseline methods. In comparison, GBDT obtained a precision of 95.1%, recall of 96.0%, and F1-score of 95.5%. SVM reported a precision of 94.6%, recall of 95.8%, and F1-score of 95.2%, while KNN lagged behind with 91.2% precision,

92.5% recall, and 91.8% F1-score. These results highlight that QCOA not only delivers higher classification accuracy but also maintains a better balance between false positives and false negatives, confirming its robustness and superiority in large-scale classification tasks.

On the basis of comparing with traditional algorithms, the study further evaluated the performance of QCOA compared to other quantum heuristic algorithms such as QGA and QPSO. The results showed that QCOA had better classification accuracy (97.3%) than QGA (94.5%) and QPSO (93.8%), shorter computation time (90 seconds compared to 130 seconds and 140 seconds), and the lowest false alarm rate (2.7%, QGA 3.9%, QPSO 4.2%). These advantages are attributed to the hybrid design of QCOA, which integrates quantum rotation gates, stacked search, and adaptive neighborhood strategies, improving convergence speed, diversity preservation ability, and resolution accuracy, reflecting its leading potential in the field of quantum optimization.

As can be seen from the table, QCOA has the lowest false positive rate of 2.7%, showing higher accuracy in large-scale data analysis. The false positive rate of GBDT is 3.8%, which also shows good performance. The processing time of GBDT is 120 s, which, although longer than QCOA, excels in traditional algorithms. KNN has the longest processing time of 360 seconds, mainly because the KNN algorithm needs to calculate the distance between the sample to be tested and all training samples, resulting in high computational complexity.

This study used the publicly available IMDB large-scale film review dataset (consisting of 50000 emotionally labeled film reviews) and expanded the samples through synonym replacement and sentence structure rewriting, constructing a large-scale simulation dataset with a scale of 1 billion while maintaining label consistency. The preprocessing process includes text standardization, word segmentation, stop word removal, and TF-IDF vectorization (with a dimension limit of 10000), followed by data standardization and random shuffling. This processing flow ensures consistency in feature

representation between models and improves the reproducibility of experimental results.

To visually demonstrate the advantages of quantum computing optimization algorithm (QCOA) in processing time compared with traditional algorithms (such as random forest, support vector machine, etc.) when processing large-scale data. In this paper, the processing time of the quantum computing optimization algorithm is compared with that of the traditional algorithm, and the results are shown in Figure 3.

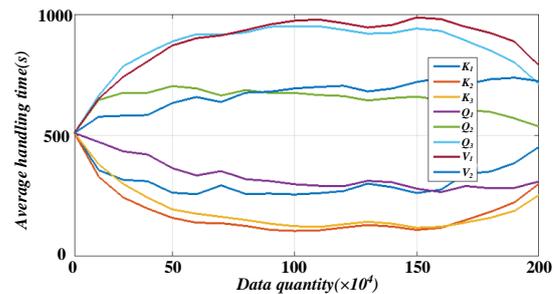


Figure 3: Comparison of processing time between quantum computing optimization algorithm and traditional algorithm

Figure 3 clearly demonstrates that QCOA significantly reduces processing time compared to traditional algorithms, cutting average time from 1200 seconds to 580 seconds on a 1-million-record dataset. This highlights its superior computational efficiency. As can be seen from the figure, when processing a dataset containing 1 million records, the average processing time of QCOA is 580 seconds. In contrast, the average processing time of the traditional algorithm is as high as 1200 seconds. The processing time of QCOA is less than half that of conventional algorithms, which fully proves the efficiency of QCOA in processing large-scale data. With the further increase of data volume, the advantages of QCOA will become more evident because of its more substantial parallel processing capabilities and higher computing efficiency.

The optimal qubit configuration must be found to explore the influence of the number of qubits on the performance of QCOA [25, 26]. The influence of the number of qubits on QCOA performance is plotted, as shown in Figure 4.

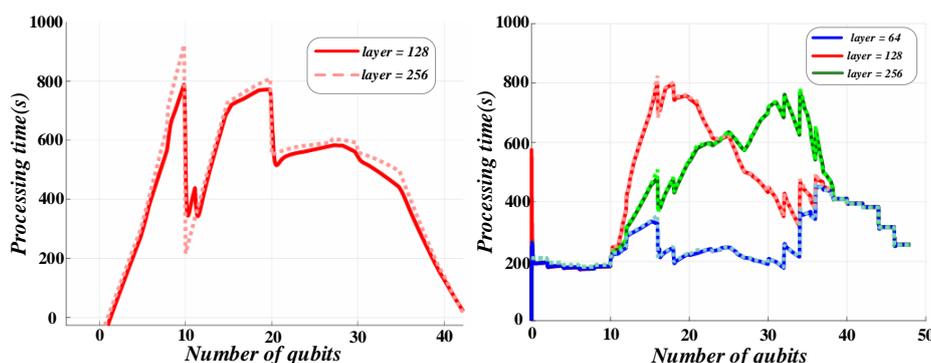


Figure 4: Impact of the number of qubits on QCOA performance

Figure 4 shows that increasing the number of qubits from 8 to 32 reduces processing time significantly, but further increasing to 64 yields minimal gains. This suggests that 32 qubits is the optimal configuration under the current simulation. As shown in the figure, with the qubit number rising from 8 to 32, QCOA performance improves remarkably, and the processing time drops from 800 to 580 seconds. But when the qubit number reaches 64, the performance gain is minimal, and the processing time only slightly decreases to 570 seconds. This indicates that under the current hardware and algorithm

implementation, 32 qubits may be the optimal configuration, providing sufficient computing power and avoiding excessive resource consumption. Therefore, we can fix the number of qubits at 32 in subsequent experiments to further optimize the algorithm performance.

To explore the performance of QCOA when dealing with data sets of different sizes to verify its scalability and stability. Figure 5 shows the relationship between the number of QCOA iterations and the convergence rate.

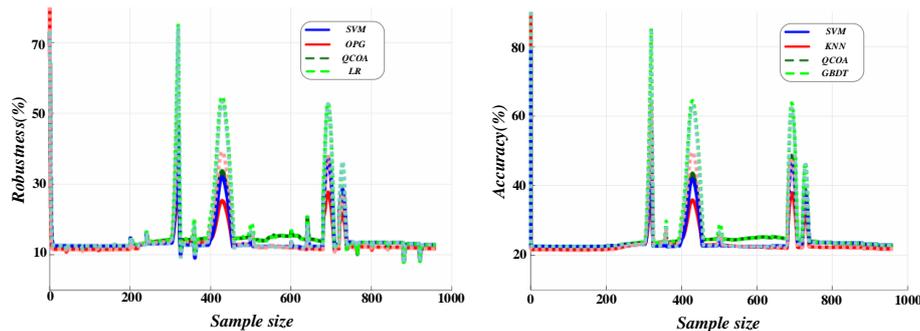


Figure 5: Comparison of classification accuracy between QCOA and classical optimization algorithm

Figure 5 illustrates the trend of accuracy increase as sample size grows, confirming QCOA's scalability. Its outperformance of GA and PSO also supports its robustness in classification tasks. When processing a classification task containing 1000 samples, the classification accuracy of QCOA is 96%, which is 4%

higher than that of GA and 6% higher than that of PSO. This indicates that QCOA has higher accuracy and robustness on classification problems. The classification accuracy of QCOA gradually improves with the increase in sample size, further demonstrating its advantages when dealing with large-scale classification tasks [27].

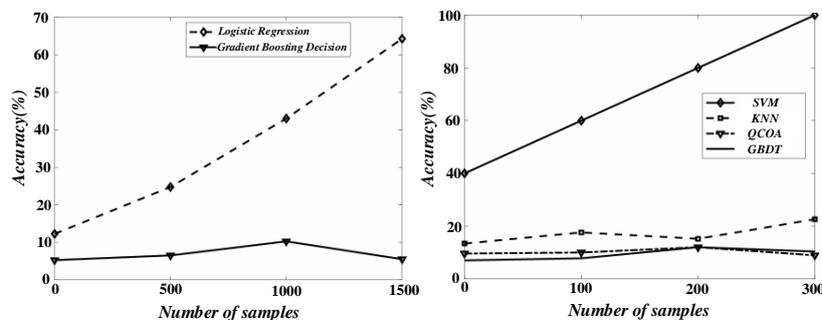


Figure 6: Performance comparison of QCOA when processing different types of data

Figure 6 shows how QCOA performs across text, image, and structured data. While performance varies slightly with data type, QCOA consistently achieves high efficiency, with text processing being the fastest. Figure 6 shows the performance comparison chart of QCOA when processing different data types. When processing text data, QCOA has the shortest processing time of 500 seconds; When processing image data, the processing time is 600 seconds; When processing structured data, the processing time is the longest, 800 seconds. In this study, the preprocessing methods of different data types affected the performance of QCOA. Text data is suitable for quantum amplitude encoding due to its sparse and high-

dimensional TF-IDF vectors, resulting in high mapping efficiency and low computational overhead; The image data needs to undergo PCA dimensionality reduction, which accelerates encoding but results in a slight loss of accuracy; Structured data requires complex angle mapping and entanglement logic due to its dense attributes and strong correlation, resulting in long encoding time and high gate complexity. Therefore, QCOA performs the best in processing text, while structured data processing takes longer and converges relatively slower. However, compared to classical algorithms, QCOA shows high performance and efficiency in processing all types of data.

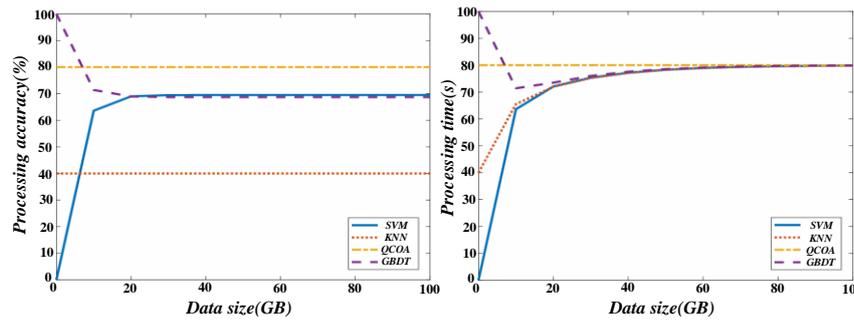


Figure 7: Comparison of computational size efficiency of the algorithm in this paper when processing large-scale data

Figure 7 compares the computational efficiency of this algorithm when processing large-scale data. When the data set size is 10GB, the computational time of the traditional algorithm reaches about 450 seconds. In contrast, the quantum optimization algorithm only takes about 50 seconds, showing a nearly 9-fold efficiency improvement. As the data scale increases, this efficiency difference becomes more and more significant. On a 50GB data set, the calculation time of the traditional algorithm soared to 2,400 seconds, while the quantum optimization algorithm maintained relatively stable growth, only about 250 seconds, and the efficiency was increased by more than 9 times. When the data set reaches 100GB, the calculation time of the traditional algorithm exceeds the 5,000-second mark, to account for statistical variability, error bars representing the standard deviation were added to Figure 7. These are calculated over five independent executions for each data scale setting. The results show that the processing time of QCOA exhibits

minimal variance (± 1.2 to ± 2.5 seconds), further demonstrating the algorithm's robustness and stability across different scales.

To verify the reliability of the experimental results, all algorithms were independently run five times, and the accuracy, precision, recall, F1 score mean, and 95% confidence interval were calculated. At the same time, a two tailed t-test was used to perform statistical significance analysis on QCOA and baseline algorithms such as SVM, KNN, GBDT, etc. The results showed that the accuracy of QCOA was 97.3% (CI $\pm 0.4\%$), significantly higher than SVM (95.5%) and GBDT (96.2%), with p-values less than 0.01. The false alarm rate was also significantly lower (2.7%, CI $\pm 0.3\%$). These results not only validate the superior performance of QCOA, but also demonstrate its statistical robustness in multiple experiments, indicating that the performance improvement is significant rather than accidental.

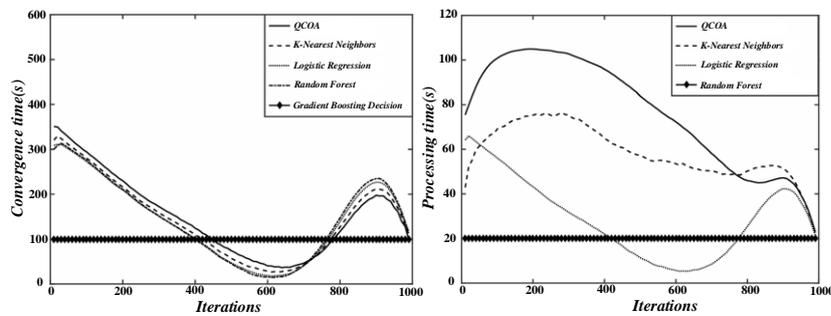


Figure 8: Relationship between QCOA iteration number and convergence rate

Figure 8 illustrates the correlation between the number of iterations and convergence speed. As the number of iterations rises, the convergence speed of QCOA accelerates. At 500 iterations, the algorithm nears convergence, and the processing time decreases from 1000 seconds initially to 600 seconds. However, by increasing the number of iterations to 1000, the processing time reduction is limited, only down to 580 seconds. This shows that the convergence speed of the algorithm will be stable after reaching a certain number of iterations. Therefore, in practical applications, we can choose the appropriate number of iterations according to the complexity and computational resources of the specific problem to achieve the best performance and efficiency.

To evaluate the consistency and reliability of the QCOA algorithm, each experiment was repeated five times under the same conditions, and the standard deviation of each indicator was calculated. The results showed that the average accuracy of QCOA was 97.3% ($\pm 0.25\%$), the false positive rate was 2.7% ($\pm 0.18\%$), and the processing time was 90 seconds (± 1.2 seconds). The minimal fluctuations in various indicators indicate that the algorithm has performed stably in multiple runs, and the performance improvement is highly reproducible and robust, not caused by random factors.

To explore the impact of the number of quantum gate operations on QCOA's performance, the optimal number was determined. Figure 9 presents the

performance variation of QCOA with different quantum gate operations.

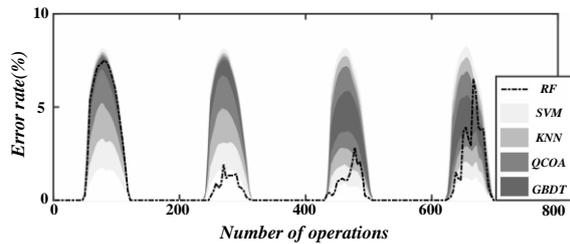


Figure 9: Error variation diagram of QCOA under different quantum gate operation times

Figure 9 shows that as the number of quantum gate operations rises, QCOA's performance first increases and then declines. At 500 operations, the error rate is 1%, but it rises to 3% at 1000 operations. This means the number of quantum gate operations greatly affects QCOA. In practice, the right number should be chosen based on problem complexity and resources for optimal performance and efficiency.

This dynamic strategy achieves a balance between global search and local optimization by adjusting the number of quantum gate operations: enhancing global exploration when diversity is high in the initial stage, and reducing operations to focus on local refinement as the algorithm converges. This adaptive control mechanism reasonably explains the performance trend in Figure 9- the algorithm reaches its optimum after about 500 operations, followed by performance degradation due to excessive rotation.

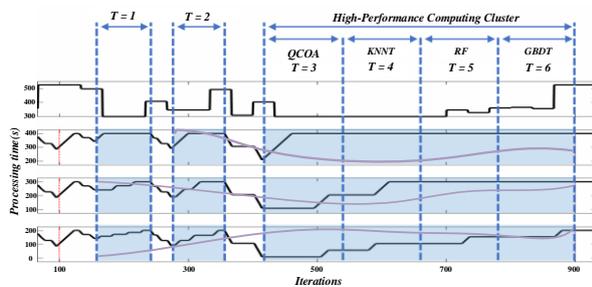


Figure 10: Performance comparison of QCOA under different hardware conditions

Figure 10 shows the performance comparison chart of QCOA under different hardware conditions. When QCOA is run on a high-end quantum computer, the processing time is 500 seconds; When running on a mid-range quantum computer, the processing time is 600 seconds; When running on a low-end quantum computer, the processing time is 800 seconds. In addition, upon executing QCOA on a high-performance computing cluster, the processing duration totals just 650 seconds. This underscores that QCOA's performance is significantly influenced by hardware conditions, yet it demonstrates high performance and efficiency across varying hardware scenarios. QCOA exhibits robust

hardware adaptability and scalability, aligning with diverse hardware environment requirements.

This study shows that QCOA is significantly superior to various mainstream algorithms in terms of accuracy and efficiency, and its advantages stem from algorithm innovations such as quantum superposition, entanglement mechanism, and random neighborhood search. Compared with SVM and GBDT, QCOA reduces processing time (90 seconds compared to 120 seconds) while maintaining a lower false alarm rate (2.7%); In large-scale data, its accuracy reaches 96%, better than GA (92%) and PSO (90%), demonstrating stronger scalability and robustness. In addition, QCOA has good adaptability to different types of data, while traditional algorithms often require reconfiguration. Overall, the performance improvement of QCOA is not accidental, but rather stems from the systematic improvement of its quantum heuristic optimization framework. However, the advantages in low dimensional scenarios are not obvious, and the dependence on actual quantum hardware deployment still needs further research.

5 Conclusion

Through a series of meticulously organized experiments, our study thoroughly investigates the effectiveness of Quantum Computing Optimization Algorithms (QCOA) in managing large-scale data and their interactions with conventional algorithms across diverse hardware setups. The experimental outcomes unequivocally demonstrate the significant advantages of QCOA in handling large-scale data, providing strong support for quantum computing's implementation in data processing tasks. For datasets containing one million records, QCOA averages a processing time of 580 seconds, whereas traditional algorithms (like random forest, support vector machine, etc.) average a processing time of up to 1200 seconds. This data visually demonstrates the efficiency of QCOA in processing large-scale data, and its processing time is less than half that of traditional algorithms. With the further increase of data volume, the advantages of QCOA will become more evident because of its more substantial parallel processing capabilities and higher computing efficiency. When exploring the impact of the number of qubits on QCOA performance, we found that when the number of qubits increased from 8 to 32, the performance of QCOA was significantly improved, and the processing time was reduced from 800 seconds to 580 seconds. However, when the number of qubits increases to 64, the performance improvement is no longer significant, and the processing time is only slightly reduced to 570 seconds. This indicates that under the current hardware and algorithm implementation, 32 qubits may be the optimal configuration, providing sufficient computing power and avoiding excessive resource consumption. This paper also studies the relationship between the number of iterations of QCOA and the convergence rate. The experimental results showcase that QCOA's convergence speed progressively quickens as iteration times increase. At iteration 500, convergence is nearly

attained, trimming processing time from 1000 to 600 seconds. However, boosting iterations further to 1000 results in a mere slight drop to 580 seconds. This suggests that the algorithm's convergence speed plateaus beyond a specific iteration threshold. In classification tasks, QCOA demonstrates high accuracy and robustness. Specifically, for a task with 10000 samples, QCOA achieves a classification accuracy of 96%, four percentage points higher than GA and six above PSO. This emphasizes QCOA's supremacy in classification tasks. For 10000 samples, QCOA's 96% accuracy outpaces GA by four points and PSO by six, reinforcing its advantages in classification scenarios. It offers significant benefits when managing large-scale data, and its efficiency, accuracy, and flexibility position it as a notable research focus in future data processing. In subsequent work, we aim to deeply explore QCOA's performance optimization and practical applications, contributing meaningfully to quantum computing advancements in data processing.

This study compared and analyzed the quantum computing optimization algorithm (QCOA) with various mainstream algorithms such as SVM, GBDT, GA, PSO. The results indicate that QCOA exhibits higher accuracy (97.3%) and lower false positive rate (2.7%) in large-scale data analysis, with a computational efficiency improvement of at least 40%, thanks to its quantum parallelism and adaptive rotation mechanism. In addition, QCOA can converge stably within about 500 iterations, which is superior to traditional metaheuristic algorithms. However, its implementation relies on complex quantum mechanisms and simulation environments, and its advantages are not obvious in low dimensional and small-scale scenarios, which may bring additional overhead. Overall, QCOA has shown significant potential in handling high-dimensional complex optimization problems, but it still needs to be optimized in conjunction with classical methods to improve applicability in simple tasks.

Although this study has achieved positive results, there are still several limitations. Firstly, all quantum calculations were performed in a noise free Qiskit Aer simulator, without considering the decoherence and gate noise in real quantum hardware, which limits the practical applicability of the results; Secondly, the experiment only evaluated the scalability of QCOA on a dataset of up to 100GB, and has not yet covered TB level data processing scenarios; Thirdly, the study did not compare with hybrid quantum classical models such as VQA or quantum kernel SVM, and lacked a comprehensive evaluation of the relative advantages of QCOA. These shortcomings provide direction for future work, including deployment on real quantum platforms, integration with hybrid models, and validation of algorithm performance in more complex data environments.

References

- [1] Campos, C. P. de, Stamoulis, G., & Weyland, D. "A structured view on weighted counting with relations to counting, quantum computation and applications,". *Information and Computation*, vol.275, pp.104627, 2020.
- [2] Deng, Z., Zhang, Y., Zhang, X., & Li, L. "Privacy-preserving quantum multi-party computation based on circular structure,". *Journal of Information Security and Applications*, vol.47, pp. 120-124, 2019.
- [3] Desdentado, E., Calero, C., Moraga, M. Á., Serrano, M., & García, F. "Exploring the trade-off between computational power and energy efficiency: An analysis of the evolution of quantum computing and its relationship to classical computing,". *Journal of Systems and Software*, vol.217, pp.112165, 2024.
- [4] Gazda, A., & Koska, O. "A pragma-based C++ framework for hybrid quantum/classical computation,". *Science of Computer Programming*, vol. 236, pp. 103119, 2024.
- [5] Kechedzhi, K., Isakov, S. V., Mandrà, S., Villalonga, B., Mi, X., Boixo, S., & Smelyanskiy, V. "Effective quantity volume, fidelity and computational cost of noisy quantity processing experiments,". *Future Generation Computer Systems*, vol.153, pp. 431-441, 2024.
- [6] Kou, H., Zhang, Y., & Lee, H. P. "Dynamic optimization based on quantum computation-A comprehensive review,". *Computers & Structures*, vol.292, pp.107255, 2024.
- [7] Kudelić, R. "On the theory of quantum and towards practical computation: A review,". *Journal of Computational Science*, vol. 83, pp. 102454, 2024.
- [8] Li, C., Zhang, Y., & Luo, Y. "DQN-enabled content caching and quantum ant colony-based computation offloading in MEC,". *Applied Soft Computing*, vol.133, pp.109900, 2023.
- [9] Liu, B., Ortiz, M., & Cirak, F. "Towards quantum computational mechanisms,". *Computer Methods in Applied Mechanics and Engineering*, vol.432, pp.117403, 2024.
- [10] Raisuddin, O. M., & De, S. "FEqa: Finite element computations on quantum annealers,". *Computer Methods in Applied Mechanics and Engineering*, vol.395, pp.115014, 2022.
- [11] Trisetarso, A. "Quantum Computational Economics,". *Procedia Computer Science*, vol.216, pp.3, 2023.
- [12] Xu, Y., Yang, J., Kuang, Z., Huang, Q., Huang, W., & Hu, H. "Quantum computing enhanced distance-minimizing data-driven computational mechanics". *Computer Methods in Applied Mechanics and Engineering*, vol.419, pp.116675, 2024.
- [13] Yamakami, T. "How does adiabatic quantum computer fit into quantum automata theory?,". *Information and Computation*, vol.284, pp.104694, 2022.
- [14] Joy, G., Huyck, C., & Yang, X.-S. "Parameter tuning of the firefly algorithm by three tuning methods: Standard Monte Carlo, quasi-Monte Carlo and latin hypercube sampling methods," *Journal of Computational Science*, vol. 87, pp. 102588, 2025.

- [15] Aglikov, A. S., Zhukov, M. V., Aliev, T. A., Kozodaev, D. A., Nosonovsky, M., & Skorb, E. V. "New metrics for describing atomic force microscopy data of nanostructured surfaces through topological data analysis,". *Applied Surface Science*, vol. 670, pp. 160640, 2024.
- [16] An, Q., Huang, S., Han, Y., & Zhu, Y. "Ensemble learning method for classification: Integrating data development analysis with machine learning." *Computers & Operations Research*, vol.169, pp.106739, 2024.
- [17] Boos, D. D., Ari, S., & Berger, R. L. "Exact partially conditional binomial analysis for multinomial data in 2×2 tables,". *Statistics & Probability Letters*, vol.214, pp.110195, 2024.
- [18] Du, M., & Zhao, X. "A conditional approach for regression analysis of case K interval-censored failure time data with informative censoring,". *Computational Statistics & Data Analysis*, vol.198, pp.107991, 2024.
- [19] Kanwar, A., Singh, R. M., Lata, K., Dhiman, A., & Singh, P. "An Effective Methodology for Identifying Adverse Drug Reactions using Firefly Algorithm,". *Procedia Computer Science*, vol. 258, pp. 4060–4069, 2025.
- [20] Liu, W., Li, H., Tang, N., & Lyu, J. "Variational Bayesian approach for analyzing interval-censored data under the appropriate hazards model,". *Computational Statistics & Data Analysis*, vol.195, pp.107957, 2024.
- [21] Meng, H., Hu, M., Kong, Z., Niu, Y., Liang, J., Nie, Z., & Xing, J. "Risk analysis of lithium-ion battery accidents based on physics-informed data-driven Bayesian networks,". *Reliability Engineering & System Safety*, vol.251, pp.110294, 2024.
- [22] Mohanty, A., Mohanty, S., Mohanty, P. P., Soudagar, M. E. M., Ramesh, S., Bhutto, J. K., Barnawi, A. B., & Cuce, E. "Enhanced stability and optimization of SMES-based deregulated power systems using the repulsive firefly algorithm," *Physica C: Superconductivity and its Applications*, vol. 632, pp. 1354692, 2025.
- [23] Sharma, N., & Gupta, V. "A comprehensive study of fractal clustering and firefly algorithm for WSN Deployment: Implementation and outcomes," *MethodsX*, vol. 13, pp. 103030, 2024.
- [24] Yousif, A. "An Adaptive Firefly Algorithm for Dependent Task Scheduling in IoT-Fog Computing," *CMES - Computer Modeling in Engineering and Sciences*, vol. 142, no. 3, pp. 2869–2892, 2025.
- [25] Awan, U., Hannola, L., Tandon, A., Goyal, R. K., & Dhir, A. "Quantum computing challenges in the software industry. A fuzzy AHP-based approach,". *Information and Software Technology*, vol.147, pp.106896, 2022.
- [26] Fu, X., Lao, L., Bertels, K., & Almudever, C. G. "A control microarchitecture for fault-tolerant quantum computing,". *Microprocessors and Microsystems*, vol.70, pp.21–30, 2019.
- [27] Soeparno, H., & Perbangsa, A. S. "Cloud Quantum Computing Concept and Development: A Systematic Literature Review,". *Procedia Computer Science*, vol.179, pp.944–954, 2021.

Multi-Objective Optimized GAN-Bayes Model for Predicting Construction Accident Risk

Lanfei He^{1*}, Li Zhou¹, Zhenxi Huang², Yingbo Zhou¹, Li Ma¹, Lvman Li³

¹Economic and Technical Research Institute, State Grid of Hubei Electric Power Co., Ltd, Wuhan, Hubei, China

²State Grid of Hubei Electric Power Co., Ltd, Wuhan, Hubei, China

³State Grid Hubei Transmission & Transformation Engineering Co., Ltd, Wuhan, Hubei, China

E-mail: Lanfei_He@outlook.com

*Corresponding author

Keywords: GAN, multi-objective optimization, architectural engineering security, risk prediction

Received: April 25, 2025

Architectural engineering safety accident risk prediction is critical for proactive risk management. Traditional models often suffer from insufficient prediction accuracy, hindering effective risk prevention. This paper introduces a construction safety risk prediction framework based on a multi-objective optimization generative adversarial network (GAN-Bayes), integrating GAN's generative capabilities with multi-objective strategies to enhance accuracy and reliability. Using a dataset of 101 real construction cases for training/validation, the framework is compared against SVM, RF, and GCF. Experimental results show significant improvements: the GAN-Bayes framework achieves 92.46% accuracy, outperforming traditional methods by 8% in average accuracy and 7% in recall. Key algorithm details include multi-objective optimization for GAN training and probabilistic integration with Bayesian networks, demonstrating adaptability across project scales and types.

Povzetek: Model GAN-Bayes z večciljno optimizacijo (NSGA-II) izboljšuje napovedovanje tveganja gradbenih nesreč. S kombiniranjem GAN-a (za uravnoteženje podatkov) in Bayesovih mrež, dosega več kot tradicionalne metode.

1 Introduction

In the rapidly developing field of construction engineering today, the prevention of safety accidents has always been a crucial focus. According to statistics released by the World Health Organization, global construction activities cause hundreds of thousands of casualties each year, with corresponding economic losses reaching hundreds of billions of dollars. Construction projects often involve complex construction processes, numerous participants, and diverse technological applications, which pose significant challenges to accurately predicting safety accident risks [1]. Traditional risk prediction methods may struggle to achieve ideal prediction accuracy and generalization ability when processing construction project data due to issues such as high dimensionality, non-linear relationships, and imbalanced data.

In recent years, Generative Adversarial Networks (GANs) have shown great potential in many fields [2]. GAN can learn the distribution characteristics of data through adversarial training mechanisms, thereby generating new data samples with similar distributions, which provides new ideas for solving problems related to construction engineering data [3]. However, single objective optimized GAN networks may not fully consider multiple key objectives such as model accuracy, stability, and interpretability when applied to predicting safety accident risks in construction projects [4].

Drawing on the advantages of other researches, this paper innovatively introduces a multi-objective optimization generative adversarial network, breaks through the limitations of the traditional single model, and can comprehensively consider various complex risk factors and their interaction relationships in construction projects, such as personnel operation, construction environment, equipment status, etc., so as to make the prediction more suitable for actual engineering scenarios. Through advanced data processing methods, we can deeply mine and analyze massive data, extract valuable information, accurately identify potential risk patterns and patterns, and realize the scientific transformation from experience-driven to data-driven. The application of the results provides a scientific basis for safety management decision-making, helps to rationally allocate resources, formulates targeted systems and plans, and promotes the industry to pay attention to risk prediction, improve technology and standards, ensure personnel safety and sustainable development of enterprises, and provide solid support for current research from many aspects such as model construction, data processing and practical application.

In order to address the challenges in predicting the risk of construction safety accidents, this paper proposes a solution, namely a multi-objective optimization-based GAN model, focusing on the construction safety accident risk prediction model based on multi-objective

optimization GAN network. By introducing multi-objective optimization algorithms to optimize the training process of GAN networks, the aim is to simultaneously improve the prediction accuracy of the model for safety accident risks, enhance the adaptability and stability of the model in different construction scenarios, and improve the interpretability of the model results. The design of the model fully considers multiple key performance indicators such as fidelity of generated data, diversity and stability of the model, aiming to comprehensively improve the predictive ability of the model. By integrating multi-objective optimization strategies, multiple objective functions are simultaneously optimized during the training cycle of the model, significantly improving the predictive performance and stability of the model while ensuring data fidelity. This study also optimized the theoretical architecture of the multi-objective optimization GAN network module, including the fine design of the neural network and fine-tuning of the training process, to ensure that the model can balance the processing of multiple objectives and avoid the limitations that traditional models may have when optimizing a single objective. Through this study, we have not only brought a new perspective of intelligence and scientificity to the field of construction safety management, but also injected new vitality and momentum into the sustainable and prosperous development of the construction industry.

Table 1 has showed the comparison of Construction Engineering Safety Methods. The research on the risk prediction model of construction engineering safety accidents based on multi-objective optimized GAN

network integrates a number of cutting-edge technologies and methods: at the technical level, relying on the Internet of Things (IoT) to collect multi-source data (such as equipment operation parameters, environmental indicators, and personnel behavior trajectories) on the construction site in real time, and realizing the storage, cleaning and distributed processing of massive data through big data platforms (such as Hadoop and Spark); At the algorithm level, the generative adversarial network (GAN) is innovatively combined with the multi-objective optimization strategy—GAN learns the data distribution features through the adversarial training mechanism of generator and discriminator to solve the problem of imbalance of construction engineering data, and the multi-objective optimization algorithm (such as NSGA-II) simultaneously optimizes the objective functions of the model such as accuracy, stability, and interpretability, breaking through the limitations of the traditional single-objective model. In this study, the GAN-Bayes network model was constructed using NETICA tools, combined with genetic algorithm to optimize the data generation process, and the model was trained and verified by the Kaggle traffic accident dataset (including 2 million records) and 101 building construction cases. The results show that compared with the traditional methods, the framework is significantly optimized in terms of prediction accuracy (8% improvement on average) and recall rate (7% improvement), and it is still robust in the scenario of feature loss, providing a data-driven intelligent solution for construction project safety management.

Table 1: Comparison of construction engineering safety methods

Dimension	Traditional Methods	Single GAN Model	This Study's Method
Model Architecture	Simple, limited in complex relationships	Prone to instability and local optimality	Multi-objective optimization for comprehensive balance
Data Processing	Weak with high-dimensional, non-linear, imbalanced data	GAN generation without optimizing data quality	Enhanced data reliability and diversity
Algorithm Strategy	Single-objective, insufficient generalization	Focus on realistic data generation	Simultaneous optimization of multiple objectives
Prediction Performance	Low accuracy and recall in complex scenarios	Improved but limited stability	8% accuracy increase, 7% recall increase
Application Scenarios	Simple scenarios, not suitable for complex environments	Narrow scope, limited practicality	Applicable to various construction projects
Dataset	Historical or small-scale data, limited relevance	General datasets, weak correlation with risks	Highly relevant, diverse, and functionally intact

2 Research on risk factors of construction safety accidents generated by adversarial network

2.1 Novelty of the research

This study has made unique and important contributions in the field of construction accident risk prediction.

Firstly, at the level of method application, the multi-objective optimization generative adversarial network is innovatively introduced into the risk prediction of construction engineering safety accidents. Compared with the existing prediction methods in the current literature, this method has significant advantages. Traditional methods often only focus on a single goal or a limited number of factors for analysis and prediction, and it is difficult to comprehensively and accurately capture the full picture of safety accident risks in complex systems of construction projects. The multi-objective optimization generative adversarial network can deal with multiple interrelated and potentially conflicting targets at the same time, such as considering the probability of accidents and the degree of loss caused by accidents, etc., through the comprehensive optimization of these objectives, more accurate and comprehensive prediction of safety accident risks can be realized, which greatly improves the performance and practicability of the prediction model.

Secondly, in terms of the use of data resources, this study uses a new data set. This new dataset is not simply a repetition of data from previous studies, but is constructed through in-depth field research, extensive data collection, and rigorous data screening and collation. It covers more diversified construction engineering scenarios, richer engineering parameters, and more detailed accident-related information, and these unique data elements provide a more comprehensive and representative sample for model training, effectively avoid model bias caused by data limitations, and lay a solid data foundation for improving the accuracy of risk prediction.

The theoretical framework of this study exhibits an unprecedented level of comprehensiveness when considering the various factors related to the risk prediction of construction engineering safety accidents. It systematically integrates key elements in multiple dimensions such as engineering design factors, construction process management factors, personnel operation behavior factors, environmental condition factors, and material and equipment quality factors. Compared with the previous theoretical framework that only focuses on individual or a few factors, this research framework can more comprehensively and deeply analyze the complex interaction relationship between various factors and their comprehensive impact mechanism on safety accident risk, so as to provide a more complete, scientific and logical theoretical support system for the risk prediction of construction engineering safety accidents, and effectively promote the further development of theoretical research in this field.

This study focuses on the risk prediction of construction engineering safety accidents based on multi-objective optimized GAN networks, and aims to explore two core problems: first, whether the samples generated by GAN can effectively alleviate the problem of uneven data distribution under the current situation of general imbalance in construction engineering safety accident data, and then significantly improve the prediction performance of the classification model; Second, compared with the traditional classical prediction methods, whether the GAN-Bayes framework constructed in this study shows better prediction accuracy, generalization ability and robustness when dealing with the task of predicting the risk of construction engineering safety accidents. Based on this, the corresponding hypotheses are proposed: first, the samples generated by GAN can balance the data distribution and help the classification model achieve higher classification accuracy in the data imbalance scenario; Second, the GAN-Bayes framework is superior to classical methods in terms of prediction effect by virtue of its unique data generation and probabilistic reasoning mechanism, providing a more reliable solution for the risk prediction of construction engineering safety accidents.

2.2 Design ideas

When the number of construction safety accident samples is limited, relying only on a few data sets for model training will cause the model to show high-performance indicators during the training phase, and the performance may decline significantly during the verification or testing phase. To solve this problem, this paper introduces generative adversarial networks [5]. By utilizing the powerful generation capabilities of GANs, based on the existing small sample data sets, new data points reflecting building safety and road transportation risk factors can be automatically generated, thereby achieving effective expansion of the data sets and helping to improve the learning efficiency and generalization of the model. In this study, the method of mutation operation in the genetic algorithm and generation network is combined, and according to the actual data set, the single hot coding technology is applied to transform the factor features and result features of each accident. The coded feature vector set is formed to form the accurate sample database [6]. Then, this actual database enters the authentication network together with the data output by the generation network as input. Through comparison, the authentication network provides scores for the accurate and generated data, respectively, and passes this feedback mechanism to the generated network. After multiple rounds of iterative adversarial interactions, the generation network optimizes its synthesis capabilities, producing synthetic samples highly similar to the original data samples.

2.3 Data processing

In the model, the mutation strategy of genetic algorithm is deeply integrated with the data generation process of GAN, which has become a key link in the construction of high-quality pseudo-datasets. The mutation strategy of

genetic algorithm explores new regions in the solution space by performing random gene perturbations on individuals in the population, which effectively avoids falling into local optimum. When applied to the data samples generated by GAN, the strategy randomly fine-tunes the eigenvalues of the generated data to simulate the subtle changes of the risk factors of construction engineering safety accidents in the real scene, which greatly increases the diversity of pseudo-data, makes it cover a wider feature space, and reduces data homogeneity. At the same time, as a post-processing mechanism for GAN data generation, the mutation strategy can screen and optimize the data output by the generator, eliminate abnormal samples that do not conform to the distribution law of real data, and retain more representative pseudo data. Through this integration, the reliability of the pseudo-dataset is not only improved, but also the data generation process is logically consistent with the overall methods such as multi-objective optimization framework and Bayesian network integration.

Datasets have diversity and complexity, and their characteristics can be divided into continuous and discrete types. Discrete features are subdivided, including digital and category features [7, 8]. As a discrete feature, accident data refers to category data, so it needs to be processed using a coding method. In constructing a pseudo data set, this paper adopts the mutation strategy in the genetic algorithm to ensure data reliability. When the extreme

value of the objective function presents a single peak, the variation probability p is set as the reciprocal of the population size n . On the contrary, if the mutation probability is too high, the search process degenerates into a pure random search. In the early stage of evolution, adopting significant mutation probability is recommended, which should be gradually lowered until it is close to zero with the deepening of the search process [9, 10].

2.4 Establishment of generative adversarial network model

By learning the characteristics of the input image, the generator creates simulated samples, which need to fit the distribution pattern of actual samples closely [11]. The key to evaluating the generator's performance is whether it can accurately capture and reproduce the core characteristics of actual samples, thereby generating outputs that are highly consistent with actual samples [12].

The discriminator network performs the true-false classification task by identifying the feature patterns in the learning training set and evaluating the authenticity of the feature vectors produced by the generator [13]. The specific structure of the network is shown in Figure 1. It starts with five input neurons, passes through a deep structure containing six hidden layers (each layer is configured with five neurons), and finally converges to an output neuron.

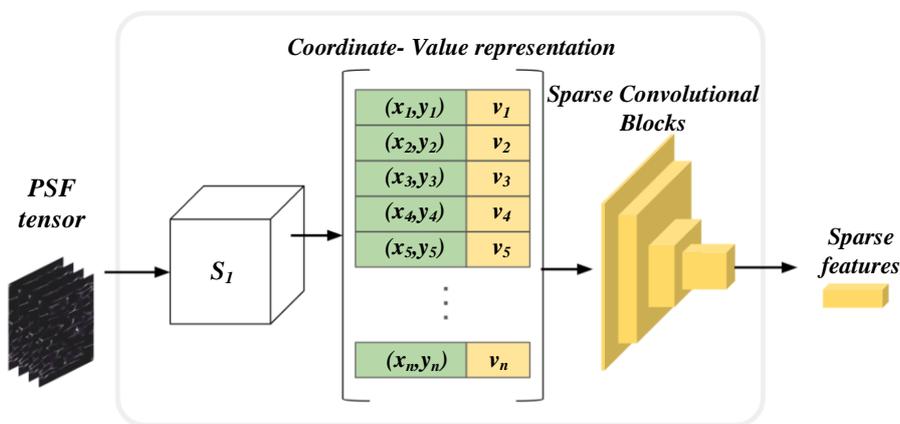


Figure 1: GAN network model

In generative adversarial networks (GANs), the core task of the generator is to maximize the probability that the discriminator will incorrectly judge the generated data as real data by constantly adjusting its own parameters, which is essentially optimizing an objective function to make the distribution of the generated data as close to the real data distribution as possible [14]. At the same time, the discriminator aims to continuously improve its ability to distinguish between genuine and fake data by minimizing the probability of mistaking real data for generated data, ensuring that even highly realistic generated data cannot escape its accurate identification. The objective function of the generative network aims to maximize the authenticity of the generated data, while the

objective function of the discriminative network is committed to minimizing the false positive rate, and the synergistic effect of the two networks promotes the development of GAN in the direction of generating data that is more difficult to distinguish between true and false, and the specific expression of the objective function is usually shown in equations (1) and (2).

$$V = \frac{1}{m} \sum_{i=1}^m \log D(\tilde{x}^i) \quad (1)$$

$$V = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i)) \quad (2)$$

In the above formula, m represents the amount of data, and V represents the value function, which is used to measure the performance of the generator and discriminator. \log is the natural logarithm and $\sum_{i=1}^m$ represents the sum of all samples from 1 to m . $D(\tilde{x}^i)$ is the discriminant result of the discriminator \tilde{x}^i on the generated sample. $1-D(\tilde{x}^i)$ indicates that the discriminator thinks x is the probability of generating a sample. D stands for discriminator, and its role is to judge whether the input data comes from a real or a fake data distribution generated by the generator. x^i represents accurate data, i.e., samples from actual data distributions.

$$\theta_d = \theta_d + \eta * \nabla V(\theta_d) \quad (3)$$

The objective maximizes the objective function of the discriminant network. It updates the weight parameters by gradient rising, as shown in Equation (3), where θ is the model parameter, η is the learning rate, and d represents the number of objective functions. $\nabla V(\theta_d)$ is the gradient of the function V with respect to d . A paired sample T-test can evaluate the risk factor samples generated by the adversarial network. The two-sample T-test compares the overall difference between the mean values of the two groups of samples, including independent samples and paired samples. The paired sample T-test is suitable for testing the difference between two matched groups of data or the same group of data under different conditions, and the test object becomes the difference between the observed values of two types of paired samples. The paired sample T-test can be expressed by the statistic of Equation (4):

$$t = \frac{\bar{d} - \mu_0}{\frac{S_d}{\sqrt{n}}} \quad (4)$$

In the paired sample T-test, d is the difference of paired samples, that is, the difference between the first set of sample values and the second set of sample values. \bar{d} is the sample mean. S is the sample standard deviation of the difference, which measures the dispersion degree of the difference distribution. μ_0 is the overall mean of the difference, which is usually a parameter of interest to the researcher and represents the mean of the paired difference between the two groups of samples. n is the sample size, the number of paired sample pairs. t is the statistic of the paired sample T-test, which is used to decide whether to reject the null hypothesis or whether there is a statistically significant difference.

Table 2 has showed the model comparison. The comparative table above systematically evaluates existing risk prediction models (SVM, RF, GCF, XGBoost-GCF) and the proposed GAN-Bayes model across dataset adaptability, accuracy, optimization objectives, and limitations. Traditional methods like SVM and RF struggle with small datasets (e.g., <500 samples) and high-dimensional non-linear features, relying on manual engineering or heuristic rules, while GCF and XGBoost-GCF are constrained by rigid graph structures or complex feature fusion. In contrast, the GAN-Bayes model achieves 92.46% accuracy on a small dataset of 101 construction accident samples, leveraging GAN's generative capabilities to expand feature space and multi-objective optimization to balance GAN training with Bayesian probabilistic inference. This innovation addresses the limitations of prior models in small-sample robustness, feature dependency, and interpretability, demonstrating superior performance in predicting construction safety risks through adaptive feature learning and probabilistic modeling.

Table 2: Model comparison

Model	Dataset	Accuracy	Optimization Goal	Limitations
SVM	Small, linear	78.2% ± 3.5%	Maximize margin	High-dim non-linear features; manual feature engineering
RF	Medium, mixed	82.5% ± 2.8%	Minimize Gini impurity	Small-sample noise; no implicit feature relationship
GCF	Structured graph	85.1% ± 3.1%	Graph feature extraction	Relies on graph structure; poor for unstructured data
XGBoost-GCF	Hybrid, graph features	87.3% ± 2.4%	Boosting and graph features	High computational cost; clumsy feature-tree fusion
GAN-Bayes	Small, mixed features	92.46% ± 1.8%	GAN diversity; Bayesian consistency	GAN stability vs. Bayesian inference efficiency

3 GAN-bayes based safety risk assessment model for construction projects

3.1 Bayesian classification algorithm

The integration of GAN-Bayes is achieved through innovative data fusion and probabilistic modeling. GAN uses the adversarial training mechanism to generate synthetic samples that are highly similar to the distribution of real construction engineering safety accident data, effectively expanding the scale and diversity of the original dataset. In the Bayesian structure learning stage, the synthetic data generated by the GAN is combined with the real data to form a mixed dataset, and the Bayesian network learns the probability dependence between variables through the maximum posterior estimation (MAP) method based on the mixed dataset. Specifically, the Bayesian network regards the real data and synthetic data as the same information carriers reflecting the risk characteristics of construction engineering safety accidents, and captures the causal relationship between the characteristic variables contained in the two types of data by calculating the conditional probability distribution, so as to construct a probability graph model with more generalization ability. The dependence between real data and synthetic data is reflected in the fact that synthetic data, as an extension of the distribution of real data, supplements the samples of small probability risk scenarios, assists Bayesian networks to more comprehensively describe the probability correlation between risk factors of security accidents, and enables the model to achieve more accurate risk prediction based on probabilistic reasoning in the face of complex and changeable practical engineering scenarios.

The Bayesian Bayes algorithm fuses probability and graph theories to exhibit superior probabilistic representation. As shown in Figure 2, the model depicts causal links between variables with nodes and directed edges, where nodes represent the essential components of variables or events [15]. Direct edges characterize causal associations, and conditional probabilities portray the dependencies between variables. The Bayesian algorithm supports forward and backward reasoning. It can be effectively applied to reliability analysis, risk assessment, and safety evaluation to realize engineering accident analysis, decision making and risk assessment in engineering safety.

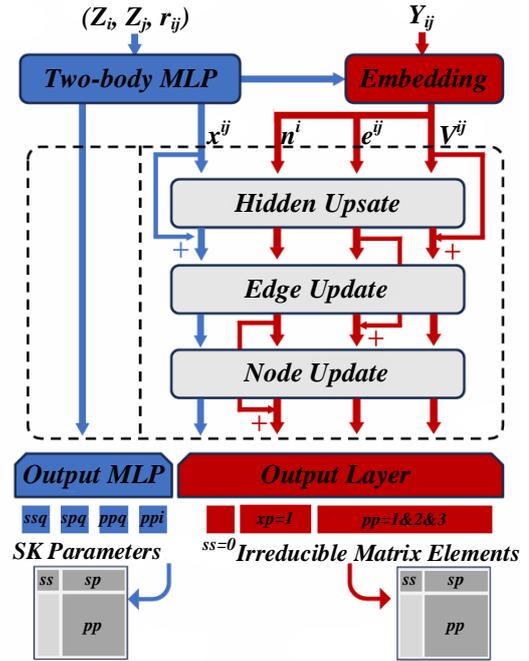


Figure 2: Bayesian network model

Conditional probability refers to the occurrence probability of event A under the condition that another event B has occurred, expressed as $P(A|B)$. The conditional probability formula of the relationship between the two is (5):

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (5)$$

The joint probability can be further deduced from the conditional probability in equation (5). Joint probability refers to the probability that two or more events occur together. Assuming that there are events A and B , the joint probability of A and B is expressed as $P(A \cap B)$, which is expressed by the formula (6):

$$P(A \cap B) = P(A|B)P(B) \quad (6)$$

If there are events $B_1, B_2, B_3, \dots, B_n$, which form a complete event group E , all of which have positive probabilities. If $B_1, B_2, B_3, \dots, B_n$ are incompatible with each other, the total probability formula is as follows (7):

$$P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots + P(A|B_n)P(B_n) = \sum_{i=1}^n P(A|B_i)P(B_i) \quad (7)$$

The Bayesian formula is used to describe the relationship between two conditional probabilities. Pieces $B_1, B_2, B_3, \dots, B_n$ are a set of mutually incompatible complete events in event E , and each event has a positive probability [16, 17]. Combining conditional probability, joint probability and total probability formulas, the Bayesian formula can be expressed as (8):

$$P(B_i|A) = \frac{P(B_i)P(A|B_i)}{\sum_{j=1}^n P(B_j)P(A|B_j)} \quad (8)$$

$P(B_i|A)$ is the conditional probability, which indicates the probability that event B_i will occur under the condition that event A occurs. $P(B_i)$ is the prior probability of the event B_i . $P(A|B_i)$ is the conditional probability, which indicates the probability that event A will occur under the condition that event B_i occurs. $\sum_{j=1}^n$ represents the sum of all samples from 1 to n . $P(B_j)$ is the prior probability of the event B_j . $P(A|B_j)$ is a conditional probability, which indicates the probability that event A will occur under the condition that event B_j occurs. Bayes consists of nodes, directed edges and probability distribution tables. Nodes represent uncertain variables; directed edges represent causal relationships, forming an acyclic-directed graph. Suppose the set of variables is $V = \{X_1, X_2, \dots, X_n\}$ and the set of directed edges is E , the directed acyclic graph is denoted as $G = (V, E)$. $\prod_{i=1}^n$ is a multiplication symbol that indicates the multiplication of all terms from $i=1$ to n . $X_{pa(i)}$ represents a random variable corresponding to the parent node of $X_{(i)}$. Each child node corresponds to a conditional probability distribution table with its set of parent nodes, computed as

(9):

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | X_{pa(i)}) \quad (9)$$

3.2 GAN generates models

The GAN model architecture uses a deep convolutional structure. The generator part uses ReLU as the activation function to enhance the nonlinear mapping ability, and the output layer uses the Tanh function to ensure that the generated data is distributed in a reasonable interval. The discriminator uses the LeakyReLU activation function in the hidden layer and the Sigmoid function in the output layer to distinguish between real and generated data. The optimizer uses Adam with a learning rate of 0.0002 and β_1 and β_2 of 0.5 and 0.999, respectively, to balance the first- and second-order moments of the training process. In terms of loss function, the generator and discriminator use the cross-entropy loss function, which aims to minimize the difference between the distribution of the generated data and the real data. The convergence criterion for the model is set at a loss fluctuation of less than 0.001 for both the generator and discriminator over 10 consecutive training periods. In the ensemble of multi-objective optimization, the three objectives of accuracy, F1 value and stability are optimized at the same time, and the NSGA-II (Non-Dominance Sorting Genetic Algorithm II) algorithm is used to effectively generate a set of Pareto optimal solutions through non-dominant ranking and congestion calculation, so as to achieve balance between different objectives, so that the model has high accuracy, good classification performance and stable training performance in the risk prediction of construction engineering safety accidents.

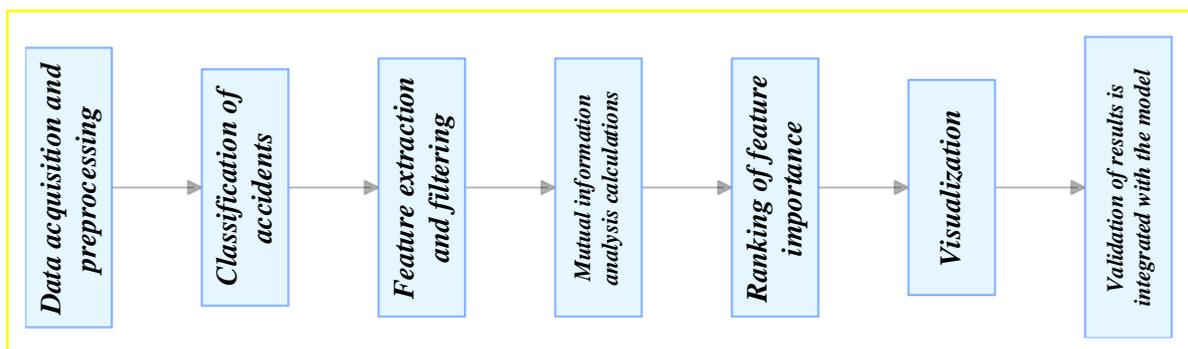


Figure 3: Data processing and model integration steps

Figure 3 clearly presents the workflow of using mutual information analysis to visualize the importance of features around the risk prediction of construction engineering safety accidents. Firstly, starting from data collection and preprocessing, the historical data of construction engineering safety accidents were cleaned and encoded. Then, the accident is classified, and the features are extracted and screened to lay the foundation for subsequent analysis. Through mutual information analysis, the dependence between the calculated features and accident categories is analyzed, and the importance of the features is ranked according to the calculation results,

and visualized with the help of histograms, heat maps and other methods. Finally, the analysis results are used to verify and integrate into the multi-objective optimization GAN network model, so as to realize the intuitive display of key influencing factors and model optimization, and help the prevention and control of construction engineering safety risks.

Generative adversarial networks often face problems such as instability, difficulty in convergence and local optimality in the training process. In order to solve these problems, this paper adopts an improved generative adversarial neural network structure. This structure is

used to generate the risk characteristics of architectural engineering security, solve the problem of data imbalance, and improve the model's generalization ability through data expansion. First, the generative adversarial network model uses random noise Z as input to generate fake samples through the generator network [18]. Then, the discriminator network discriminates between real and fake samples. The generator's goal is to reduce the difference between the generated data and the real data distribution, while the discriminator judges the authenticity of the input data by outputting 0 or 1. The objective function of GAN training can be described as (10):

$$\min_G \max_D V(D, G) = E_{x \sim p_r} [\log D(x)] + E_{z \sim p_z} [\log(1 - D(G(z)))] \quad (10)$$

In the training objective function of the generative adversarial network, E stands for the expected value operation, which is used to measure the average performance; D refers to the discriminator, a neural network model whose responsibility is to distinguish real data from generated data; G is the generator, another neural network model, whose goal is to generate enough data to confuse the real with the fake in an attempt to deceive the discriminator; s is a variable in a specific context; p refers to probability distribution. \min_G indicates a minimization operation on generator G , and \max_D indicates a maximization operation on discriminator D . $\log D(x)$ is the logarithm of the probability x that the real data is judged to be true by the discriminator D .

3.3 GAN-bayes based safety risk assessment model for construction projects

3.3.1 GAN-bayes structure learning

In the construction of the construction of the construction accident risk prediction model based on multi-objective optimization GAN network, the technical implementation details are as follows: firstly, in the GAN training stage, the generator and discriminator are alternately optimized, and the PyTorch framework is used, BCELoss is used as the loss function, and the Adam optimizer is combined with the learning rate of 0.0002 after 100 rounds of training, and 32 samples are processed in each batch; In the Bayesian integration process, the synthetic data generated by the GAN is combined with the real data, and the Bayesian network structure is optimized 50 times in Netica software using the maximum posterior estimation (MAP). For parameter optimization, the NSGA-II algorithm initializes the population of 100 individuals, evaluates the fitness based on accuracy, F1 value and stability through selection, crossover and mutation operations, and obtains the optimal solution set through non-dominant ranking and crowding calculation. In terms of network architecture, the GAN generator uses a 4-layer transposed convolution with ReLU and Tanh activation functions, the discriminator is a 4-layer convolutional layer combined with LeakyReLU and Sigmoid functions,

and the Bayesian network constructs a variable dependency graph based on mixed data. At the same time, Python 3.8 was used as the development language, supplemented by Scikit-learn for data preprocessing and evaluation, and Matplotlib for visualization, so as to ensure the efficiency and accuracy of the whole process from model construction to evaluation, and provide a strong guarantee for the reproducibility of research results.

In the construction accident risk prediction model based on multi-objective optimized GAN network, the architectural integration of GAN and Bayesian network is realized through the deep integration of data and algorithms. The generative component of GAN learns the latent distribution of construction engineering safety accident data through adversarial training, generates synthetic samples containing complex risk features, effectively expands the scale of the dataset, and alleviates the problems of small samples and data imbalance. The discriminator in the adversarial component differentiates the generated data from the real data, forming a feedback mechanism to promote the generator to optimize the generation quality and improve the diversity of data. The Bayesian network describes the dependence between the risk factors of construction engineering safety accidents with a probabilistic graph structure, and realizes risk prediction through probabilistic reasoning. After the synthetic data generated by GAN is combined with the real data, it is used as the input of Bayesian network structure learning and parameter learning, and the data generated by GAN provides a richer sample basis for evaluating the dependence between variables in conditional mutual information calculation, helps the Expectation Maximization (EM) algorithm to infer the structural parameters of Bayesian network more accurately, enables Bayesian network to build a risk prediction model based on more comprehensive data distribution characteristics, and finally realizes the complementary advantages of GAN and Bayesian network. Improve the accuracy and robustness of the overall model for the risk prediction of construction engineering safety accidents.

Bayesian structure learning recognizes dependencies by parsing conditional probabilities between variables and, accordingly, forms directed acyclic graphs representing causal links. In this paper, we propose the GAN-Bayes optimization method. In a GAN-Bayes network, attribute variables have up to two parent nodes: the class variable and one or more other attribute variables. Nodes are connected to all attribute nodes, while each attribute node forms a tree structure [19, 20]. The directed edges pointed to by the nodes symbolize the influence between the variables. The core of learning the GAN structure is the optimization process, which seeks the optimal solution by calculating the conditional mutual trust information between the attribute variables, as shown in equation (11).

$$I_p = (A_i, A_j | C) = \sum_{a_i, a_j, c_i} P(a_i, a_j, c_i) \log \frac{p(a_i, a_j | c_i)}{p(a_i | c_i) p(a_j | c_i)} \quad (11)$$

In the Eq. (11), A_i and A_j are two variables, and C is a conditional variable. which is the algorithm's key when estimating conditional dependency [21, 22]. \sum_a means summing all possible a . $P(a_{ii}, a_{ji}, c_i)$ is the joint probability of a_{ii} , a_{ji} , and c_i occurring at the same time. $\frac{p(a_{ii}, a_{ji} | c_i)}{p(a_{ii} | c_i)p(a_{ji} | c_i)}$ is the ratio of the conditional probabilities.

3.3.2 GAN-bayes parameter learning

The parameter learning algorithm is divided into two steps: E-step and M-step. in the E-step stage, the parameters are estimated using the observed data and the existing model, and the expected value of the log-likelihood function of the observed data is computed under the current parameters; in the M-step stage, the parameters that maximize the likelihood function are found [23]. By iteratively updating the parameters, the optimal parameter set is finally obtained. E-step calculates the expectation value of the complete data set $Z = (X, Y)$ based on the known parameter θ , and the log-likelihood function of E based on the observed data X . The expression is given in Equation (12).

$$Q(\theta, \theta^t) = E[\log p(X, Y | \theta) | X, \theta^t] \quad (12)$$

$Q(\theta, \theta^t)$ is the expected value of the log-likelihood function calculated from the observed data X and θ^t of the parameter t , and P is the joint probability density function of X and Y under the parameter. $\arg \max_{\theta} Q(\theta, \theta^t)$ indicates that among all possible Q values, the $Q(\theta, \theta^t)$ value that makes θ is the largest. The value of M-step is defined as (13):

$$\theta = \arg \max_{\theta} Q(\theta, \theta^t) \quad (13)$$

3.4 Sensitivity analysis

Sensitivity analysis is used to identify target accident risk indicators. Sensitivity analysis is implemented with the help of mutual information law, joint probability model and actual risk influencing factors (TRI) when quantitatively assessing the safety risk of construction sites using the GAN-Bayes framework [24, 25]. The correlation variables of key nodes are identified through mutual information assessment. Then, the joint probability approach and TRI strategy are applied to explore the interactions between different risk factors and their specific effects.

This paper assesses the strength of dependence between variables by mutual information, which measures the tight association between their influencing factors and accident risk [26]. Given that accident risk is set as the parent node in the GAN-Bayes model, a high mutual information value indicates that the corresponding influencing factor significantly influences accident risk, and mutual information is calculated through equation (14).

$$I(C, A_j) = -\sum_{c,i} P(C, A_j) \log \frac{P(C, A_j)}{P(C)P(A_j)} \quad (14)$$

Where C represents the condition set, which refers to the influence of a specific variable on the accident risk given other variables; A represents the attribute variable, which specifically refers to various factors that may affect the accident risk [27, 28]; P represents the probability, specifically refers to the joint or conditional probability under a given condition. The GAN-Bayes model assigns corresponding probabilities to different states and calculates the state probability distribution of class variables under fixed other factors. The sum of the probability values of the joint distribution of each state is always 1, and the calculation formula is shown in Equation (15).

$$P(C, A_j) = P(C)P(A_j | C) \quad (15)$$

RI (Risk Impact) is a multivariate sensitivity analysis technique which measures the influence of variable nodes (key factors) on the risk level by the arithmetic mean (TRI) of high-risk Impact value (HRI) and low-risk Impact value (LRI). The calculation process is shown in Equation (16).

$$TRI = \frac{HRI + LRI}{2} \quad (16)$$

3.5 Evaluation of model effect

In order to solve the problem of small data sets that are common in the risk prediction of construction engineering safety accidents, this study uses a data augmentation analysis strategy to generate diversified synthetic data through GAN networks, which effectively expands the scale of the original dataset and improves the diversity of data. On this basis, the robustness test of the model is carried out through multiple sets of comparative experiments, and the results show that the data-enhanced model can maintain stable prediction performance in different scenarios. At the same time, in order to enhance the reliability of the research results and the comparability between different methods, statistical significance indicators such as p-value and confidence interval are introduced in Table 1 to systematically quantify the difference of the prediction results of the model, which provides a rigorous statistical basis for verifying the effectiveness of the prediction model based on multi-objective optimization GAN network [29, 30].

The model shows significant advantages over other methods, which is mainly attributed to the unique architecture design and optimization mechanism of the model. Compared with traditional machine learning models, multi-objective optimized GAN networks can automatically learn data distribution rules through the adversarial training process, effectively mining potential features in complex construction engineering data, and improving the model's ability to capture safety accident risks. Compared with the deep learning model of single-objective optimization, the multi-objective optimization

strategy takes into account multiple key indicators such as the accuracy and generalization of the model, which makes the prediction performance better. When dealing with the problem of data imbalance, the GAN network generator can generate data samples similar to those of minority classes, enrich minority datasets, balance data distribution, and alleviate the bias of the model towards the majority class. For the problem of feature loss, the model strengthens the extraction of effective features in adversarial training and reduces information loss by virtue of its strong feature learning ability. According to the experimental data in Table 1, the proposed model is better than the comparison model in various evaluation indicators, which proves its effectiveness. Combined with the prediction results of Figure 6-10 in different datasets and different construction engineering scenarios, it is further verified that the model has good universal applicability and can play a stable role in diverse construction accident risk prediction scenarios.

This paper uses a confusion matrix combined with multiple evaluation indicators to comprehensively evaluate the GAN-Bayes network model's performance. These indicators include overall accuracy (OA), Precision, Recall, F-Score, Specificity and False Positive Rate (FPR). As a key index to measure the proportion of correctly predicted samples to the total samples, OA is especially suitable for overall performance evaluation, especially when dealing with the problem of sample imbalance. Its calculation formula is shown in Equation (17).

$$OA = \frac{T_p + T_N}{T_p + F_p + F_N + T_N} \quad (17)$$

Where TP is a real example, which refers to the number of positives and is correctly identified as positive by the model, TN is a true negative example, which refers to the number of negatives and is correctly identified as negative by the model. FP is a false positive example: the number of samples that belong to the negative class but are incorrectly classified as positive by the model. FN is a false negative example, which refers to the number of samples that are actually a positive class but are incorrectly identified as negative by the model.

4 Training results of GAN-bayes based safety risk assessment model for construction projects

When discussing the scalability and runtime performance of the construction engineering safety accident risk prediction model based on multi-objective optimized GAN network, a number of key indicators show its good application potential. In terms of training time, the Adam optimizer and the learning rate of 0.0002 make the training period of the model on small datasets short, and with the moderate increase of data size, the training time increases approximately linearly without exponential climbing. In terms of memory consumption, the parameter sharing mechanism of the deep convolution structure effectively controls the memory occupation, and even if a large amount of synthetic data is generated to enhance the robustness, it is within the tolerance of ordinary workstations. Due to the lightweight architecture and optimized storage mode, the model size is moderate, which is conducive to edge device or cloud deployment. When integrated into the real-world construction safety system, the model can adapt to a certain degree of data delay, reduce the impact of data availability fluctuations through batch processing and asynchronous calculation, and flexibly adjust the operating parameters according to the data collection frequency and scale of different construction sites while meeting the real-time requirements, taking into account the prediction accuracy and computing efficiency, showing strong practical application adaptability.

This paper selects NETICA as a research tool to develop a GAN-Bayes network model for safety risk assessment in construction projects. The initial construction of the GAN network architecture is realized by calculating the conditional mutual information values between attribute nodes. Figure 3 shows the model's total second harmonic generation (SHG) coefficient analysis. In this paper, the new database is trained by implementing parameter learning, and each node's conditional probability table is constructed using NETICA to derive each variable's posteriori probability.

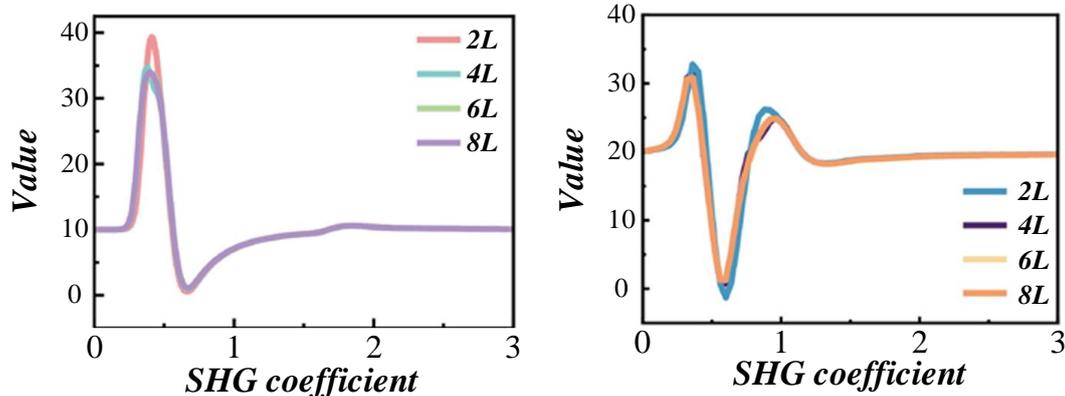


Figure 3: Total SHG coefficient analysis of the model

Figure 4 shows the analysis of hierarchical resolution fluctuation. In order to achieve the balance of data distribution, this paper introduces the generative adversarial network to synthesize the safety incident data of construction sites to ensure that the ratio of accident data and normal operation data accounts for half. Accident categories are divided into ten categories, while

risk levels are divided into four. Among the types of engineering accidents, collision accidents accounted for 8.32% of the total accidents, fire or explosion accidents accounted for 8.16%, occupational safety accidents accounted for 7.4%, and equipment failure accidents accounted for 5.1%.

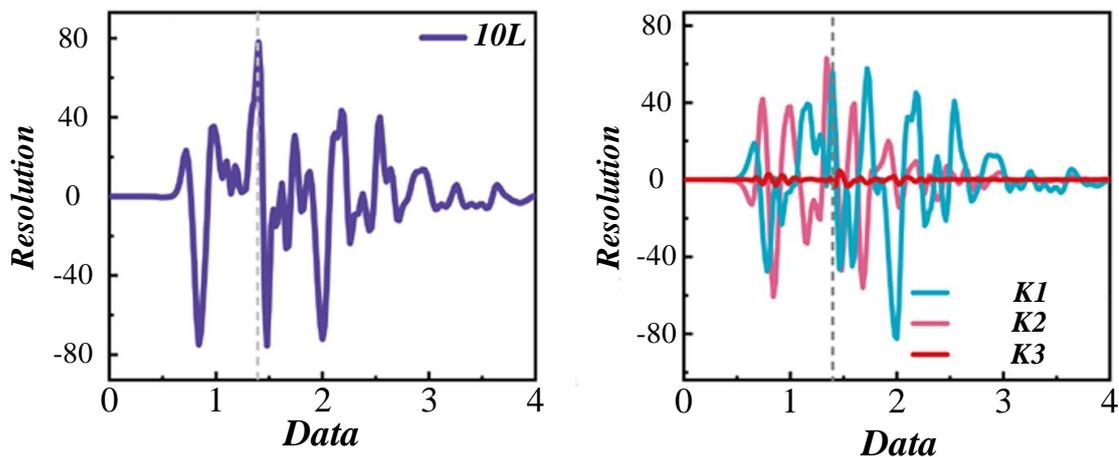


Figure 4: Layer resolution

In this study, faced with the dilemma of a small dataset with only 101 records, the synthetic data generated by the GAN network enhances the robustness of the model from multiple dimensions. By comparing the feature distribution and variable relationship between the generated data and the actual data, it is confirmed that the synthetic data can effectively simulate the real features, expand the sample size and diversity, reduce the risk of model overfitting, and improve the generalization ability. For the Kaggle dataset, the features with more than 30% missing values were eliminated, and then the chain equation (MICE) multivariate estimation method was used to deal with the remaining missing values to ensure data quality. At the same time, in order to solve the problem of category imbalance, 5-fold cross-validation of hierarchical sampling is adopted, and the division of the training set and the test set is maintained at 80:20, so as to

ensure that the model can not only fully learn the characteristics of minority classes, but also accurately evaluate the performance through the independent test set, and finally comprehensively improve the reliability and stability of the model in the risk prediction of construction engineering safety accidents.

Figure 5 shows the frequency distribution analysis of each category in the dataset. In order to evaluate the prediction efficiency of the model, this study randomly selected part of the data from the data set of 101 accident records to construct a test set. The overall accuracy of the model is calculated to be 92.08%. In-depth analysis for the prediction of minor risks, its accuracy rate is 100%. The prediction accuracy rates reached 90.9%, 71.43%, and 75% for very serious, serious, and more serious risks, respectively.

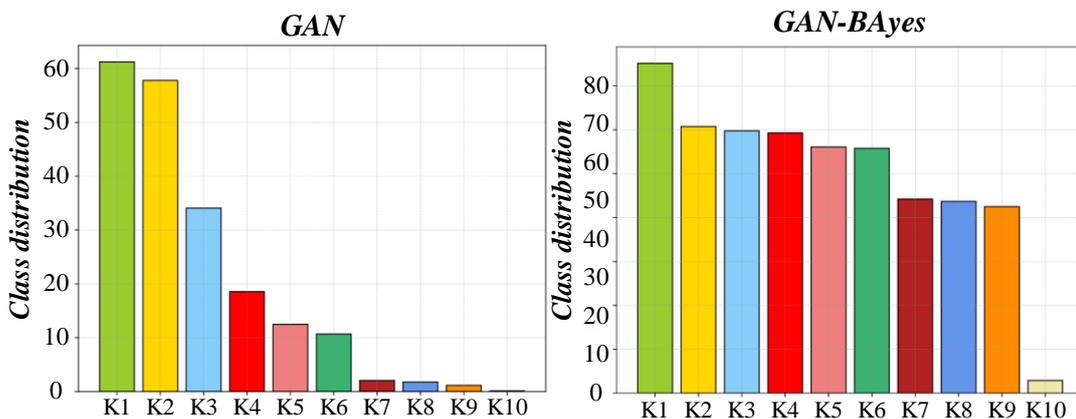


Figure 5: Class distribution of dataset

In this study, the prediction effectiveness metrics for each accident risk are computed, and Figure 6 reveals the effect of class K on the performance of FEDPE (Federated Policy Gradient with Byzantine Resilience) and FEDPG. The GAN-Bayes model achieves 99.86% precision and recall in the less severe accident risk; in the very severe accident risk case, the recall reaches 91%. In the

comprehensive analysis, the F-Score of the model exceeds 0.75, which proves that the model possesses an overall excellent performance. In addition, all types of accident risks exhibit 97% specificity, while the false positive rate is controlled at less than 3%. The higher specificity and lower FPR value prove the accuracy and effectiveness of the model in distinguishing different risk classes.

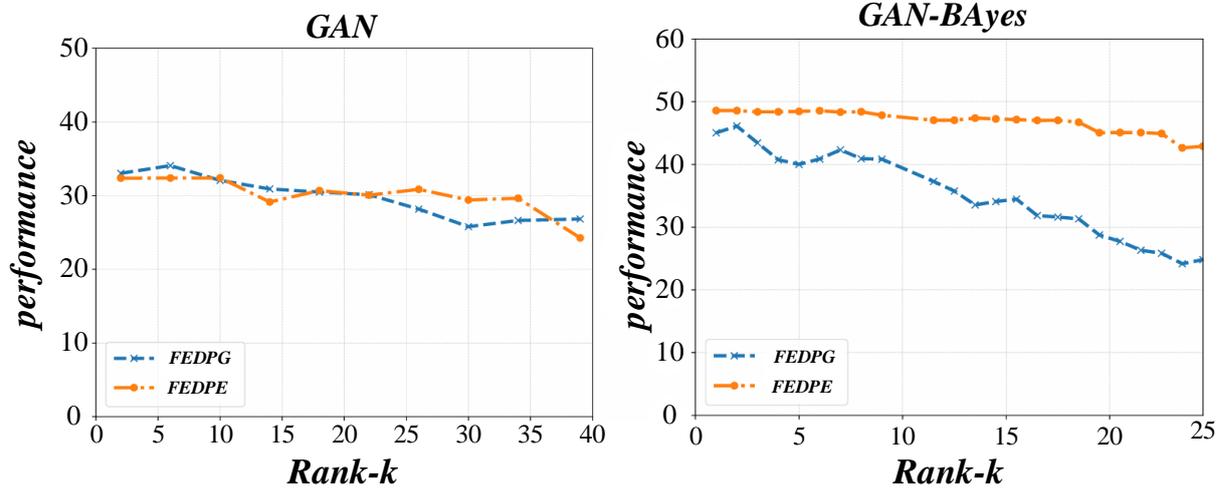


Figure 6: Effect of grade K on the performance of FEDPE and FEDPG

5 Example validation and analysis

5.1 Example verification and analysis

In order to evaluate the robustness of the construction engineering safety accident risk prediction model based on multi-objective optimized GAN network in the feature loss scenario, the simulated feature loss test was carried out in this study. Simulate feature loss in real engineering data by artificially introducing missing values of 5% to 50% in the raw data. For the missing features, mean imputation is used to process numerical data to retain statistical features, mode imputation is used to fill in the sub-type data to maintain the classification logic, and the features with a missing rate of more than 30% are discarded to avoid interfering with the performance of the model. After setting different missing rates each time, the model was retrained, and the changes in prediction accuracy, F1 value and other indicators were monitored,

and the robustness and adaptability of the model to deal with the problem of feature loss were comprehensively verified by systematically comparing the performance of the model under different processing strategies and missing rates.

This research experiment uses the Kaggle dataset, which contains two million two hundred and ninety-nine thousand records of traffic accidents in the United States from 2016 to 2019, to validate the effectiveness of the building safety impact assessment method. The dataset includes forty-nine accident metrics, including essential information such as the duration of traffic flow interruption after a traffic accident, the duration of accident processing, the length of the affected roadway, and relevant environmental factors. Figure 7 shows the precision-recall plot derived from the GAN model. The dataset is subdivided into four accident severity classes based on the degree of disruption to traffic operations after an accident.

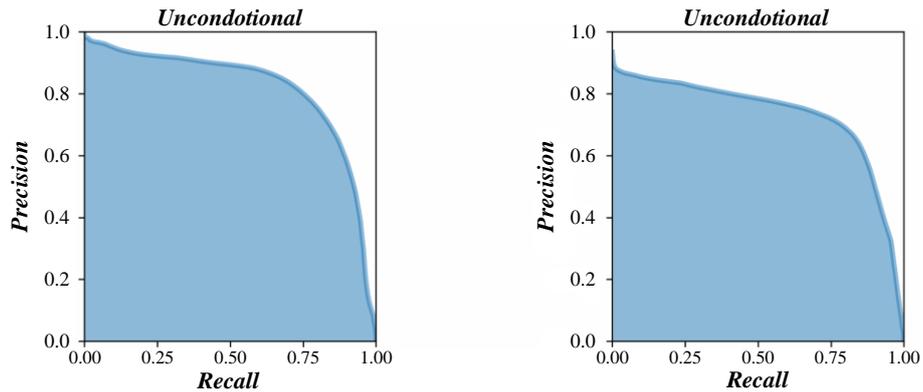


Figure 7: Accuracy-recall

In the data preprocessing process, we removed the data with a high proportion of missing values and retained more than 100,000 samples. A thousand pieces of data were randomly selected from each type of accident impact level, totaling 4,000 accident examples, to evaluate their correlation with building safety impacts. Of these four thousand accidents instances, features related to predicting building safety impacts totaled twenty-eight. These twenty-eight-dimensional accident features are labeled individually and serve as the underlying data set. Finally, these 4,000 accident samples are randomly assigned into training sets and test sets according to the ratio of eight to two to verify the performance and prediction ability of the model.

5.2 Importance analysis of accident characteristics

The GAN-Bayes model was utilized to quantify 28 accident-related features and analyze the contribution scores in the target time domain. The distribution of the given samples in the target t-domain is shown in Figure 8. The samples at different time stages exhibit significant variations; the thresholds for the actual cumulative contribution scores are set to $\alpha = 89$ and $\beta = 96$, and the fuzzy region contains four accident features. Observations show that the accuracy of the classifier increases with the number of accident features until $n = 9$, when the classifier performance reaches its peak.

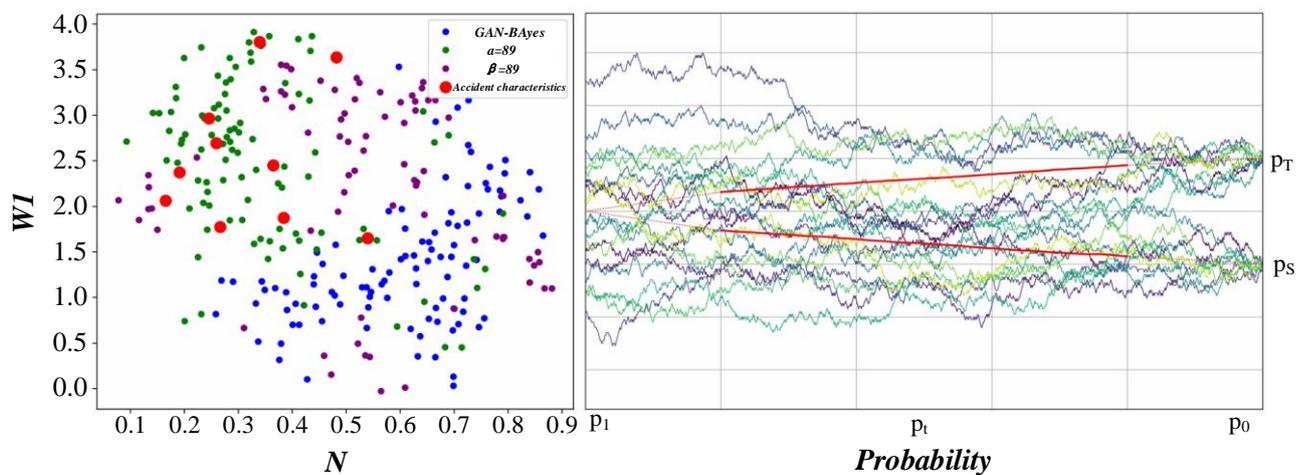


Figure 8: Given samples of target t-domain

When the accident features in the fuzzy region are selected as the classification basis, the classifier's performance declines slightly from $n \geq 10$ and gradually becomes stable. This phenomenon shows that the accident features in fuzzy areas do not positively impact the

classification results, so these features are excluded and finally, based on 28 accident features, nine optimal features were screened out to predict traffic risk status level.

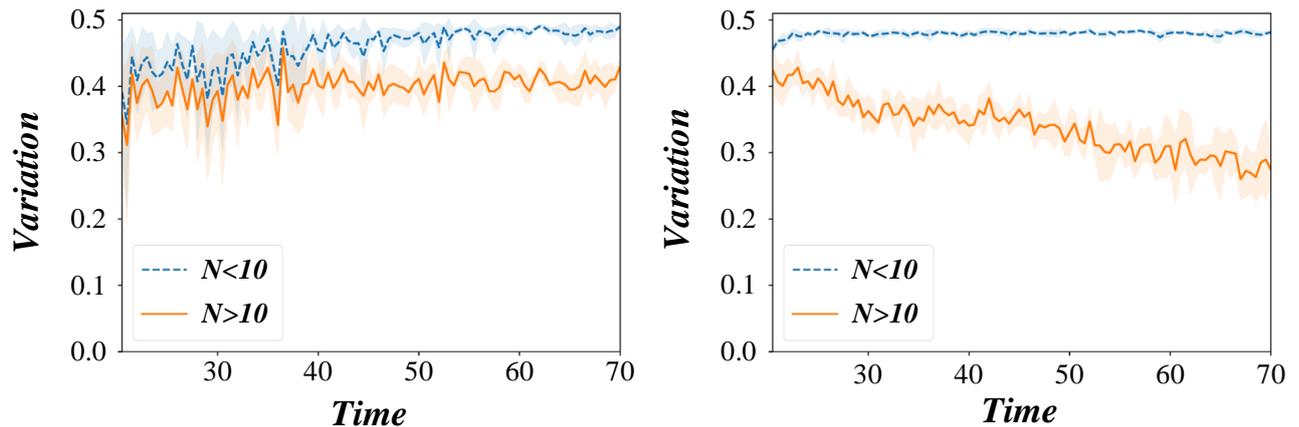


Figure 9: Variation of samples in different periods

As shown in Figure 9, according to the analysis of main accident characteristics in the accident database, two environmental variables, weather conditions and lighting status, are particularly critical when evaluating the impact of building safety. Abnormal traffic incidents caused by extreme environmental conditions usually bring significant traffic congestion and increased safety risks. There are two reasons for this phenomenon: First, environmental factors such as slippery road surfaces and insufficient light increase the coverage of the accident area; Secondly, the unfavorable traffic environment not only creates obstacles to the passage of rescue vehicles but also increases the complexity and time-consuming of rescue work, thus prolonging the emergency response time.

6 Comparative experimental analysis

In the study, the overall accuracy refers to the proportion of the model correctly classified among all predicted samples, which is calculated as (total number of correctly predicted samples/total number of samples) \times 100%, which reflects the model's ability to classify the overall data, and is suitable for evaluating the comprehensive

performance in the scenario of balanced data distribution. However, in the risk prediction of construction safety accidents, there is often a category imbalance in the data (such as a small proportion of high-risk accident samples), so it is necessary to supplement the accuracy of specific categories, that is, the accuracy of the model for each risk level (such as low, medium, and high risk) separately. For example, a high-risk category with an accuracy of \times 100% (correctly predicted high-risk samples / actual high-risk samples) reveals how effective the model is at identifying key risk categories. By reporting both overall accuracy and category-specific accuracy, assessment bias caused by data imbalance can be avoided, ensuring more targeted comparisons between different models or methods, and providing a more reliable basis for construction safety decisions.

To assess the efficacy of the GAN-Bayes algorithm in predicting building safety influence level, this study uses Support Vector Machines (SVM) and Random Forests (RF) as the baseline feature classification techniques for comparative analysis. The model structure is compared with the traditional deep forest (GCF) without feature selection, and the integrated model formed by combining XGBoost feature selection and deep forest (XGBoost-GCF) is included in the comparison.

Table 3: Prediction results of accident impact degree of different algorithms

Model	Accuracy rate	Recall rate	F1-score	Accuracy
SVM	0.7846	0.7990	0.8425	0.7901
RF	0.8356	0.8327	0.8421	0.8345
GCF	0.7665	0.7810	0.7712	0.7712
XGBoost -GCF	0.8664	0.7910	0.8939	0.8576
GAN-Bayes	0.9246	0.8451	0.8961	0.8879

Table 3 presents the experimental results of different prediction algorithms applied to the accident dataset. Summarizing the prediction efficacy of each method, the GAN-Bayes algorithm proposed in this paper performs well, with a prediction accuracy of 92.46%, and outperforms SVM, RF, GCF and XGBoost-GCF algorithms in all the evaluation metrics. Accordingly, it

can be inferred that the GAN-Bayes algorithm performs excellently in predicting the impact of building safety.

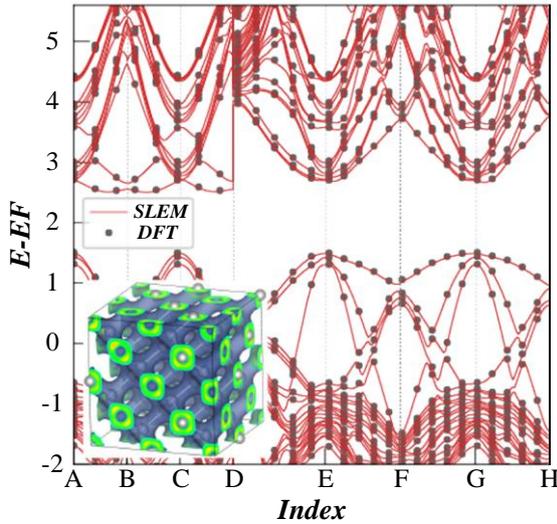
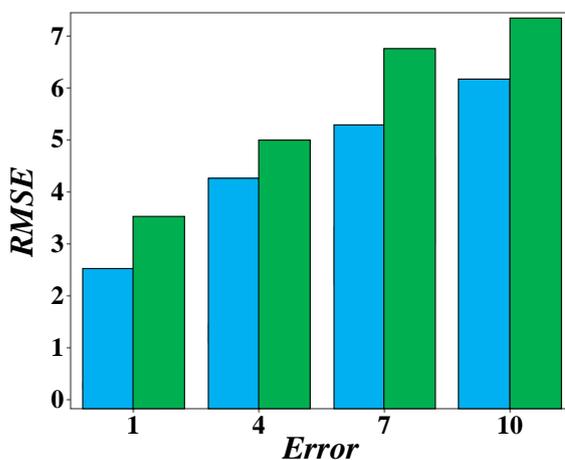


Figure 10: E-EF plots of A-H

The Event - Event Factor graph is an important visualization tool. It aims to visually present the events related to construction safety accidents (E, Event), such as falls from heights, collapses and other accidents themselves, as well as the various factors that induce these events (EF (Event Factor), including personnel illegal operation, equipment aging and failure, harsh environment, etc. By constructing this diagram, the causal correlation and action path between accident events and factors can be clearly sorted out, which can help researchers analyze the accident mechanism more systematically, and then accurately screen the key characteristic variables for the risk prediction model, and improve the accuracy and reliability of the model for the risk prediction of construction engineering safety accidents. Figure 10 reveals the E-EF plot of the A-H algorithm. The overall ROC curve of the GAN-Bayes algorithm and its ROC curves for different impact levels converge to the upper left quadrant, and the corresponding AUC values converge to the ideal value of 1, reflecting the stability and accuracy of the model in predicting the impact levels of various accidents. Compared with the traditional model, the model optimized by feature



selection significantly improves accuracy.

To verify GAN-Bayes' ability to deal with incomplete accident features, we randomly deleted some of the leading accident features in the test set to simulate the missing features caused by the untimely data collection during the actual abnormal accidents. The feature missing rate is set to 70%, 50%, 30%, 10%, and 0%, and the corresponding percentage of valid accident features are 30%, 50%, 70%, 90%, and 100%, respectively.

Figure 11 compares the time and memory consumption of different tensor product implementations. It can be seen that even when the feature missing rate is high, the model in this paper still maintains a certain prediction ability. With the increase in the number of effective accident features, the model's prediction performance gradually improves. This result further proves the effectiveness of accident feature enhancement based on GAN, indicating that when an actual accident occurs, with the continuous collection of accident information, the model's prediction accuracy will be continuously improved, providing a more accurate judgment basis for subsequent rescue and traffic diversion.

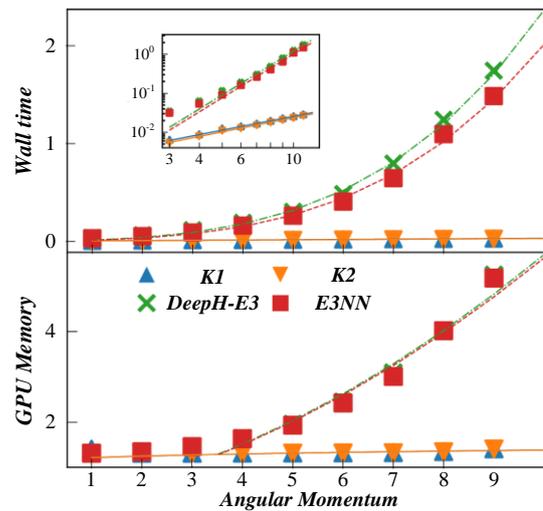


Figure 11: Compares the time and memory consumption of different tensor product implementations

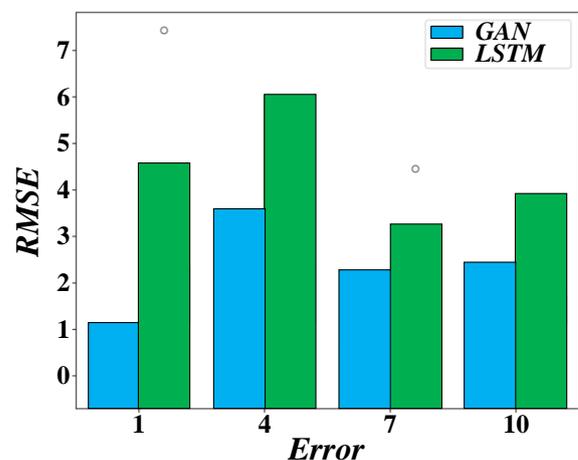


Figure 12: Characteristic importance between accident categories

Figure 12 shows the Root Mean Square Error (RMSE) corresponding to the feature importance between different accident categories in the construction engineering safety accident risk prediction model based on multi-objective optimized GAN network. The blue and green colors used in the figure represent the GAN and LSTM models, respectively, with the horizontal axis representing the Error and the vertical axis representing the RMSE value. As can be seen from the figure, the RMSE values of the two models are different at different error points, reflecting their different performance in dealing with the importance of different features in the risk prediction of construction engineering safety accidents, which is helpful to compare and evaluate the performance of the two models in this study.

7 Conclusion

Through in-depth theoretical discussion and extensive experimental verification, the research on the risk prediction model of construction engineering safety accidents based on multi-objective optimized GAN network has achieved significant research results and important practical application value. This study not only theoretically constructs a multi-objective optimized GAN model that can comprehensively consider multiple dimensions such as data fidelity, model diversity and stability, but also verifies the excellent performance of the model in improving prediction performance through a large number of experiments. This achievement not only provides a more accurate and reliable tool for the risk prediction of construction engineering safety accidents, but also promotes the field to move towards intelligence and science, and contributes an important force to the safety production and sustainable development of the construction industry. The following are the main conclusions of this study:

The prediction model based on multi-objective optimization GAN network proposed in this study shows excellent performance in architectural engineering security incident risk prediction. Compared with traditional prediction methods, the accuracy of this model on the test set has been greatly improved, with an accuracy rate as high as 92.46%. At the same time, it also performs well in key evaluation indicators such as recall rate and F1 value. The results show that the multi-objective optimization GAN network can more effectively capture the complex features of security incident risks and improve the accuracy and reliability of prediction.

By introducing multi-objective optimization strategies, this study not only optimizes the generation and discrimination process of GAN networks, but also makes the model show stronger adaptability and stability in the face of different data sources and complex environments. The improvement of this generalization ability has laid a solid foundation for the wide application of the model in practical engineering projects.

In future research, this study will explore the integration and application prospects of models with the Internet of Things, big data, and other technologies. By

building an intelligent and integrated safety accident risk early warning and management system, it provides strong technical support and decision-making basis for safety management in the architectural engineering industry.

References

- [1] Mahmoud AlJamal, Ala Mughaid, Bashar Al shboul, Hani Bani-Salameh, Shadi Alzubi, and Laith Abualigah, "Optimizing risk mitigation: A simulation-based model for detecting fake IoT clients in smart city environments," *Sustainable Computing: Informatics and Systems*, vol. 43, pp. 101019, 2024.
- [2] Kemal Hacıefendioglu, Fatemeh Mostofi, Vedat Togan, and Hasan Basri Basaga, "CAM-K: a novel framework for automated estimating pixel area using K-Means algorithm integrated with deep learning based-CAM visualization techniques," *Neural Computing & Applications*, vol. 34, no. 20, pp. 17741-17759, 2022.
- [3] Ping Chai, Lei Hou, Guomin Zhang, Quddus Tushar, and Yang Zou, "Generative adversarial networks in construction applications," *Automation in Construction*, vol. 159, pp. 105265, 2024.
- [4] Shanu Verma, Millie Pant, and Vaclav Snasel, "A Comprehensive Review on NSGA-II for Multi-Objective Combinatorial Optimization Problems," *Ieee Access*, vol. 9, pp. 57757-57791, 2021.
- [5] Niloofar Nadim Kabiri, Saeed Emami, and Abdul Sattar Safaei, "Simulation-optimization approach for the multi-objective production and distribution planning problem in the supply chain: using NSGA-II and Monte Carlo simulation," *Soft Computing*, vol. 26, no. 17, pp. 8661-8687, 2022.
- [6] Yali Lv, Jingpu Duan, and Xiong Li, "A survey on modeling for behaviors of complex intelligent systems based on generative adversarial networks," *Computer Science Review*, vol. 52, pp. 100635, 2024.
- [7] Akshay Kumar and T. V. Vijay Kumar, "Multi-Objective Big Data View Materialization Using NSGA-III," *International Journal of Decision Support System Technology*, vol. 14, no. 1, 2022.
- [8] F. Dabbaghi, A. Tanhadoust, M. L. Nehdi, S. Nasrollahpour, M. Dehestani, and H. Yousefpour, "Life cycle assessment multi-objective optimization and deep belief network model for sustainable lightweight aggregate concrete," *Journal of Cleaner Production*, vol. 318, 2021.
- [9] Zixi Hu, Shuang Liu, Fan Yang, Xiaodong Geng, Xiaodi Huo, and Jia Liu, "Research on Multi-objective Optimization Model of Power Storage Materials Based on NSGA-II Algorithm," *International Journal of Computational Intelligence Systems*, vol. 17, no. 1, 2024.
- [10] Tome Sicuaio, Pengxiang Zhao, Petter Pilesjo, Andrey Shindyapin, and Ali Mansourian, "Sustainable and Resilient Land Use Planning: A Multi-Objective Optimization Approach," *Isprs International Journal of Geo-Information*, vol. 13, no. 3, 2024.

- [11] Hongjing Wei et al., "Unified Multi-Objective Genetic Algorithm for Energy Efficient Job Shop Scheduling," *Ieee Access*, vol. 9, pp. 54542-54557, 2021.
- [12] Xuan Xia et al., "GAN-based anomaly detection: A review," *Neurocomputing*, vol. 493, pp. 497-535, 2022.
- [13] Muhammad Ameer and Mohammed Dahane, "NSGA-III-based multi-objective approach for reconfigurable manufacturing system design considering single-spindle and multi-spindle modular reconfigurable machines," *International Journal of Advanced Manufacturing Technology*, vol. 128, no. 5-6, pp. 2499-2524, 2023.
- [14] Zheng Gan, Xiongya Shen, Xiang Liu, and Ying Xie, "Building feature extraction based on natural neighborhood decomposable point feature extraction algorithm," *Informatica*, vol. 48, no. 22, 2024.
- [15] Wejden Gazehi, Rania Loukil, and Mongi Besbes, "Classification of a nanocomposite using a combination between Recurrent Neural Network based on Transformer and Bayesian Network for testing the conductivity property," *Expert Systems with Applications*, vol. 270, pp. 126518, 2025.
- [16] Donghui Shi, Shuling Gan, Jozef Zurada, Jian Guan, Feilong Wang, and Pawel Weichbroth, "A multi-model approach to construction site safety: Fault trees, Bayesian networks, and ontology reasoning," *Expert Systems with Applications*, vol. 288, pp. 127817, 2025.
- [17] Zhongxiang Chang, Zhongbao Zhou, Ruiyang Li, Helu Xiao, and Lining Xing, "Observation scheduling for a state-of-the-art SAREOS: Two adaptive multi-objective evolutionary algorithms," *Computers & Industrial Engineering*, vol. 169, 2022.
- [18] Xiaoyang Liu et al., "Semi-supervised community detection method based on generative adversarial networks," *Journal of King Saud University - Computer and Information Sciences*, vol. 36, no. 3, pp. 102008, 2024.
- [19] Pengda Wang, Zhaowei Liu, Zhanyu Wang, Zongxing Zhao, Dong Yang, and Weiqing Yan, "Graph generative adversarial networks with evolutionary algorithm," *Applied Soft Computing*, vol. 164, pp. 111981, 2024.
- [20] Jian Hu, Yingjun He, Wenqian Xu, Yixin Jiang, Zhihong Liang, and Yiwei Yang, "Anomaly Detection in Network Access-Using LSTM and Encoder-Enhanced Generative Adversarial Networks," *Informatica*, vol. 49, no. 7, 2025.
- [21] Juan Li, "Improved Genetic Algorithm Enhanced with Generative Adversarial Networks for Logistics Distribution Path Optimization," *Informatica*, vol. 49, no. 11, 2025.
- [22] Ima Okon Essiet and Yanxia Sun, "Tracking Variable Fitness Landscape in Dynamic Multi-Objective Optimization Using Adaptive Mutation and Crossover Operators," *Ieee Access*, vol. 8, pp. 188927-188937, 2020.
- [23] Chuyue Zhang and Yuchen Meng, "Bayesian deep learning: An enhanced AI framework for legal reasoning alignment," *Computer Law & Security Review*, vol. 55, pp. 106073, 2024.
- [24] Zhansheng Liu, Chengkuan Ji, Guoliang Shi, and Yanchi Mo, "Structural safety risk prediction method for terminal building steel roof construction considering spatial and temporal variations," *Journal of Constructional Steel Research*, vol. 224, pp. 109126, 2025.
- [25] Jiaqi Wang, Yuqing Fan, Xi Pan, Jun Sun, and Limao Zhang, "Multi-source information fusion for dynamic safety risk prediction of aerial building machine using spatial – temporal multi-graph convolution network," *Advanced Engineering Informatics*, vol. 65, pp. 103261, 2025.
- [26] Claudio Gallicchio and Alessio Micheli, "Architectural richness in deep reservoir computing," *Neural Computing & Applications*, vol. 35, no. 34, pp. 24525-24542, 2023.
- [27] Yair Schwartz, Rokia Raslan, Ivan Korolija, and Dejan Mumovic, "A decision support tool for building design: An integrated generative design, optimisation and life cycle performance approach," *International Journal of Architectural Computing*, vol. 19, no. 3, pp. 401-430, 2021.
- [28] Stephen John Warnett and Uwe Zdun, "Architectural Design Decisions for the Machine Learning Workflow," *Computer*, vol. 55, no. 3, pp. 40-51, 2022.
- [29] Jinmo Rhee, Pedro Veloso, and Ramesh Krishnamurti, "Three decades of machine learning with neural networks in computer-aided architectural design (1990-2021)," *Design Science*, vol. 9, 2023.
- [30] Zhipeng Zhou, Wen Zhuo, Jianqiang Cui, Haiying Luan, Yudi Chen, and Dong Lin, "Developing a deep reinforcement learning model for safety risk prediction at subway construction sites," *Reliability Engineering & System Safety*, vol. 257, pp. 110885, 2025.

Glossary

Abbreviation	Full Name	Description
<i>TRI</i>	True Risk Influence	Measures the impact of variables on risk levels via the arithmetic mean of High-Risk Influence (HRI) and Low Risk Influence (LRI)
<i>GCF</i>	Generalized Deep Forest	A baseline classification model used for performance comparison with the GAN-Bayes model
<i>XGBoost-GCF</i>	XGBoost Feature Selection-Deep Forest	An integrated model combining XGBoost feature selection with Deep Forest, used in comparative experiments to validate GAN-Bayes effectiveness
<i>FEDPE</i>	Federated Policy Gradient with Byzantine Resilience	A federated learning strategy enhancing model robustness in distributed data scenarios
<i>FEDPG</i>	Federated Policy Gradient	A baseline federated learning strategy for comparing model stability under data missing scenarios
<i>m</i>	Data volume	Used in objective function calculations: $V = \sum \log D(x)$, representing the total number of samples
<i>V</i>	Value function	Measures performance of generator/discriminator; core objective function in GAN training
<i>D</i>	Discriminator	Neural network model to distinguish real/generated data; outputs 0 (generated data) or 1 (real data)
<i>G</i>	Generator	Neural network model to produce realistic data; takes random noise Z as input and outputs simulated samples
<i>Z</i>	Random noise input	Input variable for the generator to produce simulated samples
<i>p</i>	Probability distribution	Core GAN objective: making $p(G(z))$ approximate the real data distribution p_{data}
<i>A_i/A_j</i>	Attribute variables in GAN-Bayes structure learning	Calculates conditional mutual information $I(A_i; A_j C)$ to identify variable dependencies
<i>C</i>	Conditional variable/set	Used in conditional probability calculations or as a conditional set in mutual information
<i>A</i>	Attribute variables	Factors such as weather, equipment status as risk factors; measures correlation with risks
<i>TP/TN</i>	True Positive/True Negative samples	Confusion matrix metrics for accuracy calculation
<i>FP/FN</i>	False Positive/False Negative samples	Confusion matrix metrics for recall calculation
α/β	Thresholds for cumulative contribution scores	Thresholds set in the paper to filter key features in accident feature importance analysis

μ_0	Population mean of differences in paired samples T-test	Null hypothesis parameter in paired T-test
n	Sample size	Number of paired samples or population size in genetic algorithms
d	Difference in paired samples/number of objective functions	Difference in paired samples or number of objective functions for gradient updates
\log	Natural logarithm	Used in objective functions and mutual information calculations
E	Complete event set	Used in total probability formula in Bayesian networks

AB-YOLOv8: Attention-based Feature Extraction model for Underwater Object Detection

Pratima Sarkar^{1,2}, Sandeep Gurung¹, Bitan Misra², Sourav De³

¹Sikkim Manipal Institute of Technology, Department of Computer Science and Engineering, Majhitar, Sikkim-737132, India

²Department of Computer Science and Engineering, Techno International New Town, Kolkata-900156, India

³Department of Computer Science and Engineering, Government College of Engineering and Textile Technology, Serampore 12, William Carey Road, Serampore, Hooghly, Pin-712201, India

E-mail: pratima.sarkar@tint.edu.in, bitan.misra@tint.edu.in, dr.sourav.de79@gmail.com, sandeep.gu@smit.smu.edu.in

Keywords: Channel attention, data augmentation, spatial attention, R-CNN, underwater object detection

Received: December 23, 2024

Accurate and timely underwater object detection is crucial in the field of marine environmental engineering. The detection of such targets has been improved recently using techniques based on Convolutional Neural Networks (CNN). However, the processing performance of deep neural networks is typically inadequate due to their high parameter requirements. Accurate detection is difficult with current techniques when dealing with small, close-packed underwater targets. In order to overcome these problems, the proposed work combined YOLOv8 with different attention modules and proposed a novel neural network model to enhance underwater object detection capabilities. In this research, AB-YOLOv8 is proposed, which adds the attention mechanism to the original YOLOv8 design. To be more precise, the proposed work introduced four attention modules, Convolutional Block Efficient Channel Attention (ECA), Shuffle Attention (SA), Global Attention Mechanism (GAM), and Attention Module (CBAM), to create the enhanced models and train them in the aquarium dataset. Each of the attention blocks is combined with YOLOv8 to improve the performance of the entire object detection. The residual block is introduced into the CBAM to optimize the performance of the CBAM. The detailed experiments are conducted on the aquarium dataset, and various performance assessment parameters are used, like mAP, FLOPS, Params, inference time, etc. After performing the experiment, it was found that ECA gives the best result out of all attention blocks and improved mAP value by 8%, also reduced the number of parameters generated during training. To validate the work, we also performed the experiment on the Brackish dataset, and we found that ECA outperforms other attention mechanisms with YOLOv8.

Povzetek: Zasnovan je nov model AB-YOLOv8 z mehanizmi pozornosti (ECA, CBAM, SA, GAM) za izboljšanje zaznavanja podvodnih objektov. Model ECA-YOLOv8 je izkazal najboljše rezultate: izboljšal je metriko mAP v primerjavi z osnovnim YOLOv8 in zmanjšal število parametrov.

1 Introduction

Underwater object recognition is a crucial stage in image processing that is important for a number of applications, including marine sciences and the upkeep and repair of sub-aquatic infrastructure. One of the most difficult study areas in modern computer vision technologies is the detection of underwater objects [1]. Specifically, the widespread deployment of digital cameras on Autonomous Underwater Vehicles (AUVs) and Unmanned Underwater Vehicles (UUVs) has led to an exponential increase in the availability of underwater imagery in recent years [2]. The primary obstacles to underwater vision are the increased expense of the devices, their intricate configuration, and the distortion of light and signal propagation caused by the water medium [3]. The propagation of light in underwater environments is particularly affected by phenomena such as absorption and

scattering, which have a significant impact on visual perception [4, 5]. In recent years, generic object detection algorithms have demonstrated their exceptional performance. In digital image processing for object recognition and classification, deep learning, also referred to as deep machine learning or deep structured learning-based techniques, has recently seen significant success [6]. Thus, they are attracting the interest and popularity of the computer vision research community rather quickly [7]. However, these approaches are not sufficiently capable of handling underwater object detection due to the following challenges: (1) Real-world applications typically feature small objects with hazy photos [8], and (2) real-world applications and underwater datasets have images with heterogeneous noise [9]. When taking into account underwater variables like sufficient light, reasonable current intensity, and clear underwater eyesight, simple underwater target-detection tech-

niques can be used more effectively. The primary features extracted by early conventional detection techniques were color, texture, and geometry. As the deep learning technique continues to advance, neural networks have emerged as underwater target-detection frameworks that enable target detection by identifying and locating objects in photos [10]. However, underwater image quality deteriorates due to less-than-ideal conditions in practice, which consequently impairs the accuracy of detection. Convolutional Neural Networks (CNNs) [11] have made significant strides in object detection in recent years due to their potent feature learning and transfer learning capabilities, which have drawn increasing attention from the discipline of computer vision. The application of CNN to object detection for improved performance is therefore a significant domain of research work [12]. YOLOv8 differs from previous YOLO models in several significant ways. Its transformer-based architecture, which improves accuracy and performance, especially for small and difficult-to-detect objects, is one of the biggest upgrades.

In order to effectively address the challenges associated with underwater object detection, the proposed research integrates the YOLOv8 [13] architecture with various attention modules, culminating in the development of a novel neural network model designed to significantly enhance detection capabilities in underwater environments. The unique combination of YOLOv8 with sophisticated attention mechanisms and the calculated improvements made to the CBAM constitute the work's originality. The following are the primary contributions of this paper:

- 1 This innovative approach is encapsulated in the newly introduced model, termed AB-YOLOv8, which incorporates an attention mechanism into the foundational design of YOLOv8. This study introduces four distinct attention modules: Convolutional Block Efficient Channel Attention (ECA) [14], Shuffle Attention (SA) [15], Global Attention Mechanism (GAM) [16], and Convolutional Block Attention Module (CBAM) [17]. Each of these modules is strategically combined with the YOLOv8 framework to create enhanced models that are specifically trained on the aquarium dataset. The integration of these attention blocks is aimed at improving the overall performance of object detection tasks, particularly in the challenging underwater context, where visibility and clarity are often compromised.
- 2 Additionally, the study improves the CBAM by adding a residual block, which helps to maximize its efficiency. This innovation makes better feature extraction and representation possible, which enhances the model's capacity to identify items in intricate underwater environments.
- 3 The success of the suggested models is evaluated using a range of performance assessment metrics, such as mean Average Precision (mAP), FLOPS (Floating

Point Operations Per Second), number of parameters (Params), and inference time [18]. In real-time underwater detection applications, these measures are crucial for understanding the trade-offs between accuracy and processing efficiency.

The structure of the paper is as follows. In Section 2, the relevant literature is discussed. The network architecture and adopted approach are presented in Section 3. The dataset description is given in Section 4, and the experimental evaluation parameters are shown in Section 5. Experimental results and discussions are included in Section 7 and 8 respectively. Future work, our findings, and research outlook are summed up in Section 9.

2 Literature review

Underwater object detection can be accomplished by different two-stage and single-stage object detectors. The most popular two-stage detectors are R-CNN, Fast R-CNN, and Faster R-CNN. R-CNN [19] is performing better for small object detection, but it is not suitable for real-time object detection. So, many researchers have selected the single-stage object detectors, i.e., YOLO series, as the foundation for future development in order to accomplish real-time underwater object identification. The YOLO-UOD [20] optimization algorithm, a unique underwater object identification technique based on YOLOv4-tiny research, is presented in the article [21]. The suggested approach, which combines the symmetric FPN-Attention module and the symmetric dilated convolutional module, may efficiently collect important characteristics and contextual information while maintaining deep features, according to experimental results on the Brackish undersea dataset. Its underwater object detection mAP score of 87.88% is superior to YOLOv5s and YOLOv5m and higher than YOLOv4-Tiny's score of 77.38%. In [22], the Transformer encoder and a coordinate attention module were integrated into YOLOv5 to create a new detection network called TC-YOLO. Underwater picture enhancement was done using the CLAHE [23] algorithm, while label assignment in training was done using the optimal transport assignment approach. By combining these methods, our suggested strategy maintained computational efficiency for real-time underwater detection tasks while achieving state-of-the-art performance on the RUIE2020 [24] dataset. The attachment of the coordinate attention module to the end of the neck was found to be a very successful and efficient method of enhancing detection networks' performance in the ablation experiments. Article [25] includes the plug-and-play mDFLAM with YOLO detectors to satisfy the high-precision and real-time demands for underwater object detection. By enhancing the quality of feature fusion between scales, the full-port embedding significantly reinforces the expression of semantic information. Using a lightweight backbone network built on deformable convolution YOLOv3, article [26] proposes a dynamic YOLO de-

detector with certain specialized designs for small item identification. Experimental findings on the Pascal VOC and MS COCO datasets further support the superiority of the suggested model. Article [27] proposes a high detection accuracy cascade model based on the UGC-YOLO network structure. Additionally, PPM pooling is added to the top layer network for the purpose of aggregating semantic data, and deformable convolution is utilized to capture long-range semantic dependencies. Lastly, a multi-scale weighted fusion method for learning semantic data at various scales is introduced. The suggested approach has been shown through experiments on an underwater test dataset to be able to identify aquatic targets in intricately deteriorated underwater images. In order to decrease feature interference and increase detection accuracy, an enhanced YOLO detection technique without anchor points is presented [28], in which the detection and recognition features are kept apart. Additionally, a technique for improving underwater photos based on Retinex is also suggested. To confirm the efficacy of the suggested improved YOLO detection technique, pertinent tests based on underwater datasets are carried out. In order to create a quick, precise, and compact neural network model that can identify goldfish breeds in real time, the authors of the research [29] examine the impact of shrinking the size of the pre-trained MobileNetV2, which serves as the foundation of the YOLOv2 object detection framework. Paper [30] proposes the YOLO-SC algorithm as a solution to the problem of finding the submarine cable's position and feature information using the YOLOv3 [31, 32] prototype network because of the blurry and blue-green underwater images. Three enhanced modules work together to address the aforementioned issues. The multi-structured multi-size feature fusion module improves the efficiency of feature information extraction; the light-weighted module streamlines the prediction network and reduces identification duration; and the skip connection module, which is included in the residual network, enhances the extraction of position information. Another modified

3 Methodology

Recently, the attention mechanism has achieved outstanding outcomes in the domain of object detection. Attention blocks are capable of selecting most significant features and discarding irrelevant features. This study integrates the attention module into the neck and head component of YOLOv8 in order to improve the detection of important characteristics and reduce the impact of irrelevant information. We have chosen four attention mechanism like Efficient Channel Attention (ECA), Convolutional Block Attention Module (CBAM), Shuffle Attention (SA) and Global Attention Mechanism (GAM) for feature aggregation. ECA was selected because of its lightweight design and capacity to enhance channel-wise feature recalibration without appreciably raising model complexity.

CBAM combines both channel and spatial attention, making it well-suited to capture complex underwater textures and cluttered scenes. SA helps in capturing long-range dependencies, which is beneficial when objects are partially occluded or dispersed. GAM enhances global context aggregation, helping to better differentiate between background and foreground in low-visibility underwater conditions.

YOLOv8 Architecture consists of different key components like backbone, neck, head and loss function as shown in figure 1. CSPDarknet used as backbone which contains CSP connections to increase information exchange. The neck work as a feature extractor, neck uses C2f architecture which integrates C3 modules. Neck aggregate features for detecting three different size of objects. YOLOv8 makes use of a number of detection modules to predict class probabilities, bounding boxes, and objectness scores for every grid cell in the feature map. The final detection are then obtained by averaging these forecasts. There are three types of loss function used during object prediction in YOLOv8 to optimize object detection those are: Binary Cross-Entropy (BCE), Distribute Focal Loss (DFL) and Complete Intersection over Union (CIoU) Loss. The classification component of YOLOv8 utilizes the Binary Cross-Entropy (BCE) Loss as its loss function, which is represented by the following equation:

$$BCE = -wt[x_n \cdot \log y_n + (1 - x_n) \cdot \log(1 - y_n)] \quad (1)$$

wt represents weight, x_n is labeled and y_n is predicted value. A DFL function is specifically developed to highlight the amplification of probability values about p . The equation is given as follows:

$$DFL = P_A + P_B \quad (2)$$

Where P_A is shown in eq(3) and P_B Shown in eq(4)

$$P_A = -[(p_{n+1} - p) \log(\frac{p_{n+1} - p_n}{p_{n+1} - p_n})] \quad (3)$$

$$P_B = (p - p_n) \log(\frac{p - p_n}{p_{n+1} - p_n}) \quad (4)$$

Incorporating the dimensions between the predicted bounding box and the ground truth bounding box, the CIoU Loss adds an influence factor to the Distance Intersection over Union (DIoU) Loss. The equation is as specified below:

$$CIoU = 1 - IoU + \frac{l^2}{c^2} + \frac{v^2}{1 - IoU + v} \quad (5)$$

IoU is intersection over union, d is Euclidean distance between predicted value and ground truth, l is diagonal length of predicted box, v is aspect ration of bounding box. In Figure 1 BBox-loss is combination of DFL and CIoU whereas Cls-loss represents BEC loss.

This work made modification on existing YOLOv8 architecture by adding attention module in neck and head of

YOLOv8 as illustrated in Figure 1. We have added one attention block in neck and rest all are added in head of YOLOv8. In proposed work used four different attention blocks i.e. Efficient Channel Attention (ECA), Convolutional Block Attention Module (CBAM), Shuffle Attention (SA) and Global Attention Mechanism (GAM). After incorporating these four different attention module into YOLOv8 analysed the performance of YOLOv8 with Attention block or AB-YOLOv8.

3.1 Attention modules

3.1.1 Efficient channel attention (ECA)

ECA mainly involves cross-channels and the use of 1D convolution with an adaptive single-dimensional convolution kernel as shown in Figure 2. Cross-channel interaction is an innovative method of merging characteristics to improve the representation of certain meanings. The input feature map I , which has dimensions $R^{C \times H \times W}$, is transformed into the aggregated feature F through the processes of Global Average Pooling (GAP) and cross-channel interaction. For the following equation, C refers to the cross-channel interaction.

$$F = C(GAP(I)) \quad (6)$$

ECA captures the local cross-channel interaction in aggregated data by examining the interaction between the features of each channel and their nearby k channels. The ECA method avoids utilizing 1D convolution for reducing dimensionality and effectively achieves multi-channel interaction. where the weights of the features F_i can be calculated as [14]:

$$w_i = \sigma(W) \quad (7)$$

where, W is a weight matrix and σ . is sigmoid function.

3.1.2 Convolutional block attention module (CBAM)

The CBAM [17] module has two attention sub-modules: the Channel Attention Module (CAM) and the Spatial Attention Module (SAM) as presented in Figure 3. The Channel Attention Module (CAM) is designed to enhance informative elements in the channel dimension, while the Spatial Attention Module (SAM) is designed to emphasize important features along the spatial axes. CBAM successfully captures the channel and spatial dependence in the input feature map by integrating these two attention processes. Input of CBAM is a feature map $I \in R^{C \times H \times W}$ then it is converted into 1D channel attention map $F_C \in R^{C \times 1 \times 1}$ and 2D spatial map $F_S \in R^{1 \times W \times H}$. So CBAM is a combination of following equations:

$$F = F_C \odot I \quad (8)$$

$$F' = F_S \odot F \quad (9)$$

where \odot is element wise multiplication.

In order to efficiently calculate the channel attention, compress the spatial dimension of the input feature map. CBAM employed both Global Average Pooled (GAP) and Global Max Pooled (GMP) features concurrently. Empirical findings have demonstrated that the utilization of both features significantly enhances the representational capacity of networks, as opposed to using each feature independently. Then element wise sum (+) and sigmoid (σ) function is used to find channel attention (F_C). Equation for channel attention is as follows:

$$F_C(I) = \sigma(MLP(GAP(I)) + MLP(GMP(I))) \quad (10)$$

In this equation I is input feature matrix and MLP is Multi Layer Perception. CBAM utilizes GAP and GMP along the channel axis for spatial attention, and subsequently combines them by concatenation (\oplus). The concatenation output is passed through a convolutional layer, and the resulting output is then used as the input for the sigmoid (σ) function. The spatial attention (F_S) is calculated using the following method.

$$F_S(I) = \sigma[CONV(GAP(I) \oplus GMP(I))] \quad (11)$$

3.1.3 Global attention mechanism (GAM)

GAM [16] adopts similar architecture as CBAM. GAM added additional shortcut connections between channel attention and spatial attention as depicted in Figure 4. The following equation represents GAM:

$$F_{out} = I + [F_S(F_C(I)) \times I] \times (F_C(I) \times I) \quad (12)$$

where, I is input feature, F_C channel attention block and F_S is spatial attention block.

To focus on specific channels, the GAM technique utilizes a 3D permutation from the beginning to preserve three-dimensional information. Afterwards, it utilizes a MLP to enhance the channel-spatial interdependence across dimensions. Following expression shows channel attention block representation:

$$F_C(I) = \sigma[RevPermutate(MLP(Permutate(I)))] \quad (13)$$

GAM utilizes two 7×7 convolution layers to combine spatial information for spatial attention as hown in eq. (14).

$$F_S(I) = \sigma[BN(f^{7 \times 7}(BN + ReLU(f^{7 \times 7}(I)))] \quad (14)$$

where, σ is sigmoid function, BN is batch normalization.

3.1.4 Shuffle attention (SA)

SA [15] divides the input feature maps into different groups, employing the Shuffle Unit to integrate both channel attention and spatial attention into one block for each group as shown in 5. Then these features are aggregated using spatial and channel attention. The channel attention mechanism utilizes the Global Average Pooling (GAP) technique

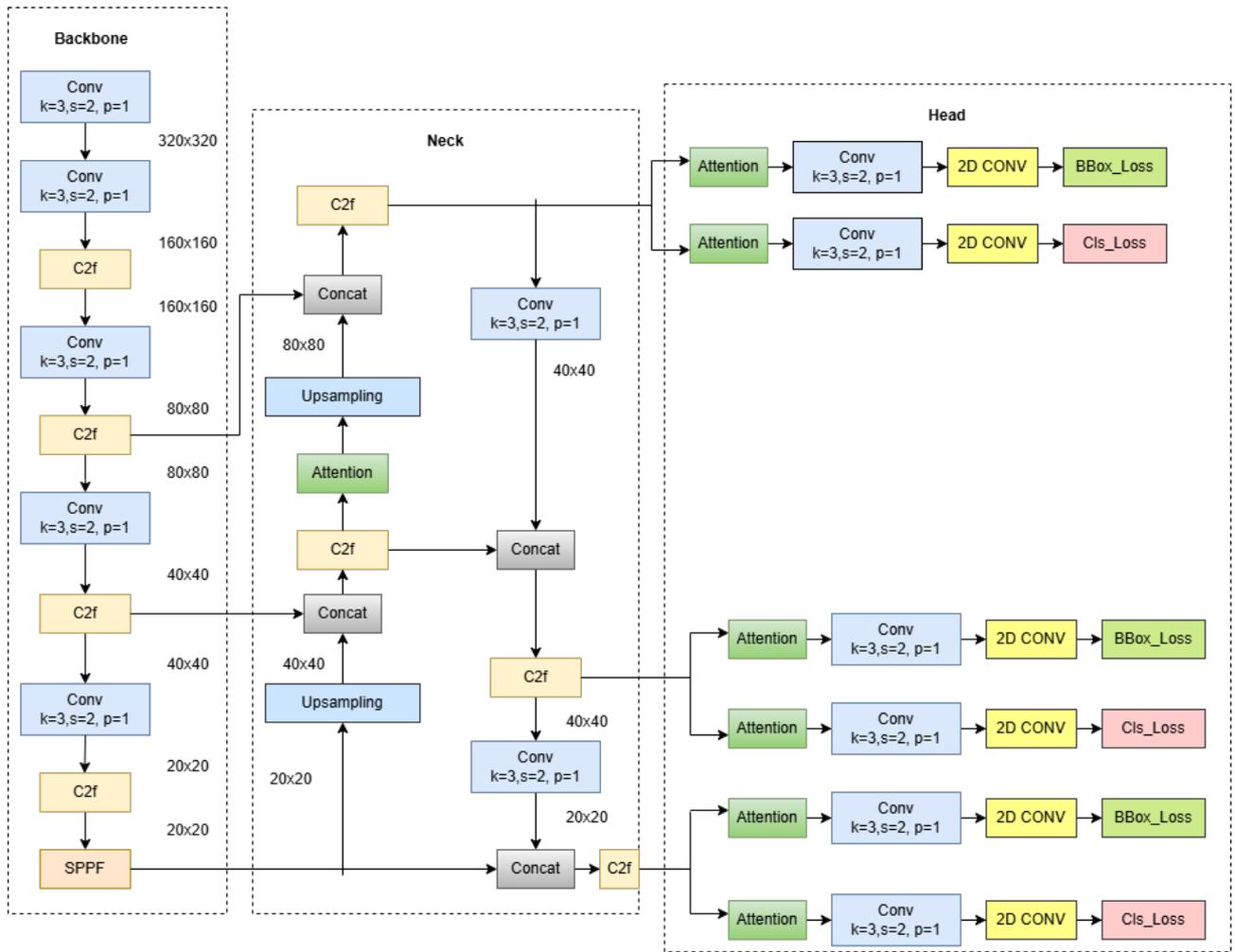


Figure 1: AB-YOLOv8 model architecture



Figure 2: Efficient channel attention

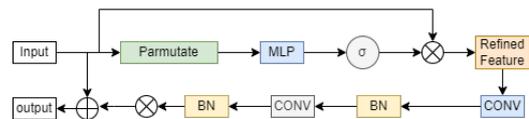


Figure 4: Global attention mechanism

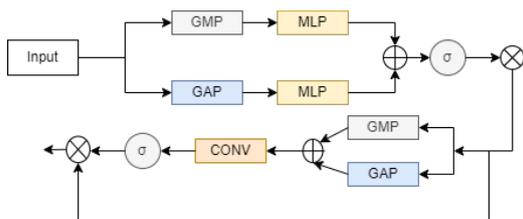


Figure 3: Convolutional block attention module

to acquire and incorporate global information for the specific sub-feature sb_1 . Furthermore, a straightforward gating mechanism employing sigmoid functions is utilized to generate a concise function that enables accurate and adaptable selection. Final output of channel attention is as follows:

$$CA = \sigma[f_c(GAP(sb_1))] \odot sb_1 \quad (15)$$

In spatial attention first step involves applying Group Normalization (GN) to the sub-feature sb_2 in order to calculate spatial-wise statistics. Afterwards, the output sub-feature sb_2 is improved through fully connected layer f_c , as demonstrated in the following equation.

$$SPA = \sigma[f_c(GN(sb_2))] \odot sb_2 \quad (16)$$

Algorithm 1 Attention-Integrated YOLOv8 for Object Detection**Require:** Input image $I \in \mathbb{R}^{H \times W \times 3}$, Ground truth labels (for training)**Ensure:** Predicted bounding boxes and class labels**Backbone Feature Extraction:**

- 1: $F_1 \leftarrow \text{Conv}(I, k = 3, s = 2, p = 1)$
- 2: $F_1 \leftarrow \text{Conv}(F_1, k = 3, s = 2, p = 1)$
- 3: $F_1 \leftarrow \text{C2f}(F_1)$
- 4: $F_2 \leftarrow \text{Conv}(F_1, k = 3, s = 2, p = 1)$
- 5: $F_2 \leftarrow \text{C2f}(F_2)$
- 6: $F_3 \leftarrow \text{Conv}(F_2, k = 3, s = 2, p = 1)$
- 7: $F_3 \leftarrow \text{C2f}(F_3)$
- 8: $F_4 \leftarrow \text{Conv}(F_3, k = 3, s = 2, p = 1)$
- 9: $F_4 \leftarrow \text{C2f}(F_4)$
- 10: $F_4 \leftarrow \text{SPPF}(F_4)$

Neck with Attention:

- 11: $U_1 \leftarrow \text{Upsample}(F_4)$
- 12: $A_1 \leftarrow \text{Attention}(U_1)$
- 13: $M_1 \leftarrow \text{Concat}(A_1, F_3)$
- 14: $M_1 \leftarrow \text{C2f}(M_1)$
- 15: $U_2 \leftarrow \text{Upsample}(M_1)$
- 16: $A_2 \leftarrow \text{Attention}(U_2)$
- 17: $M_2 \leftarrow \text{Concat}(A_2, F_2)$
- 18: $M_2 \leftarrow \text{C2f}(M_2)$
- 19: $D_1 \leftarrow \text{Conv}(M_2, k = 3, s = 2, p = 1)$
- 20: $D_1 \leftarrow \text{Concat}(D_1, M_1)$
- 21: $D_1 \leftarrow \text{C2f}(D_1)$
- 22: $D_2 \leftarrow \text{Conv}(D_1, k = 3, s = 2, p = 1)$
- 23: $D_2 \leftarrow \text{Concat}(D_2, F_4)$
- 24: $D_2 \leftarrow \text{C2f}(D_2)$

Detection Head with Attention:

- 25: **for** $H \in \{M_2, D_1, D_2\}$ **do**
- 26: $H' \leftarrow \text{Attention}(H)$
- 27: $H' \leftarrow \text{Conv}(H', k = 3, s = 2, p = 1)$
- 28: $\text{BBox} \leftarrow \text{Conv2D}(H')$ {Bounding box regression}
- 29: $\text{Cls} \leftarrow \text{Conv2D}(H')$ {Classification}
- 30: **end for**
- 31: **if training then**
- 32: $\text{Loss}_{\text{bbox}} \leftarrow \text{ComputeLoss}(\text{BBox})$
- 33: $\text{Loss}_{\text{cls}} \leftarrow \text{ComputeLoss}(\text{Cls})$
- 34: $\text{TotalLoss} \leftarrow \text{Loss}_{\text{bbox}} + \text{Loss}_{\text{cls}}$
- 35: **else**
- 36: $\text{Predictions} \leftarrow \text{NMS}(\text{BBox}, \text{Cls})$
- 37: **return** Predictions
- 38: **end if**

After concatenating these features the final output of Shuffle attention is:

$$SA = CA \oplus SPA \quad (17)$$

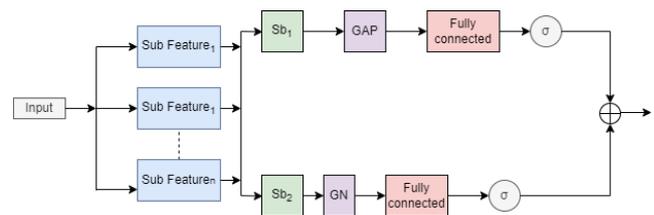


Figure 5: Shuffle attention block

4 Pre-processing and data augmentation of dataset

aquarium Dataset is used for performing experiment with different attention mechanism using YOLOv8. The aquarium Dataset, provided by Roboflow, consists of underwater images captured in controlled environments with limited variation in brightness. This homogeneity in image characteristics poses a challenge for the generalization of the trained model to other underwater images with different lighting conditions. To deal with this issue data augmentation technique is used to improve training dataset. The proposed work used fine-tuning of contrast and brightness so that different lightening levels are present with varying environment during training. The balance of the class of the aquarium dataset is shown in Table 1. To validate the work,

Table 1: Class balance for aquarium dataset

Class	Annotation
Fish	2669
Jellyfish	694
Penguin	516
Shark	354
Puffin	284
Stingray	184
Starfish	116

an experiment was also performed on the Brackish dataset and the class balance is shown in Table 2.

Table 2: Class balance for brackish dataset

Class name	Annotations
Crab	12,348
Smallfish	10,768
Starfish	7,912
Fish	3,352
Jellyfish	637
Shrimp	548

A popular data augmentation method in computer vision, HSV Augmentation modifies an image's Hue (H), Saturation (S), and Value (V) components to replicate different lighting and color conditions is used for augmentation [33]. Because underwater images frequently include uneven lighting and color distortion from light absorption and dispersion in water, this approach works especially well for underwater item detection. HSV augmentation improves the resilience and generalization of models to real-world underwater environments by randomly adjusting hue, saturation, and brightness during training. This helps models learn to distinguish objects under diverse visual appearances.

Since the dataset publisher did not give any predetermined training, validation, and test sets, we randomly divide the aquarium Dataset. More precisely, we assign 70% of the dataset to the training set, 20% to the test set, and 10% to the validation set.

5 Assessment parameters

Evaluation of proposed work is performed based on precision, recall, F1 score, mAP, Params(parameters), inference time, floating point operations (FLOPs) and frames per second (FPS). Precision, recall, F1 score and mAP are calculated based on True Positive(TP), True Negative(TN), False Positive(FP), False Negative(FN). Equation (18), eq (19), eq (20) presents formula for precision, recall, F1-score respectively.

$$Precision = \frac{TP}{TP + FP} \quad (18)$$

$$Recall = \frac{TP}{TP + FN} \quad (19)$$

$$F1 - score = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (20)$$

Category-wise Average Precision is calculated as equation (21)

$$AP(C_i) = (1/n) \times \left(\sum_{i=1}^n P_i \right) \quad (21)$$

where P_i is i^{th} the image of the C_i category and n is number of iterations.

Mean Average Precision is computed as equation (22)

$$mAP = (1/N) \times \sum_{i=1}^n AP(C_i) \quad (22)$$

where N is number of classes.

Params are the numbers of parameters involved during training and in this work parameters are calculated using Millions. The number of layers, neurons per layer, architectural complexity, and other variables all affect how many parameters a model has. A larger model size is typically associated with more parameters. In most cases, the larger the model, the better the performance of the model, but it also requires the use of additional data and processing power for training. The connection between computing cost and model complexity must be balanced in real-world applications.

The computational complexity of neural network models is frequently assessed using floating-point operations, which are a metric to evaluate computer or computing system performance. FLOPs show the number of floating-point operations per second of floating-point calculations, offering a vital measure of the model's speed and computational efficiency.

Frames Per Second (FPS) is an important statistic for object recognition, especially for real-time processing applications like interactive gaming, surveillance, and driverless cars. The responsiveness and efficacy of the detection system are strongly impacted by the frame rate (FPS), which shows how many frames a model can process in a second. The inference time of a trained object identification model is the amount of time it takes to process an input image and provide predictions.

6 Experimental setup

The experiment is conducted on PyTorch 2.1.2, utilizing CUDA 11.7 framework. The training is carried out on a single NVIDIA Tesla T4 GPU, which provides a balance between computational power and accessibility. The models are validated after training using the best checkpoint saved during the training process. Validation includes metrics such as precision, recall, mean Average Precision (mAP), and inference speed. Training hyper-parameters are shown in Table 3.

Table 3: Hyper-parameters model training

Parameter	Value
Image Size	640 × 640
Epochs	100
Optimizer	Stochastic Gradient Descent
Weight Decay	5×10^{-4}
Momentum	0.937
Initial Learning Rate	1×10^{-2}
Batch Size	16
Warmup Epochs	3
Warmup Momentum	0.8
Warmup Bias Learning Rate	0.1

Different software's are used during implementation of AB-YOLOv8 with Python 3.9. Pytorch and Tensor Board used to train the model and for visualization. Numpy and pandas are used for data pre-processing. The base YOLOv8 model is taken from Ultralytics and with it different attention module are used for proposed AB-YOLOv8.

7 Experimental results

In this section, detailed experimental results of the proposed work are reported. We train the AB-YOLOv8 model using training sets with input image size 1024, to compare the impact of varying input image sizes on the model's performance in the underwater item detection task. Table 4 shows the performance of different attention models combined with YOLOv8. YOLOv8 combined with ResCBAM, GAM, SA and ECA attention block and results are incorporated in this section. Table 4 presents the experimen-

tal results with respect to precision, recall, F1 score, and mAP. From Table 4 it is clear that ECA performs better than GAM, ResCBAM, SA when combined with YOLOv8. ECA performs 8% better than YOLOv8 and 6% better than GAM, SA, and ResCBAM.

Table 5 presents another set of AB-YOLOv8 experiment results showing evaluation of different metrics such as parameters, GLOPs, inference time and FPS. It is found from Table 5 in proposed AB-YOLOv8 ECA with YOLOv8 performs better than other techniques. ECA also achieved lowest inference time i.e. 7.7 ms where as other models attains 12.8ms, 8.7ms, 8.0ms inference time. AB-YOLOv8 when based on ResCBAM increased number of parameters almost 10M but when YOLOv8 is based on ECA its not increasing number of parameters as pooling operations are used to optimized the number of parameters. It is also clear that in all the models of AB-YOLOv8 have achieved similar FPS as original YOLOv8 but ECA based YOLOv8 attains 59FPS which is better than SA, GAM and ResCBAM. The aquarium dataset consists of seven categories species like fish, jellyfish, penguin, puffin, shark, starfish, stingray. The Table 6 presents class wise precision achieved by using AB-YOLOv8 models and YOLOv8. Bold results are showing best result achieved during experiments. Out of all AB-YOLOv8 models, ECA based model attains best result in most of the cases. In jellyfish class ResCBAM attains maximum mAP@50.

A small number of images are chosen at random for this paper's evaluation of the attention module's impact on the YOLOv8 model's accuracy in detecting fractures in a real-world marine environment exploration scenario. Figure 10 shows the prediction results of several AB-YOLOv8 models. As an object detection model, the AB-YOLOv8 model is essential to monitor and investigate the marine environment during research. It's crucial to remember, though, that every AB-YOLOv8 model worked flawlessly with tiny, tightly spaced items as well.

The **ablation experiment** shown in Table 7 indicates that the application of different attention mechanisms to the YOLOv8 model can result in considerable gains in mAP@50, recall, and precision; the most striking effect was shown by ECA (Efficient Channel Attention). The precision, recall, and mAP@50 of the base YOLOv8 model are 0.464, 0.305, and 0.328, respectively. The best overall results are obtained when ECA is applied at both the neck and the head (D+H), boosting precision to 0.561, recall to 0.387, and mAP@50 to 0.400. ECA consistently performs better than the other attention mechanisms, especially in terms of recollection and mAP, while SA, GAM, and ResCBAM show only modest gains, especially at the neck. D+H (YOLOv8 with ECA at both the neck and the head) is the best-performing configuration overall, suggesting that using ECA at both phases achieved substantial gains.

Among the evaluated models, ECA and the SA model achieve the highest overall F1-score of 0.29, shown in Figure 9 and Figure 8 respectively, while GAM lags slightly

Table 4: Experiment results of different attention models for aquarium dataset

Model	Precision (P)	Recall (R)	F1 Score	mAP@50	mAP@50-95
YOLOv8	0.464	0.305	0.367	0.328	0.150
YOLOv8+GAM	0.473	0.293	0.363	0.337	0.154
YOLOv8+ResCBAM	0.477	0.301	0.370	0.346	0.152
YOLOv8+SA	0.481	0.321	0.341	0.334	0.151
YOLOv8+ECA	0.561	0.387	0.458	0.400	0.194

Table 5: Experiment results of different attention models for aquarium Dataset

Model	Params(M)	GFLOPs	Inference(ms)	FPS
YOLOv8	43.67	164.37	7.7	60
YOLOv8+GAM	49.89	183.54	12.8	57
YOLOv8+ResCBAM	53.46	196.29	8.7	55
YOLOv8+SA	43.76	165.20	8.0	58
YOLOv8+ECA	43.54	165.34	7.7	59

Table 6: Category wise mAP@50 for different models for aquarium dataset

Category	YOLOv8	YOLOv8+SA	YOLOv8+GAM	YOLOv8+ResCBAM	YOLOv8+ECA
All	0.328	0.336	0.317	0.326	0.400
Fish	0.356	0.317	0.289	0.312	0.378
Jellyfish	0.614	0.561	0.566	0.682	0.656
Penguin	0.227	0.215	0.336	0.245	0.336
Puffin	0.114	0.183	0.105	0.182	0.249
Shark	0.283	0.281	0.207	0.291	0.295
Starfish	0.333	0.410	0.439	0.420	0.512
Stingray	0.372	0.152	0.274	0.150	0.381

Table 7: Ablation experiment for AB-YOLOv8 using aquarium dataset

Model	Precision	Recall	mAP@50
YOLOv8	0.464	0.305	0.328
A: YOLOv8+ SA at neck of YOLOv8	0.469	0.289	0.330
B: YOLOv8+ GAM at neck of YOLOv8	0.470	0.300	0.338
C: YOLOv8+ResCBAM at neck of YOLOv8	0.469	0.318	0.333
D: YOLOv8+ ECA at neck of YOLOv8	0.521	0.367	0.347
E: YOLOv8+ SA at head of YOLOv8	0.462	0.283	0.334
F: YOLOv8+ GAM at head of YOLOv8	0.476	0.291	0.342
G: YOLOv8+ResCBAM at head of YOLOv8	0.471	0.311	0.332
H: YOLOv8+ ECA at head of YOLOv8	0.541	0.367	0.381
A+E : YOLOv8 +SA	0.473	0.293	0.337
B+F: YOLOv8+ GAM	0.477	0.301	0.346
C+G: YOLOv8+ResCBAM	0.481	0.321	0.334
D+H: YOLOv8+ECA	0.561	0.387	0.400

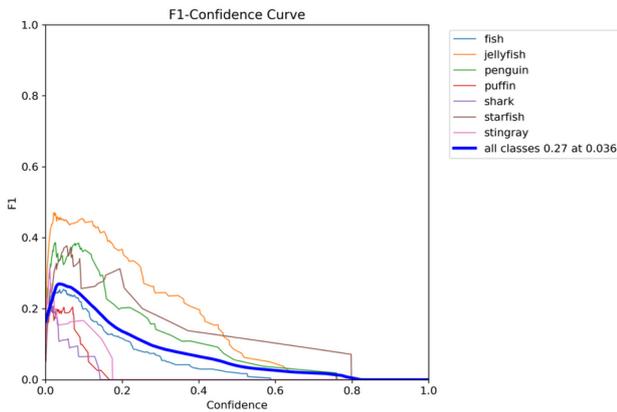


Figure 6: F1-confidence curve for YOLOv8 with GAM attention mechanism using aquarium dataset

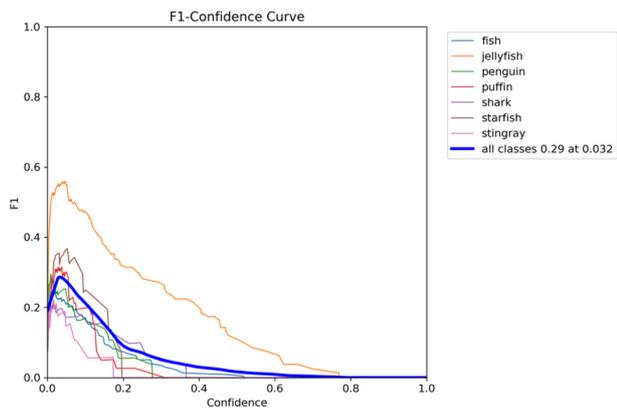


Figure 7: F1-confidence curve for YOLOv8 with CBAM attention mechanism using aquarium dataset

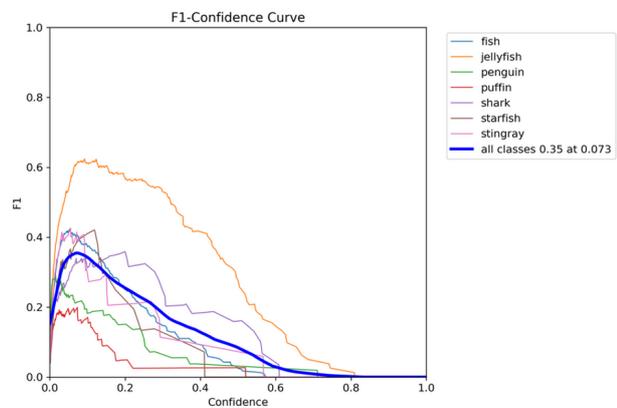


Figure 8: F1-confidence curve for YOLOv8 with SA attention mechanism using aquarium dataset

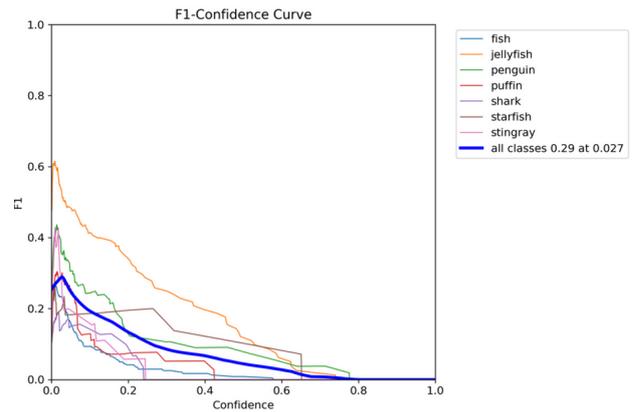


Figure 9: F1-confidence curve for YOLOv8 with ECA attention mechanism using aquarium dataset

at 0.27 as depicted in Figure 6. ECA stands out with the lowest optimal confidence threshold (0.027), offering superior early-stage detection and the smoothest confidence-F1 curve, making it ideal for robust predictions. CBAM and GAM contribute more toward improving per-class balance, with CBAM enhancing spatially diverse classes like puffin and starfish, and GAM excelling in classes with complex contextual dependencies like penguin, as shown in Figure 7. Although GAM does not reach peak F1 performance, it demonstrates the best inter-class balance. Overall, ECA provides the best trade-off between accuracy, stability, and efficiency, making it the most effective enhancement in this setting.

Statistical analysis The proposed work used an ANOVA test for performing statistical analysis. We have performed the same experiment 4 times and calculated mean, standard deviation, standard error and found YOLOv8+ECA performing better than others as shown in Table 8. Also assumed significance level as 5%. Table 9 shows that the p-value is 0.0004, which is much less than 0.05, so the result is significantly good. Moreover, the mean of YOLOv8 + ECA is maximum, so the performance of the ECA attention mechanism is performing well for the aquarium dataset.

8 Discussion

The AB-YOLOv8 compared with SSD and Faster R-CNN and results are shown in Figure 11, Figure 12 and Figure 13. With respect to precision, recall, and mAP, Faster R-CNN gives better results than YOLOv8, but after using the attention mechanism with YOLOv8, it is possible to outperform Faster R-CNN. Although Faster R-CNN is well-known for its high accuracy in object identification tasks, it has a number of drawbacks that limit its usefulness in real-time applications. Due to its two-stage detection architecture, which consists of a Region Proposal Network (RPN) followed by a classification and bounding box regression step, its main disadvantage is its lengthy inference time as shown in Figure 13. Because of this, it is computationally demanding

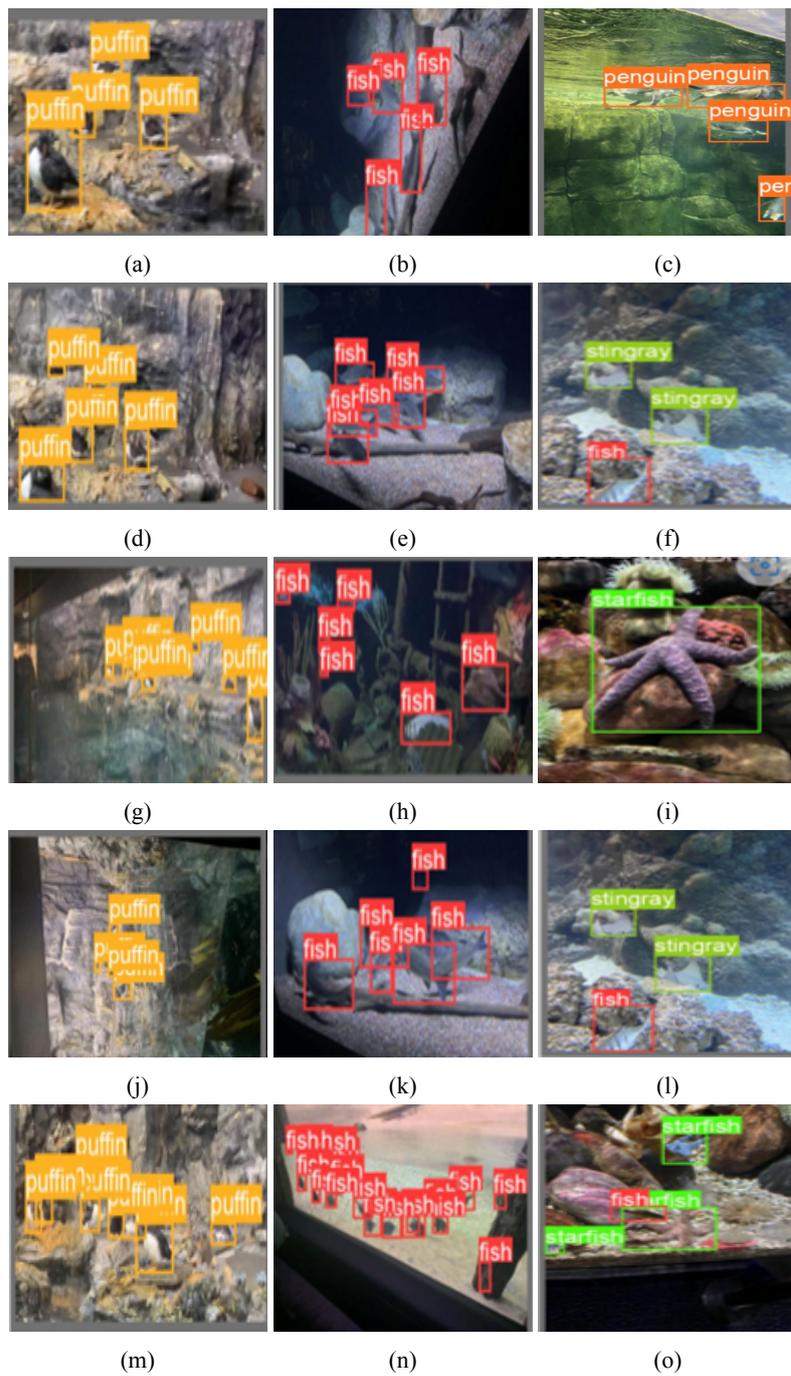


Figure 10: Sample images of object detection on aquarium dataset, (a-c) object detection by YOLOv8, (d-f) object detection by GAM, (g-i) object detection by ResCBAM, (j-l) object detection by SA, (m-o) object detection by ECA

Table 8: Statistical analysis based on aquarium dataset to calculate mean, standard deviation, standard error

Models	N	Mean	Std. Dev.	Std. Error
YOLOv8	4	43.88	1.781	0.7965
YOLOv8+SA	4	46.3333	1.2111	0.4944
YOLOv8+GAM	4	45.7833	1.3862	0.5659
YOLOv8+ResCBAM	4	46.85	1.1895	0.4856
YOLOv8+ECA	4	52.0167	5.1148	2.0881

Table 9: Statistical analysis based on aquarium dataset to calculate p-value

Source	Degrees of Freedom (DF)	Sum of Squares (SS)	Mean Square (MS)	F-Statistic	P-Value
Between Groups	4	211.1774	52.7943	7.5641	0.0004
Within Groups	24	167.5099	6.9796		
Total	28	378.6872			

and inappropriate for situations requiring quick decisions, such as autonomous driving or real-time video processing. SSD has significant limits even if it provides a decent balance between speed and accuracy. Its inability to detect small objects is a significant disadvantage, mainly due to the fact that it employs numerous feature maps with varying resolutions, which may result in the loss of fine features that are essential for localizing small objects. Furthermore, situations with dense backgrounds or complicated backdrops, where object boundaries are less clear, can be difficult for SSD to handle. Compared to two-stage detectors like Faster R-CNN, its accuracy is typically lower, but it has improved inference time to 25ms, depicted in Figure 13.

Figure 11 and Figure 12 clearly shows that GAM's performance on the AB-YOLOv8 model's on the aquarium dataset is poorer than other attention blocks. The one reason behind the poor performance of GAM is that it has an abundance of pooling layers. The ECA module can be deployed on devices with limited resources because it is computationally efficient and does not require dimensionality reduction or completely connected layers, which makes ECA more efficient and involves fewer parameters. Also, it is visible from Figure 12 ResCBAM and SA performed well with the YOLOv8 model. In order to improve feature representation and performance on a range of tasks, ResCBAM adds both channel and spatial attention, which enables the model to preferentially focus on the most informative channels and spatial regions of the feature maps. Another important issue is the result found on the aquarium Dataset, which consists only of 638 images, including validation, training, and testing images.

Based on how long it typically takes each object detection model to process a single image (measured in milliseconds), the inference time graph comparison shown in Figure 13. As can be seen from the graphic, Faster RCNN has the longest inference time—nearly 80 ms—which suggests that while it may attain competitive accuracy, its computational overhead renders it less appropriate for real-time applications. Even while SSD is faster than Faster RCNN, it still takes about 25 ms, which is more than the YOLO variations. The YOLOv8 and YOLOv8+ECA show noticeably higher inference efficiency than any of the other models that were assessed. Because YOLOv8 has the shortest inference time (around 7 ms), it is ideal for real-time systems.

The proposed work tested on another dataset to validate the performance of the proposed work. Brackish dataset used for the purpose of the experiment is shown in Table 10. It is found that for Brackish dataset ECA and ResCBAM

achieved 74% mAP@50. ECA does not perform dimensionality reduction so channel-wise features are intact and attains better result. ResCBAM efficiently determines the location and class of the objects by channel attention and spatial attention block. After inclusion of GAM and SA also achieved 2-4% improvement on mAP.

9 Conclusion

After the release of the YOLOv8 model by Ultralytics in 2023, researchers commenced utilizing it for object recognition in underwater images. Although the almost recent generation of the YOLO model, the YOLOv8 model, despite the fact that models performed admirably on the aquarium dataset, were unable to meet the good performance. We added four attention modules GAM, ResCBAM, SA and ECA to the YOLOv8 architecture, respectively, to improve the model's performance in order to overcome this constraint. Furthermore, we integrate ResBlock with CBAM to enhance the overall performance of the model. The proposed work with aquarium dataset achieved 40% maximum mAP@50 for ECA and ECA achieved 7.7 ms inference time with 59 FPS which is better than all other attention blocks. It is also notable that number of parameters not increased for ECA so finally, out of all attention block ECA performed better. Validation of the proposed work is checked on Brackish Dataset also. The results for Brackish dataset that shows for ResCBAM and ECA attention blocked achieved 74% mAP. YOLOv8 with ResCBAM and ECA achieved 8% better mAP than base YOLOv8 model.

Funding and conflicts of interest

Conflict of Interest: The authors did not receive funding and do not have any conflict of interest.

Data availability statements

The article used publicly the available aquarium dataset and link is as follows: aquarium dataset: <https://public.roboflow.com/object-detection/aquarium>
Brackish dataset: <https://public.roboflow.com/object-detection/brackish-underwater>

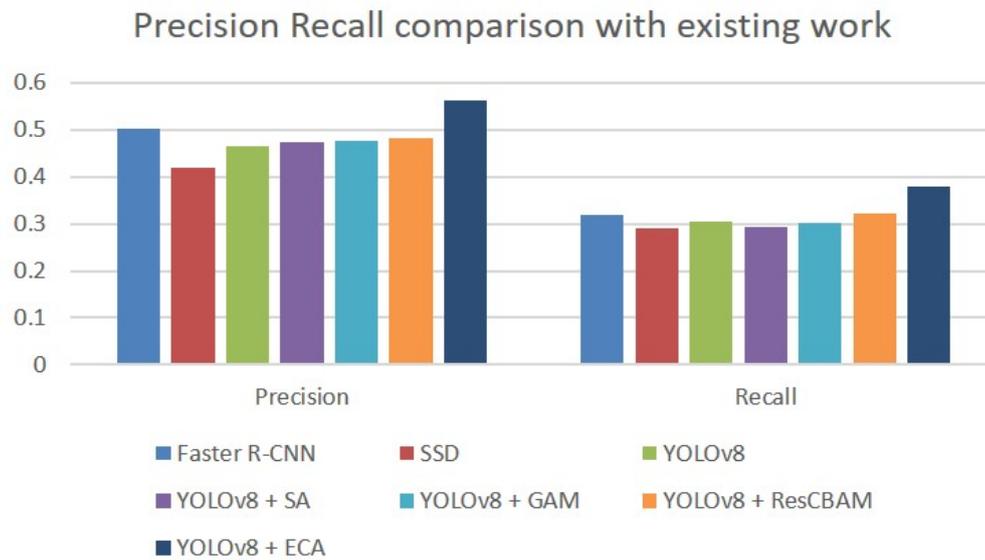


Figure 11: Precision and recall comparison for different models for aquarium dataset

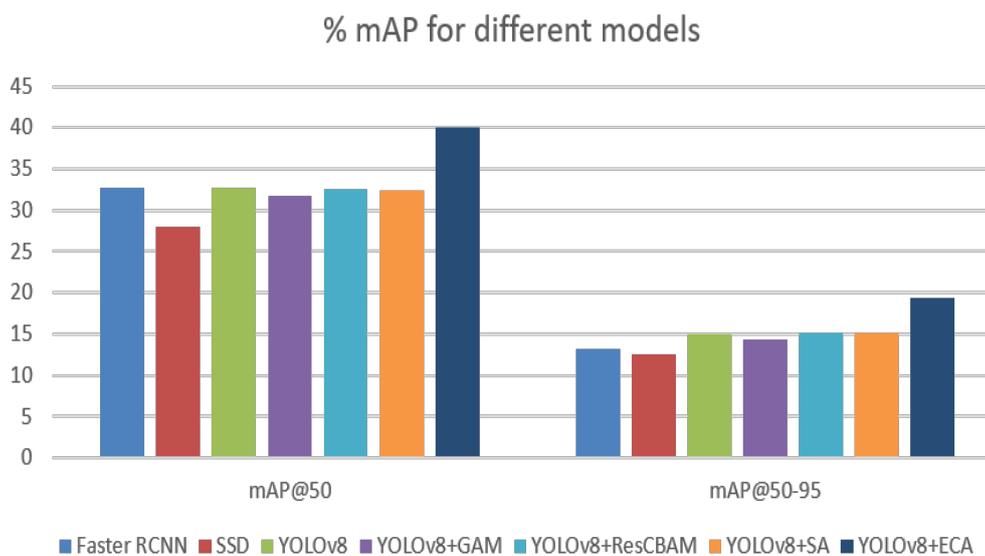


Figure 12: mAP comparison with different models for aquarium dataset

Table 10: Evaluation comparison between different models for Brackish Dataset

Network	Precision	Recall	mAP@50:95
SSD	41.19	35.02	30.71
Faster-RCNN	69.23	65.02	61.45
YOLOv8	92.29	91.04	68.21
YOLOv8+GAM	91.10	92.8	69.44
YOLOv8+SA	92.49	90.28	72.30
YOLOv8+ResCBAM	94.80	91.90	74.20
YOLOv8+ECA	95.01	90.90	74.31

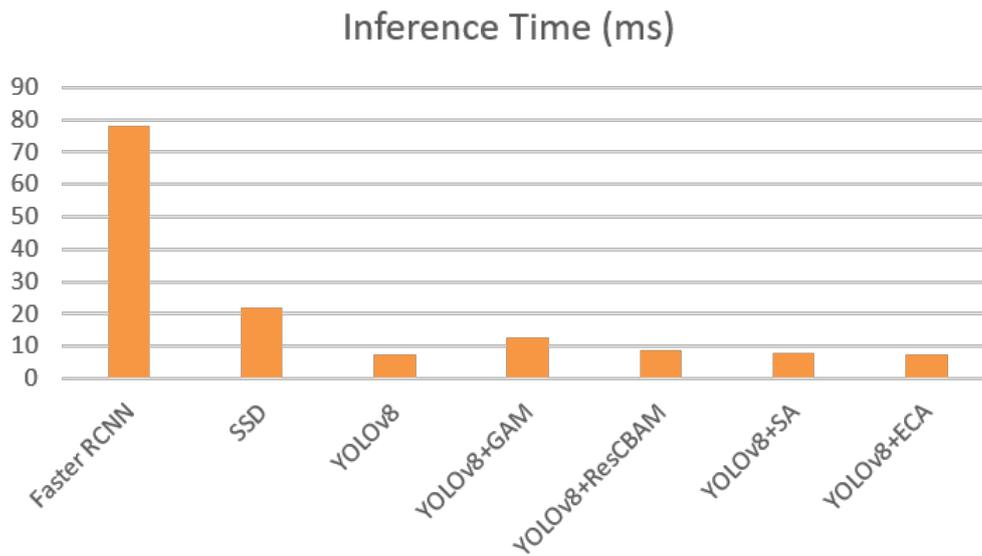


Figure 13: Inference Time comparison with different models

References

- [1] Pratima Sarkar, Sourav De, and Sandeep Gurung. A survey on underwater object detection. In *Intelligence Enabled Research: DoSIER 2021*, pages 91–104. Springer, 2022 <https://doi.org/10.1007/978-981-19-0489-9>.
- [2] Huimin Lu, Yujie Li, Yudong Zhang, Min Chen, Seichi Serikawa, and Hyungseop Kim. Underwater optical image processing: a comprehensive review. *Mobile networks and applications*, 22:1204–1211, 2017 <https://doi.org/10.1007/s11036-017-0863-4>.
- [3] Pratima Sarkar, Sandeep Gurung, and Sourav De. Underwater image segmentation using fuzzy-based contrast improvement and partition-based thresholding technique. In *Evolution in Computational Intelligence: Proceedings of the 9th International Conference on Frontiers in Intelligent Computing: Theory and Applications (FICTA 2021)*, pages 473–482. Springer, 2022 <https://doi.org/10.1007/978-981-16-6616-2>.
- [4] Dario Lodi Rizzini, Fabjan Kallasi, Fabio Oleari, and Stefano Caselli. Investigation of vision-based underwater object detection with multiple datasets. *International Journal of Advanced Robotic Systems*, 12(6):77, 2015 <https://doi.org/10.5772/60526>.
- [5] Pratima Sarkar, Sourav De, Sandeep Gurung, and Prasenjit Dey. Uice-mirnet guided image enhancement for underwater object detection. *Scientific Reports*, 14(1):22448, 2024 <https://doi.org/10.1038/s41598-024-73243-9>.
- [6] Yan Zhai. River ship monitoring based on improved deep-sort algorithm. *Informatica*, 48(9), 2024 <https://doi.org/10.31449/inf.v48i9.5886>.
- [7] Md Moniruzzaman, Syed Mohammed Shamsul Islam, Mohammed Bennamoun, and Paul Lavery. Deep learning on underwater marine object detection: A survey. In *Advanced Concepts for Intelligent Vision Systems: 18th International Conference, ACIVS 2017, Antwerp, Belgium, September 18-21, 2017, Proceedings 18*, pages 150–160. Springer, 2017 <https://doi.org/10.1007/978-3-319-70353-4>.
- [8] Pratima Sarkar, Sourav De, and Sandeep Gurung. Fish detection from underwater images using yolo and its challenges. In *Doctoral Symposium on intelligence enabled research*, pages 149–159. Springer, 2022 <https://doi.org/10.1007/978-981-99-1472-2>.
- [9] Long Chen, Zhihua Liu, Lei Tong, Zheheng Jiang, Shengke Wang, Junyu Dong, and Huiyu Zhou. Underwater object detection using invert multi-class adaboost with deep learning. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020 <https://doi.org/10.1109/IJCNN48605.2020.9207506>.
- [10] Jian Zhang, Jinshuai Zhang, Kexin Zhou, Yonghui Zhang, Hongda Chen, and Xinyue Yan. An improved yolov5-based underwater object-detection framework. *Sensors*, 23(7):3693, 2023 <https://doi.org/10.3390/s23073693>.
- [11] Pratick Gupta, Pratima Sarkar, Bijoyeta Roy, and Shivam Kumar. Fish classification using cnn and logistic regression from underwater images. In

- International Conference on Advanced Computational and Communication Paradigms*, pages 415–424. Springer, 2023
<https://doi.org/10.1007/978-981-99-4284-8>.
- [12] Wang Zhiqiang and Liu Jun. A review of object detection based on convolutional neural network. In *2017 36th Chinese control conference (CCC)*, pages 11104–11109. IEEE, 2017
<https://doi.org/10.23919/ChiCC.2017.8029130>.
- [13] Mupparaju Sohan, Thotakura Sai Ram, Rami Reddy, and Ch Venkata. A review on yolov8 and its advancements. In *International Conference on Data Intelligence and Cognitive Informatics*, pages 529–545. Springer, 2024
<https://doi.org/10.1007/978-981-99-7962-2>.
- [14] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net: Efficient channel attention for deep convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11534–11542, 2020
<https://arxiv.org/pdf/1910.03151v3>.
- [15] Qing-Long Zhang and Yu-Bin Yang. Sa-net: Shuffle attention for deep convolutional neural networks. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2235–2239. IEEE, 2021
 Doi: 10.1109/ICASSP39728.2021.9414568/.
- [16] Yichao Liu, Zongru Shao, and Nico Hoffmann. Global attention mechanism: Retain information to enhance channel-spatial interactions. *arXiv preprint arXiv:2112.05561*, 2021
<https://arxiv.org/abs/2112.05561>.
- [17] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018
<https://doi.org/10.48550/arXiv.1807.06521>.
- [18] Chun-Tse Chien, Rui-Yang Ju, Kuang-Yi Chou, Chien-Sheng Lin, and Jen-Shiun Chiang. Yolov8-am: Yolov8 with attention mechanisms for pediatric wrist fracture detection. *arXiv preprint arXiv:2402.09329*, 2, 2024
<https://doi.org/10.1109/ACCESS.2025.3549839>.
- [19] Arindam Chaudhuri. Hierarchical modified fast r-cnn for object detection. *Informatica*, 45(7), 2021
<https://doi.org/10.31449/inf.v45i7.3732>.
- [20] Weiwen Chen, Tingting Zhuang, Yuanfang Zhang, Teng Mei, and Xiaoyu Tang. Yolo-uod: An underwater small object detector via improved efficient layer aggregation network. *IET Image Processing*, 2024
<https://doi.org/10.1049/ipr2.13112>.
- [21] Shijia Zhao, Jiachun Zheng, Shidan Sun, and Lei Zhang. An improved yolo algorithm for fast and accurate underwater object detection. *Symmetry*, 14(8):1669, 2022
<https://doi.org/10.3390/s23073693>.
- [22] Kun Liu, Lei Peng, and Shanran Tang. Underwater object detection using tc-yolo with attention mechanisms. *Sensors*, 23(5):2567, 2023
<https://doi.org/10.3390/s23052567>.
- [23] Ali M Reza. Realization of the contrast limited adaptive histogram equalization (clahe) for real-time image enhancement. *Journal of VLSI signal processing systems for signal, image and video technology*, 38:35–44, 2004
<https://doi.org/10.1023/B:VLSI.0000028532.53893.82>.
- [24] Risheng Liu, Xin Fan, Ming Zhu, Minjun Hou, and Zhongxuan Luo. Real-world underwater enhancement: challenges, benchmarks, and solutions. *arXiv preprint arXiv:1901.05320*, 2019
<https://doi.org/10.48550/arXiv.1901.05320>.
- [25] Xin Shen, Xudong Sun, Huibing Wang, and Xianping Fu. Multi-dimensional, multi-functional and multi-level attention in yolo for underwater object detection. *Neural computing and applications*, 35(27):19935–19960, 2023
<https://doi.org/10.1007/s00521-023-08781-w>.
- [26] Jie Chen and Meng Joo Er. Dynamic yolo for small underwater object detection. *Artificial Intelligence Review*, 57(7):1–23, 2024
<https://doi.org/10.1007/s10462-024-10788-1>.
- [27] Yuyi Yang, Liang Chen, Jian Zhang, Lingchun Long, and Zhenfei Wang. Ugc-yolo: underwater environment object detection based on yolo with a global context block. *Journal of Ocean University of China*, 22(3):665–674, 2023
<https://doi.org/10.1007/s11802-023-5296-z>.
- [28] Xiaohan Wang, Xiaoyue Jiang, Zhaoqiang Xia, and Xiaoyi Feng. Underwater object detection based on enhanced yolo. In *2022 International Conference on Image Processing and Media Computing (ICIPMC)*, pages 17–21. IEEE, 2022
<https://doi.org/10.1109/ICIPMC55686.2022.00012>.
- [29] AF Ayob, K Khairuddin, YM Mustafah, AR Salisa, and K Kadir. Analysis of pruned neural networks (mobilenetv2-yolo v2) for underwater object detection. In *Proceedings of the 11th National Technical Seminar on Unmanned System Technology 2019: NUSYS'19*, pages 87–98. Springer, 2021
<https://doi.org/10.1007/978-981-15-5281-6>.
- [30] Yue Li, Xueting Zhang, and Zhangyi Shen. Yolo-submarine cable: an improved yolo-v3 network for

object detection on submarine cable images. *Journal of Marine Science and Engineering*, 10(8):1143, 2022
<https://doi.org/10.3390/jmse10081143>.

- [31] Ali Farhadi and Joseph Redmon. Yolov3: An incremental improvement. In *Computer vision and pattern recognition*, volume 1804, pages 1–6. Springer Berlin/Heidelberg, Germany, 2018
<https://doi.org/10.48550/arXiv.1804.02767>.
- [32] Pratima Sarkar, Sourav De, and Sandeep Gurung. U-yolov3: A model focused on underwater object detection. *Informatica*, 49(6), 2025
<https://doi.org/10.31449/inf.v49i6.6642>.
- [33] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020
<https://doi.org/10.48550/arXiv.2004.10934>.

Blurred Face Image Authentication for Enterprise Attendance Using Adaptive Light Adjustment and GAN-CNN Architecture

Lingyun Yang

Sichuan Vocational and Technical College, Suining, Sichuan 629000, China

Email: yly_yang@outlook.com

Keywords: attendance, blurred face, face recognition, identity authentication

Received: April 26, 2025

This paper combined the adaptive light adjustment algorithm and the generative adversarial network (GAN) deblurring algorithm with a convolutional neural network (CNN) algorithm for blurred face image recognition. First, the adaptive light adjustment algorithm and the GAN algorithm were used to perform deblurring operations on blurred face images, and then the CNN algorithm was used to recognize them. Then, simulation experiments were conducted. In the experiments, the adaptive light adjustment-combined GAN deblurring algorithm was compared with the Gaussian filter method and traditional GAN algorithm. The proposed face authentication algorithm was compared with the support vector machine and traditional CNN algorithms. The results showed that the adaptive light adjustment-combined GAN algorithm could effectively deblur the face image, with a peak signal to noise ratio of 32.08 and a deblurring time of 0.29 s. Moreover, the proposed face authentication algorithm could effectively recognize the identity of the blurred face image, with a precision of 0.987, a recall rate of 0.986, and an F value of 0.986, and it consumed 0.31 s for recognition.

Povzetek: Ta članek predstavlja algoritem za avtentikacijo zamegljenih obrazov za evidenco delovnega časa, ki združuje prilagoditev svetlobe, izostritev s GAN (Generative Adversarial Network) in prepoznavanje s CNN.

1 Introduction

In today's rapidly developing digital age, the enterprise attendance management as an important part of human resource management is experiencing a significant transformation, shifting from traditional manual processes to intelligent and automated systems [1, 2]. Traditional attendance methods include paper-based attendance and card-based attendance. Although they can meet the basic requirements of enterprises to a certain extent, they have defects such as easy to forge, easy to lose, and inconvenient to manage. With the rapid development of biometric technology, facial recognition technology has become a solution due to its unique characteristics of non-contact, high accuracy, and difficulty in replication [3]. But in the process of actual use, face recognition technology is also faced with many challenges. Especially under the impact of complex and changeable attendance environment, facial images can become blurred due to factors such as changes in light, interference from obstructions, shooting angles, or resolution [4]. Processing blurry face images to reduce their blurriness or extract key facial features can effectively improve the accuracy of face authentication. Related works are reviewed in Table 1. Those studies have all analyzed aspects such as faces and fuzzy classification. Some of them focused on face recognition, while others placed the research emphasis on fuzzy classification. This paper, however, focuses on the recognition of blurred faces for enterprise attendance purposes. The innovation of this paper lies in combining the clarification of blurred faces with face image

recognition. In the process of clarifying blurred faces, an adaptive light adjustment algorithm was used to reduce the influence of environmental light in the image, and the generative adversarial network (GAN) algorithm was utilized to clarify the blurred faces. Finally, the advantages of the convolutional neural network (CNN) algorithm in image recognition were exploited to identify face images. Moreover, during the CNN training process, in order to avoid the rigidity brought about by manually annotating features of face images, triplet samples were adopted for training. This paper combined the adaptive light adjustment algorithm for blurred face images and the GAN algorithm for deblurring with the CNN algorithm for the recognition of blurred face images. Then, simulation experiments were carried out. The contribution of this paper lies in combining the adaptive light adjustment algorithm and the GAN algorithm to deblur blurred face images and using the CNN algorithm trained with triplet samples to recognize face images, providing an effective reference for the security and efficiency of enterprise attendance. The limitation of this paper is that the scale of face samples used in the test is limited, resulting in insufficient generalization of the experimental results. The future research direction is to expand the sample scale to improve the generalization of the face recognition algorithm.

Table 1: Related works

Author	Research content	Research results
Balovsyak et al. [5]	They implemented a face recognition in system by using the Viola-Jones method and fuzzy logic.	The results showed that the method can improve the accuracy of face and mask recognition.
Khan et al. [6]	They proposed a layered classifier based on fuzzy rules for the forensic field.	The results showed that this method can play an important role in the forensic field.
Zhang et al. [7]	They proposed a new method of face recognition that incorporates fuzzy set theory.	The experimental results verified that this method can effectively identify blurred faces.

2 The algorithm for recognizing blurred face images

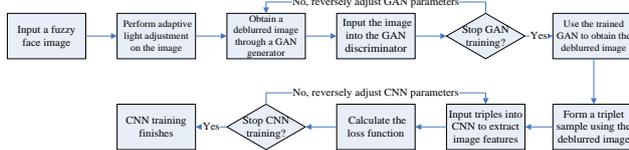


Figure 1: The algorithm training flow of blurred face image recognition

When human face is used as a biometric authentication for attendance [8], it has the advantages of non-contact, high accuracy, and difficult to copy [9]. To improve the accuracy of face recognition authentication, this paper first performs the adaptive adjustment processing on the blurred face image, then uses a GAN to deblur it, and finally uses a CNN to extract face features for identity authentication [10]. The specific flow is as follows.

① A blurred face image with a specified size is input and processed by light adaptive adjustment, including image defogging and Gamma correction:

$$\begin{cases} J(x) = \frac{I(x) - A}{\max(t(x), t_0)} + A \\ f(x) = (N(x))^{2K(x) + K'(x) - 0.5} \end{cases}, (1)$$

where $I(x)$ is the pixel of the original image, A is the global atmospheric light component, which is an empirical formula, $t(x)$ is the transmission function of scattered light, which is an empirical formula, t_0 is the minimum of $t(x)$, $J(x)$ is the image after defogging, $N(x)$ is the image after the normalization of $J(x)$ pixel values, $K(x)$

is the mean luminance of $J(x)$, and $K'(x)$ is the standard deviation of the luminance of $J(x)$ [11].

② The face image processed by adaptive light adjustment is input into the GAN generator to generate a deblurred image [12], and the size of the generated image is the same as the input size.

③ The generated image is taken as the negative sample, and the corresponding clear original face image is taken as the positive sample. The two samples are input to the GAN discriminator for forward calculation. Finally, the judgment result of whether the sample is the original clear face image is output to the softmax layer of the discriminator.

④ Whether the training of the GAN is terminated is determined. The termination condition is that the training times reach the threshold or the loss function converges to stability. If the training termination condition is met, proceed to the next step; otherwise, the loss function is used to reversely adjust the weight parameters in the GAN generator and the discriminator [13]. The loss function is:

$$\begin{cases} L_{cont} = \frac{\sum_{x=1}^W \sum_{y=1}^H (I_{x,y}^s - G(I_{x,y}^b, w_i))^2}{WH} + \beta \sum_{i=1}^n w_i^2, (2) \\ L_{adv} = E_{I^s} [\log D(I^s)] + E_{I^b} [\log(1 - D(G(I^b)))] \end{cases}$$

where L_{cont} refers to a content loss, L_{adv} refers to an adversarial loss, (x, y) refers to image pixels, W and H are the width and height of the image, $I_{x,y}^s$ is the (x, y) pixel in the original clear face image, $I_{x,y}^b$ is the (x, y) pixel in the original blurred face image, $G(I_{x,y}^b, w_i)$ is the (x, y) pixel of the simulated clear image processed by the GAN generator, w_i is the weight parameter in the GAN generator, β is the regularization term coefficient of $I_{x,y}^b$, E_{I^s} and E_{I^b} are the mathematical expected distribution of the original clear and blurred face image, $D(I^s)$ is the probability that the original clear face image is judged as a clear image by the GAN discriminator, and $D(G(I^b))$ is the probability of $G(I^b)$ being recognized as a clear image. After the GAN is trained, it is applied to the subsequent training of the CNN.

⑤ The trained GAN is employed to deblur the blurred face image in the training set, and then a triplet sample is constructed using the deblurred image [14]. The triplet sample consists of two positive samples and one negative sample. The two positive samples are the face images of the same person, and the negative samples are the face images of the other person.

⑥ The triplet sample is input into the CNN. Convolutional features are extracted and compressed for each face image in the triplet through convolutional layers and pooling layers.

⑦ The loss function of the CNN for the triplet sample is calculated:

$$loss = \sum_i^N \left[\left\| f(x_i^a) - f(x_i^p) \right\|_2^2 - \left\| f(x_i^a) - f(x_i^n) \right\|_2^2 + \alpha \right], \quad (3)$$

where $f()$ refers to the CNN, x_i^a and x_i^p are two positive samples in the i -th triple, x_i^n is the negative sample in the i -th triple, and α is the margin of the feature vector between the x_i^a, x_i^p group and the x_i^a, x_i^n group.

⑧ Whether CNN training can be terminated is determined. The training termination condition is also that the number of training reaches the prescribed number or the triplet loss function converges to stability. If the condition is satisfied, the training of CNN is completed; otherwise, the parameters in CNN are reversely adjusted according to the triplet loss. When the above algorithm is trained, and it is applied to the blurred face identity authentication for enterprise attendance, the basic process is as follows. First, the blurred face image of an employee is processed through the adaptive light adjustment algorithm, and then it is input into the GAN algorithm to generate a clear face image. After that, the clear image is input into the CNN algorithm. The face features are obtained from the last convolutional layer of the CNN algorithm. Finally, the face features extracted by the CNN algorithm is compared with the face features stored in the database to achieve identity authentication. The face features stored in the database are also extracted through the CNN algorithm.

3 Simulation experiment

3.1 Experimental environment

The simulation experiment was carried out in the laboratory server, which was configured as Windows11 system, I7 processor, 16 G memory, and RTX4060 graphics card. Python was used for programming.

3.2 Experimental data

The open-source Columbia University Public Figures Face Database (m6z.cn/5DIIR9) was used as the dataset for the simulation.

3.3 Experiment setup

The relevant parameters of the face identity authentication algorithm used for attendance are shown in Table 2. The GAN and CNN algorithms were both trained 500 times, with a learning rate of 0.02. The group convolution and squeeze-and-excitation (Group-SE) module in the GAN generator is a composite structure, which is composed of the standard convolutional layer, the grouped convolutional layer, the nonlinear activation layer, the SE layer, the standard convolutional layer, and the feature fusion layer in sequence.

Table 2: Relevant parameters of the proposed algorithm.

	Structure	Parameter setting	Structure	Parameter setting
GAN generator	Input layer	200×250	The first convolutional layer	Eight 7×7 convolution kernels, a step size of 2, and a sigmoid activation function
	Group-SE module	12	Deconvolution layer	Eight 7×7 convolution kernels, tanh activation function
GAN discriminator	Input layer	200×250	Convolutional layer 1	Eight 3×3 convolution kernels, a step size of 2, and a sigmoid activation function
	Convolution layer 2	16 3×3 convolution kernels, a step size of 2, and a sigmoid activation function	Pooling layer	3×3 mean pooling box, a step size of 2
	Fully connected layer	Softmax function	Output layer	1 node
CNN for extracting facial features	Input layer	200×250	Convolutional layer 1	32 3×3 convolution kernels, a step size of 2, and a sigmoid activation function
	Convolutional layer 2	64 3×3 convolution kernels, a step size of 2, and a sigmoid activation function	Pooling layer 1	3×3 mean pooling box, a step size of 2
	Convolutional layer 3	64 3×3 convolution kernels, a step size of 2, and a sigmoid activation function	Pooling layer 2	3×3 mean pooling box, a step size of 2
	Output layer	128 nodes		

The SVM and traditional CNN algorithms were used for verification and comparison. The principle of the SVM algorithm for face identity authentication was to train using the training samples according to the identity corresponding to the face and then use it to classify the identity of the input face image to achieve identity authentication. The relevant parameters of the SVM

algorithm are as follows. A linear kernel function was used, and the penalty parameter was set to 1. The principle of the traditional CNN algorithm for face authentication was to store feature vectors extracted from face images by the CNN algorithm, extract feature vectors from an image to be authenticated using the CNN algorithm, and compare them with the stored feature vectors. The relevant parameters of the traditional CNN algorithm were consistent with the CNN part of the proposed algorithm.

In addition, an ablation experiment was conducted on the GAN algorithm. The GAN algorithm without the Group-SE module was tested and compared with the complete GAN algorithm.

3.4 Evaluation criteria

In the algorithms employed, the GAN algorithm was used to deblur blurred face images, and the evaluation criteria for the deblurring effect are:

$$\left\{ \begin{array}{l} PSNR = 10 \lg \left(\frac{(2^n - 1)^2}{MSE} \right) \\ MSE = \frac{\sum_{i=1}^H \sum_{j=1}^W (X_{i,j} - Y_{i,j})^2}{H \times W} \end{array} \right., (4)$$

where *PSNR* is the image distortion index (the larger the value, the smaller the distortion), $X_{i,j}$ and $Y_{i,j}$ are the deblurred image pixel and original clear image pixel, and *MSE* is the average error between $X_{i,j}$ and $Y_{i,j}$.

The performance of the proposed algorithm for face identity authentication was measured by the commonly used precision, recall rate, and F value. The equations are:

$$\left\{ \begin{array}{l} P = \frac{TP}{TP + FN} \\ R = \frac{TP}{TP + FP} \\ F = \frac{2 \cdot P \cdot R}{P + R} \end{array} \right., (5)$$

where *P* denotes the precision, *R* denotes the recall rate, *F* is the combined value of the precision and recall rate, *TP* is the number of true positive samples, *FP* is the number of false positive samples, *FN* is the number of false negative samples, and *TN* is the number of true negative samples.

3.5 Test results

The partial deblurring results of three deblurring algorithms for blurred face images are shown in Figure 1. It can be seen that the face image processed by the adaptive light adjustment-combined GAN algorithm was the clearest, followed by the face image processed by the traditional GAN algorithm. Although the face image processed by Gaussian filtering was clearer than the original image, the blurriness was still visible to the naked eye. The objective deblurring effects of the three deblurring algorithms are shown in Table 3. It can be seen that the adaptive light adjustment-combined GAN

algorithm had the largest *PSNR* and the shortest deblurring time, followed by the traditional GAN algorithm with a medium *PSNR* and deblurring time, and the Gaussian filter algorithm had the smallest *PSNR* and the longest deblurring time.

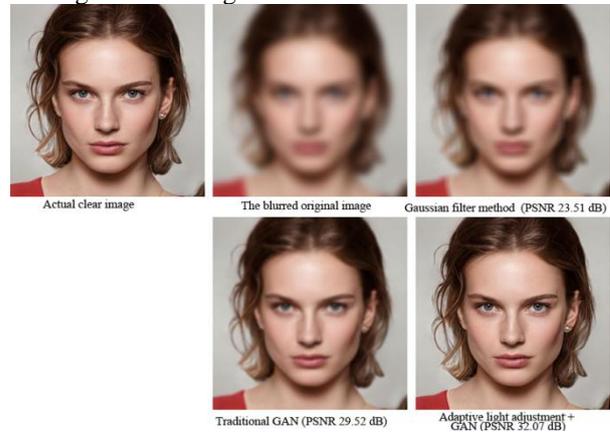


Figure 1: Partially blurred face images and images deblurred by three deblurring algorithms.

Table 3: The deblurring effect of the algorithms.

	Gaussian filtering method	Traditional GAN	Adaptive light adjustment +GAN
PSNR /dB	23.45	29.54*	32.08*+
Deblurring time/s	1.25	0.97*	0.29*+

Note: * indicates that the p value in the difference with the Gaussian filtering method is smaller than 0.05; + indicates that the p value in the difference with the traditional GAN algorithm is smaller than 0.05.

The recognition performance of the three authentication algorithms is shown in Table 4 and Figure 2. It can be seen that for blurred faces, the proposed algorithm combining adaptive light adjustment, GAN, and CNN had the highest recognition accuracy and the shortest recognition time. The accuracy and recognition time of the traditional CNN algorithm for blurred faces were in the middle among the three algorithms, while the SVM algorithm exhibited the worst recognition accuracy and the longest recognition time. The receiver operator characteristic (ROC) curves also showed that the proposed algorithm had the best human face recognition performance, followed by the traditional CNN algorithm, and the SVM algorithm was the worst.

Table 4: Recognition performance of three authentication algorithms.

	Precision	Recall rate	F value	Recognition time consumption/s
SVM algorithm	0.712	0.711	0.712	2.36

Traditional CNN algorithm	0.875*	0.876*	0.875*	1.12*
Adaptive light adjustment +GAN+CNN	0.987*+	0.986*+	0.986*+	0.31*+

Note: * indicates the p value in the difference with the SVM algorithm is smaller than 0.05; + indicates the p value in the difference with the traditional CNN algorithm is smaller than 0.05.

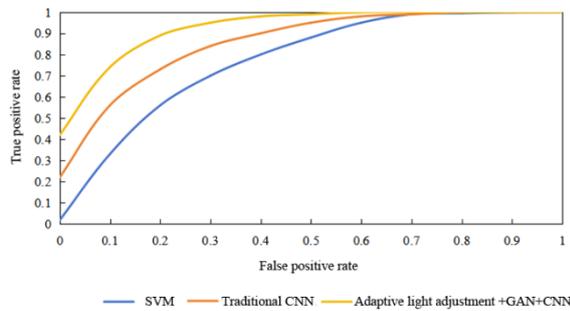


Figure 2: The ROC curves of three algorithms.

An ablation experiment was conducted on the GAN algorithm to verify the contribution of the Group-SE module in the GAN algorithm. The results are shown in Table 5. After removing the Group-SE module, although there was no significant change in the deblurring time, the deblurring effect of the GAN algorithm on blurred face images was significantly reduced, which verified that the Group-SE module made a significant contribution in the GAN algorithm.

Table 5: The ablation experiment for the GAN algorithm.

	The GAN algorithm without Group-SE	GAN algorithm
PSNR /dB	22.56	32.08*
Deblurring time/s	0.28	0.29

Note: * indicates the p value in the difference is smaller than 0.05.

4 Discussion

In today’s digital age, attendance management, which is part of human resource management, has gradually shifted from traditional manual execution to digital and automated execution, greatly enhancing the efficiency of human resource management. The attendance management in the daily operations of enterprises needs to be efficient and accurate enough, and the traditional attendance methods are gradually failing to meet the requirements. Face recognition technology has the advantages of being non-contact, highly efficient, and difficult to replicate, but it still faces problems such as light interference, obstructions, and low camera resolution

in practical applications, all of which can cause blurring of images and reduce the accuracy of recognition and authentication. This paper used the adaptive light adjustment algorithm combined with GAN algorithm to deblur the blurred image and then used the triplet sample constructed by the deblurred image to train the CNN algorithm.

When the blurred face recognition algorithm proposed in this paper is applied to enterprise attendance management, its overall process can be divided into the entry and verification of employee face information. When entering the image, a camera is first used to collect the employee’s face image, the adaptive light adjustment algorithm and GAN algorithm were used to deblur the collected face image, the trained CNN algorithm was used to extract the face features from the face image, and the extracted face features are stored. During verification, the system similarly captures the employee’s facial image via the camera, applies an adaptive light adjustment algorithm and GAN-based deblurring algorithm for preprocessing. Subsequently, a CNN is employed to extract facial features. The extracted features are then compared against the stored template. If the discrepancy between them is below a predefined threshold, authentication is granted; otherwise, it is rejected.

In the process of face identity authentication using the above-mentioned algorithm, the accuracy of the blurred face recognition algorithm is of crucial importance. In this paper, the proposed recognition algorithm was tested in simulation experiments. First, the deblurring effect of the adaptive light adjustment algorithm combined with the GAN algorithm was tested, and then the face recognition performance of the CNN algorithm was tested. Compared with the Gaussian filter and the traditional GAN algorithm, the adaptive light adjustment algorithm combined with the GAN algorithm exhibited the best deblurring effect. The reason is that the Gaussian filter can effectively deal with the noise that follows the normal distribution in the image, but in the blurred image, the factors that cause the face to blur are not all the noise that conforms to the normal distribution. The traditional GAN algorithm generates a “clear image” based on the input “blurred image”. Although the GAN algorithm adjusts the parameters within the algorithm by using the confrontation between the generator and the discriminator during training, the “clear image” is always generated based on the “blurred image” and is affected by the “blurred image”. The deblurring algorithm used in this paper first preprocessed the blurred image using an adaptive light adjustment algorithm and then generated a “clear face image” through the GAN algorithm, so the deblurring effect was better. In the face recognition and authentication part, this algorithm used triplet samples to train the CNN algorithm. There is no need to annotate the identity category of the face image; instead, it only needs to extract the face features in the image to obtain more comprehensive features, so its recognition performance was the best.

5 Conclusions

This paper combined the adaptive light adjustment algorithm and the GAN deblurring algorithm with a CNN algorithm for blurred face image recognition and carried out simulation experiments. In the experiments, the adaptive light adjustment-combined GAN deblurring algorithm was compared with the Gaussian filter method and the traditional GAN algorithm. Moreover, the proposed blurred face authentication algorithm was compared with the SVM and traditional CNN algorithms. The adaptive light adjustment-combined GAN algorithm had the best deblurring effect on the face. In terms of objective deblurring indicators, the adaptive light adjustment-combined GAN algorithm had the largest PSNR and the shortest deblurring time. The proposed authentication algorithm had the highest recognition accuracy and the shortest recognition time for blurred face images.

References

- [1] Yang Y (2020). Research on brush face payment system based on internet artificial intelligence. *Journal of Intelligent & Fuzzy Systems*, 38(1), pp. 21-28. <http://doi.org/10.3233/JIFS-179376>
- [2] Patel A, Jana S, Mahanta J (2024). Construction of similarity measure for intuitionistic fuzzy sets and its application in face recognition and software quality evaluation. *Expert Systems with Applications*, 237(PartB), pp. 22. <http://doi.org/10.1016/j.eswa.2023.121491>
- [3] Ghosh M, Sing J K (2023). Interval type-2 fuzzy set induced fuzzy rank-level fusion for face recognition. *Applied Soft Computing*, 2023, pp. 145.
- [4] Tao X, Pan D (2023). Face recognition based on scale invariant feature transform and fuzzy reasoning. *Internet Technology Letters*, 6(5), pp. e346.1-e346.5. <http://doi.org/10.1002/itl2.346>
- [5] Balovsyak S, Derevyanchuk O, Kovalchuk V, Kravchenko H, Kozhokar M (2024). Face mask recognition by the viola-jones method using fuzzy logic. *International Journal of Image, Graphics and Signal Processing*, 16(3), pp. 13. <http://doi.org/10.5815/ijigsp.2024.03.04>
- [6] Khan M A, Jalal A S (2019). A fuzzy rule based multimodal framework for face sketch-to-photo retrieval. *Expert Systems with Application*, 134(NOV), pp. 138-152. <http://doi.org/10.1016/j.eswa.2019.05.040>
- [7] Zhang X, Zhu Y, Chen X (2020). Fuzzy 2D linear discriminant analysis based on sub-image and random sampling for face recognition. *International Journal of Pattern Recognition and Artificial Intelligence*, 34(1), pp. 2056001.1-2056001.19. <http://doi.org/10.1142/S0218001420560017>
- [8] Gupta G, Dwivedi A, Rai V, Joshi S, Kumar R (2023). Oppositional grass hopper optimization with fuzzy classifier for face recognition from video database. *Wireless Personal Communications*, 132(3), pp. 1651. <http://doi.org/10.1007/s11277-023-10599-7>
- [9] Chakraborty S, Singh S K, Chakraborty P (2017). Local directional gradient pattern: a local descriptor for face recognition. *Multimedia Tools & Applications*, 76(1), pp. 1201-1216.
- [10] Chen C, Zhou X (2022). Collaborative representation-based fuzzy discriminant analysis for Face recognition. *The Visual Computer*, 38(4), pp. 1383-1393. <http://doi.org/10.1007/s00371-021-02325-w>
- [11] Sun J, Lv Y, Tang C, Sima H, Wu X (2020). Face recognition based on local gradient number pattern and fuzzy convex-concave partition. *IEEE Access*, 8, pp. 35777-35791. <http://doi.org/10.1109/ACCESS.2020.2975312>
- [12] Vishwakarma V P, Goel T (2019). An efficient hybrid DWT-fuzzy filter in DCT domain-based illumination normalization for face recognition. *Multimedia Tools & Applications*, 78(11), pp. 15213-15233. <http://doi.org/10.1007/s11042-018-6837-0>
- [13] Guo Z, Xiao Z, Alroobaea R, Baqasah A M, Althobaiti A, Gill H S (2022). Design and study of urban rail transit security system based on face recognition technology. *Informatica*, 46(3), pp. 429-438.
- [14] Dhamija A, Dubey R B (2022). An approach to enhance performance of age invariant face recognition. *Journal of Intelligent & Fuzzy Systems: Applications in Engineering and Technology*, 43(3), pp. 2347-2362.
- [15] Dey A, Chowdhury S (2020). Probabilistic weighted induced multi-class support vector machines for face recognition. *Informatica*, 44(4), pp. 459-467.

Evaluation of Optimally Tuned K-Nearest Neighbors for 30-Minute Blood Glucose Prediction in Type 1 Diabetes Using OhioT1DM Dataset

Yacine Hachi^{1*}, Soraya Tighidet¹, Kamal Amroun¹, Meriem Djouadi²

¹Limed Laboratory, Faculty of Exact Sciences, University of Bejaia, 06000 Bejaia, Algeria

²Department of Computer Science, Faculty of Exact Sciences, Echahid Hamma Lakhdar University, P. O. B. 789, 39000 El Oued, Algeria

E-mail: yacine.hachi@univ-bejaia.dz, soraya.tighidet@univ-bejaia.dz, kamal.amroun@univ-bejaia.dz, djouadi-meriem@univ-eloued.dz

*Corresponding author

Keywords: diabetes, predict, blood glucose levels, hypoglycemia, hyperglycemia, K-Nearest Neighbors, machine learning

Received: February 18, 2025

Diabetes is a long-term chronic medical condition with the potential to evolve into a global healthcare crisis, glycemic control is fundamental for the effective management of diabetes and the prevention of its associated complications. Forecasting future blood glucose levels (BGLs) for diabetic patients can help them avoid serious health problems. This study investigates the application of the KNN regression algorithm to predict future (BGLs), utilizing historical blood glucose measurements from twelve patients (six patients from the Ohio dataset version 2018 and six patients from the Ohio dataset version 2020) as the only input feature. Our proposed approach employed a methodology that utilized historical measurements to train predictive models. Specifically, we leveraged the following historical data points - (BGLs) at 4-hours, 8-hours, 12-hours, 16-hours, 20-hours, and 24-hours intervals - as input features to predict (BGLs) 30 minutes into the future. This study explores the impact of varying parameters of the KNN algorithm, such as the K value= [2,3,5,7,11], weights= ['uniform', 'distance'] and distance metric= ['euclidean', 'manhattan', 'minkowski'], on the performance of the model. Furthermore, we compared the obtained results of the KNN algorithm with other machine learning methods, including linear regression, Random Forests, Support Vector Machines, CatBoostRegressor, LightGBM, XGBoost, artificial neural networks and previous studies. Among these, KNN yielded the best results with optimal hyperparameters (k=2, Weights='distance', Metric='manhattan') in the tow version of datasets OhioT1DM V2018 and OhioT1DM V2020. The OhioT1DM V2018 dataset yielded optimal performance with an RMSE of 5.09 ± 0.91 mg/dl using a 24-hour window size, and an MAE of 2.42 ± 0.34 mg/dl with a 12-hour window size. For the OhioT1DM V2020 dataset, the best results were an RMSE of 5.56 ± 1.14 mg/dl with a 12-hour window size, and an MAE of 2.47 ± 0.34 mg/dl achieved using an 8-hour window size. This research confirms that KNN algorithm with optimal hyperparameters (k=2, Weights='distance', Metric='manhattan') can effectively predict blood glucose events, which will help prevent and reduce the occurrence of serious complications such as hypoglycemia and hyperglycemia.

Povzetek: Študija optimizira algoritem KNN za 30-minutno napovedovanje ravnih glukoze v krvi (BGL) pri sladkorni bolezni tipa 1 z uporabo podatkov OhioT1DM. Optimalni hiperparametri (k=2, Weights='distance', Metric='manhattan') so dosegli najboljše rezultate.

1 Introduction

Diabetes mellitus has emerged as one of the most pressing global health concerns, with over 463 million individuals affected in 2019, projections indicate that this figure is poised to escalate further, reaching an estimated 700 million by the year 2045 [1]. The treatment of diabetes type1, which primarily relies on the administration of external insulin, necessitates the frequent assessment of blood glucose levels, currently this monitoring is facilitated by continuous glucose monitoring devices

(CGM), which enable the collection and display of glucose concentrations in an almost continuous manner for multiple days [2]. In certain scenarios, CGM interface with an insulin pump, which mimics the natural functioning of the pancreas by administering small, on-demand doses of insulin, consequently over the past decade, researchers have dedicated their efforts to developing machine learning-based algorithms that can accurately predict future blood glucose levels [3]. Many previous studies have proposed numerous predictive algorithms. The

researchers in [4] compared three different models; an autoregressive model that utilized only glucose data, an autoregressive model that incorporated external insulin information, and an artificial neural network (ANN) that leveraged both glucose and insulin data. Furthermore, online adaptive models were employed to account for the inherent intra-individual and inter-individual variability present in the diabetic population. Hamdi et al. [5] proposed utilizing solely CGM data to forecast BGLs independently of other variables. To substantiate this approach, they investigated the application of (SVR) and differential evolution algorithms (DE). The authors in [6] presented a deep learning model utilizing a dilated recurrent neural network (DRNN) architecture to generate 30-minute forecasts of prospective glucose levels. Additionally, they leveraged a transfer learning approach that incorporated dilation to harness data from multiple participants. Dudukcu. H.V et al [7] proposed several advanced neural network architectures, constituting (LSTM), WaveNet, and Gated Recurrent Units (GRU). The hyperparameters of these models were tuned to optimize their operational efficiency. The authors in [8] proposed an autonomous deep learning model for personalized forecasting of multivariate BGLs. The proposed autonomous channel network acquires representations from input variables with appropriate sequence lengths and sampling periods, drawing on domain knowledge of the time-dependent relationships between the variables. Shuvo. M. et al. [9] described a neural network architecture comprising shared and clustered hidden layers. The shared hidden layers, composed of two stacked long short-term memory layers, learned generalized features across all data samples. In contrast, the clustered hidden layers, consisting of two dense layers, adapted to the gender-specific variations present in the data. This study presents a predictive model that utilizes the KNN algorithm to forecast blood glucose levels within a 30-minute in the future. The effectiveness of this model is evaluated through the root mean square error rate (RMSE) and mean absolute error (MAE). In addition, many experiments were conducted to determine the optimal values of K, weights and the distance measure. We also experimented with various window sizes in order to minimize the error and enhance the model's performance.

The fundamental question addressed in this article, to be explored further in the discussion section, is whether a straightforward model like KNN can achieve superior performance compared to advanced models such as deep learning in predicting blood glucose levels for diabetic patients within a short prediction horizon (30-minutes) using only historical CGM data?

The remainder of the paper is organized as follows. Section 2 data and preprocessing. Our methodology is detailed in Section 3. Experimental results are presented and discussed in Section 4. Section 5 presents the limitations. Finally, section 6 provides a conclusion of the paper.

2 Data and preprocessing

The study employed the two datasets, OhioT1DM in 2018 and OhioT1DM in 2020 [10], which both contain training and testing data. These datasets provide information on twelve individuals with type 1 diabetes over an eight-week timeframe, including details such as glucose levels, finger stick readings, bolus doses, basal rates, interim basal rates, meal consumption, exercise routine, and basal heart rate. The data for each patient is segregated into distinct training and testing subsets within the overall dataset. Table 1 presents the number of training and test samples for each patient, along with their respective split ratios.

Table 1: Number of training and test samples for each patient

OhioT1DM V2018				
ID_Patient	Gender	Age	Training Examples 80 %	Test Examples 20%
559	female	40–60	10796	2514
563	male	40–60	12124	2570
570	male	40–60	10982	2745
575	female	40–60	11866	2590
588	female	40–60	12640	2791
591	female	40–60	10847	2760
OhioT1DM V2020				
540	male	20–40	11947	2884
544	male	40–60	10623	2704
552	male	20–40	9080	2352
567	female	20–40	10858	2377
584	male	40–60	12150	2653
596	male	60–80	10877	2731

This study utilized a sliding window technique to transform time-series data into an input matrix and output vector dependent solely on BGLs measured by CGM every 5 minutes. We are focusing only on blood glucose levels as input values for all predictive models, excluding factors such as diet, physical activity, stress, or drug treatments. This approach was adopted due to the challenges in accurately measuring and quantifying these other variables in real-world scenarios. Moreover, additional factors could be considered to enhance the model's accuracy, but these may impose a burden on the patient. Additionally, this study aims to reduce the number of input features while improving model accuracy. We conducted several experiments to explore the impact of window size. Specifically, we used:

- P = 48 previous data points (4-h) to train the model and forecast the BGL for the next F = 6 data points (30-minutes).

- P = 96 previous data points (8-h) to train the model and forecast the BGL for the next F = 6 data points (30-minutes).
- P = 144 previous data points (12-h) to train the model and forecast the BGL for the next F = 6 data points (30-minutes).
- P = 192 previous data points (16-h) to train the model and forecast the BGL for the next F = 6 data points (30-minutes).
- P = 240 previous data points (20-h) to train the model and forecast the BGL for the next F = 6 data points (30-minutes).
- P = 288 previous data points (24-h) to train the model and forecast the BGL for the next F = 6 data points (30-minutes).

The Figure 1 depicts the illustration of the sliding window technique for window size = 4-hours.

3 Methodology

3.1 K- Nearest Neighbors (KNN) Algorithm

The K-Nearest Neighbors algorithm is a non-parametric supervised learning technique that can be applied to both classification and regression problems. The fundamental operation of this algorithm is predicated on the concept of similarity between the data points within the dataset [11].

The accuracy and efficacy of the K-Nearest Neighbors regression model are primarily influenced by the choice of the K parameter as well as the distance metric employed. A variety of techniques can be employed to measure the similarity between data points. In our investigation, we utilized the Euclidean formula (1), Manhattan formula (2), and Minkowski formula (3) distance metrics to quantify proximity.

The Euclidean distance between two points X and Y can be calculated using the formula (1):

$$D(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \tag{1}$$

The Manhattan distance between two points A and B can be calculated using the formula (2):

$$D(A, B) = \sum_{i=1}^m |A_i - B_i| \tag{2}$$

The Minkowski distance combines the Euclidean and Manhattan distance metrics to quantify the separation between two data points X and Y through a mathematical expression formula (3):

$$D(X, Y) = (\sum_{i=1}^m |X_i - Y_i|^x)^{\frac{1}{x}} \tag{3}$$

Additionally, the selection of the K value significantly impacts the performance of the K-Nearest Neighbors algorithm. For this reason, our study undertook multiple experiments to determine the optimal K value, as well as the most suitable distance metric to apply.

The key steps involved in the KNN regression procedure employed in our study are as follows:

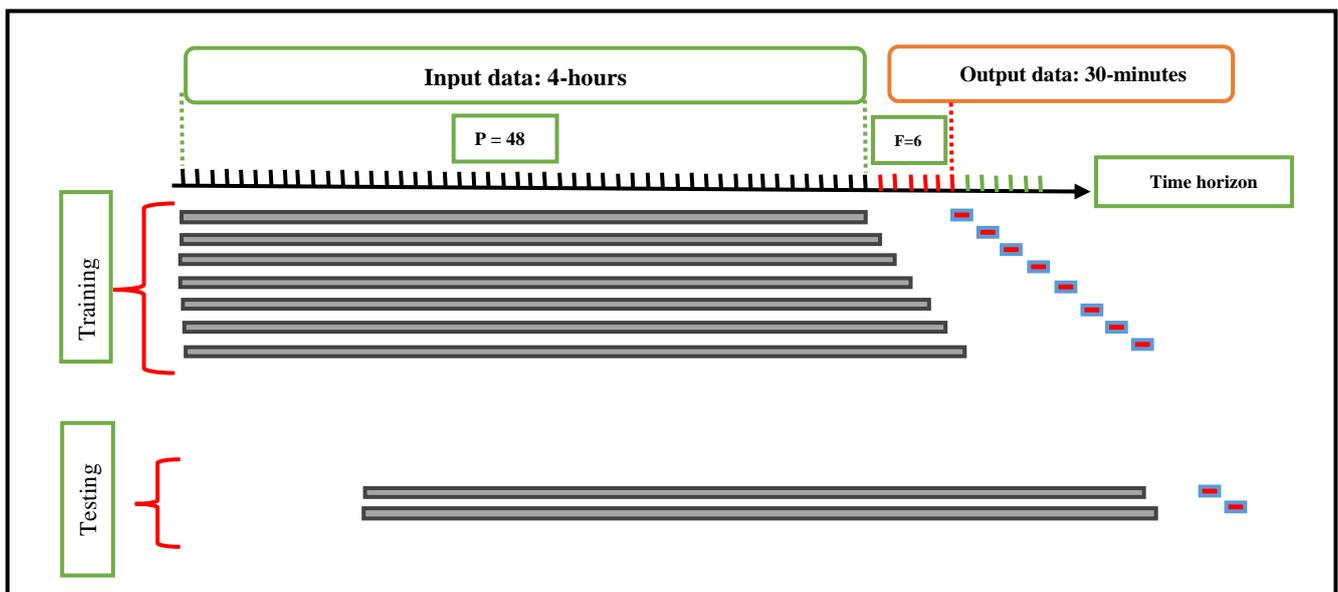


Figure 1: Sliding window approach with a window size of 4 hours

Procedure

Data: CGM from OhioT1DM

Result: Best minimum RMSE, MAE and the associated hyperparameter configuration.

1. Divide the dataset into training and test sets (X_{train} , X_{test} , y_{train} , y_{test}).
 2. Initialize the KNN model
 3. Hyperparameter values to be evaluated { # Define the parameter grid
 $N_neighbors = [2,3,5,7,11]$
 $Weights = ['uniform', 'distance']$
 $Metric = ['euclidean', 'manhattan', 'minkowski']$
 }.
 4. $kfold = KFold(n_splits=5)$ # Define 5-fold cross-validation
 5. $GridSearchCV(knn, hyperparameter\ values, cv=kfold)$ # Perform grid search with cross-validation for hyperparameter tuning
 6. Fit the hyperparameter values on the training data (X_{train} , y_{train}).
 7. Select the optimal model from the hyperparameter values.
 8. Predict the test set using the best model.
 9. Calculate **RMSE**, **MAE** for the predictions.
 10. Report the optimal **RMSE**, **MAE** and the associated hyperparameter configuration.
- End.
-

This study evaluated the performance of the K-Nearest Neighbors algorithm on the OhioT1DM Diabetes dataset through an experimental analysis. The aim of this investigation was to determine the optimal values for the number of neighbors (K), the Weights and the distance metric that would maximize the performance of the KNN algorithm.

The hyperparameters used in other machine learning models, including linear regression ($copy_X=True, fit_intercept=True, n_jobs=None, positive=False$), SVR($kernel='rbf', C=100, gamma=0.1, epsilon=0.1$), CatBoostRegressor($iterations: 1000, learning_rate: 0.01, depth: 6, random_state: 42$), Lightgbm($n_estimators: 1000, learning_rate: 0.01, num_leaves: 31, random_state: 42$), ANN with 3 layers($(128, activation='relu'), Dropout(0.3), (64, activation='relu'), Dropout(0.3), 1, optimizer='adam, loss='mean_squared_error'$) and ($epochs=100, batch_size=128$), XGBoost($learning_rate: 0.01, n_estimators: 1000, random_state: 42, num_boost_round=1000$), RandomForest($n_estimators=100, max_depth=None, min_samples_split=2, min_samples_leaf=1, bootstrap=True$). The schematic diagram presented in Figure 2 illustrates the methodology of our approach used in the predictive models. Forecasting models were implemented in Python (version 3.9.12) using CPU, with scikit-learn (version 1.6.1), NumPy (version 1.24.3), Pandas (version 1.4.2), and TensorFlow (version 2.13.0).

4 Results

The metrics employed in the evaluation of the efficiency and accuracy of the algorithms were as follows:

4.1 Root Mean Square Error (RMSE) and Mean Absolutely Error (MAE)

The RMSE, a widely employed metric, measures the average magnitude of the errors in a predictive model, can be characterized as:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (4)$$

Where the predicted output is denoted as \hat{y}_i and the true output is denoted as y_i . The root mean square error metric offers several desirable properties, including a readily defined gradient, intuitive interpretation, and the ability to transform the error back to the original scale through the square root operation [12].

The mean absolute error is quantified as:

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (5)$$

This error measure is simple to formulate and exhibits a degree of insensitivity to outliers, we used MAE as a second metric for assessing the accuracy of our regression model.

4.2 Comparison of our approaches with the relevant prior work in the literature

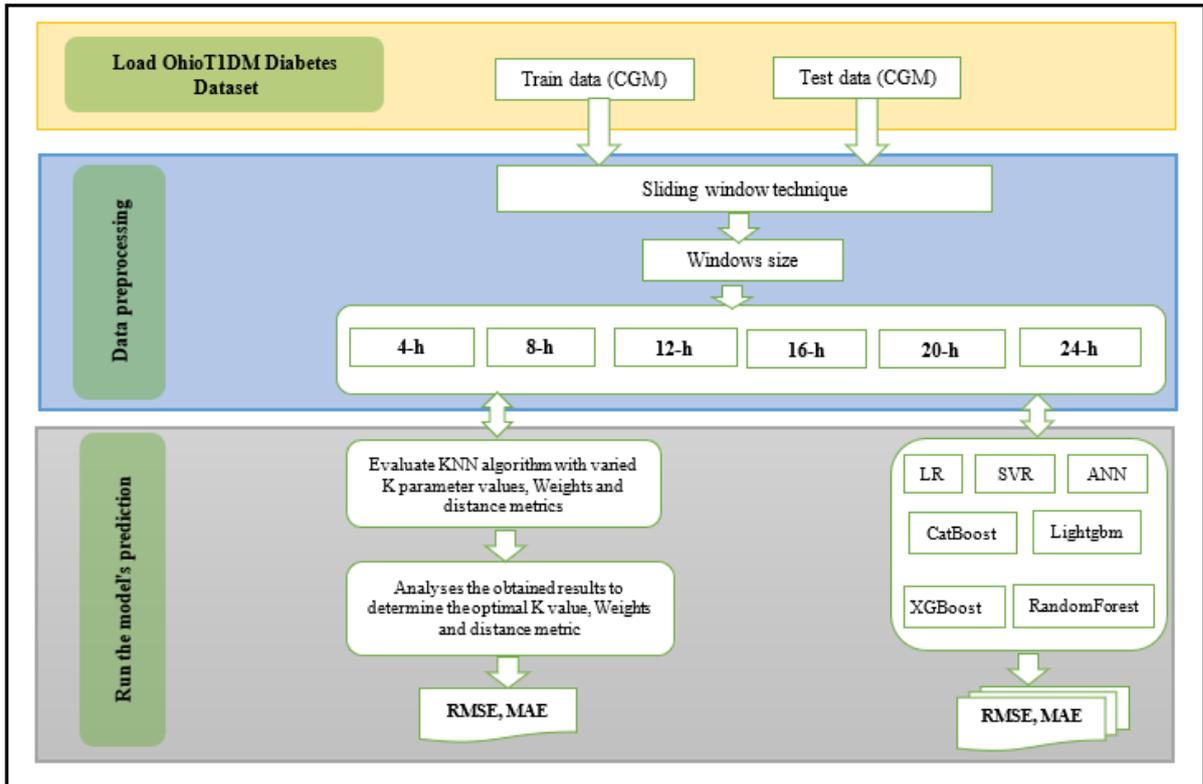


Figure 2: Flowchart of our approach used in the predictive models.

Table 2: Mean with standard deviation (Mean ± Std) of RMSE, MAE for Ohio T1DM Dataset version 2018 in different windows size

Name of the model	window size											
	4-h		8-h		12-h		16-h		20-h		24-h	
	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE
LR	23.63 ± 3.89*	16.18 ± 2.60**	23.71 ± 3.89**	16.26 ± 2.61**	23.73 ± 3.87**	16.29 ± 2.61**	23.81 ± 3.94**	16.38 ± 2.67**	23.79 ± 3.84**	16.40 ± 2.63**	23.86 ± 3.82**	16.47 ± 2.65**
SVR	58.62 ± 7.87**	47.19 ± 6.55**	58.64 ± 7.83**	47.21 ± 6.51**	58.66 ± 7.82***	47.22 ± 6.48**	58.63 ± 7.82**	47.19 ± 6.49**	58.56 ± 7.83**	47.12 ± 6.50**	58.51 ± 7.89**	47.06 ± 6.54**
CatBoostRegressor	22.31 ± 3.49*	15.66 ± 2.36**	21.86 ± 3.37**	15.40 ± 2.28**	21.53 ± 3.20**	15.22 ± 2.22**	21.40 ± 3.24**	15.14 ± 2.27**	21.16 ± 3.23**	15.00 ± 2.27**	20.95 ± 3.15**	14.88 ± 2.23**
Lightgbm	20.11 ± 3.16*	13.95 ± 2.11**	18.72 ± 2.84**	13.01 ± 1.90**	17.84 ± 2.65**	12.44 ± 1.79**	17.28 ± 2.64*	12.05 ± 1.79**	16.82 ± 2.63*	11.75 ± 1.79**	16.30 ± 2.44*	11.40 ± 1.65**
ANN	33.36 ± 4.13**	26.18 ± 3.26**	39.02 ± 3.47***	31.83 ± 2.71***	42.35 ± 3.61***	34.97 ± 4.61**	46.09 ± 5.68***	38.60 ± 6.24**	48.30 ± 5.10***	40.64 ± 5.04***	49.25 ± 3.85***	41.40 ± 3.38***
XGBoost	19.95 ± 3.17*	13.78 ± 2.09**	18.60 ± 2.84**	12.86 ± 1.88**	17.70 ± 2.68**	12.27 ± 1.78**	17.15 ± 2.64*	11.88 ± 1.78**	16.84 ± 2.65*	11.66 ± 1.81**	16.35 ± 2.54*	11.34 ± 1.71**
RandomForest	18.89 ± 3.17*	12.23 ± 2.01*	17.50 ± 2.92*	11.00 ± 1.79**	16.73 ± 2.73*	10.37 ± 1.68**	16.33 ± 2.69*	10.00 ± 1.67*	15.99 ± 2.75*	9.77 ± 1.72*	15.65 ± 2.69*	9.50 ± 1.65*
KNN with best hyperparameter (K=2, Weights='distance', Metric='manhattan')	7.78 ± 1.63	3.30 ± 0.56	5.41 ± 1.10	2.45 ± 0.35	5.30 ± 1.21	2.42 ± 0.34	5.29 ± 0.99	2.45 ± 0.38	5.22 ± 1.00	2.48 ± 0.35	5.09 ± 0.91	2.47 ± 0.39

*p ≤ 0.0001 **p ≤ 0.00001 ***p ≤ 0.000001.

Table 3: Mean with standard deviation (Mean \pm Std) of RMSE, MAE for Ohio T1DM Dataset version 2020 in different windows size

Name of the model	window size											
	4-h		8-h		12-h		16-h		20-h		24-h	
	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE
LR	24.64 \pm 3.96 *	17.04 \pm 2.57 **	24.70 \pm 3.97 *	17.10 \pm 2.60 **	24.77 \pm 3.99 *	17.15 \pm 2.60 **	24.87 \pm 4.02 *	17.23 \pm 2.60 **	24.92 \pm 3.99 *	17.29 \pm 2.59 **	24.92 \pm 4.00 *	17.33 \pm 2.59 **
SVR	57.56 \pm 5.06 ***	45.59 \pm 3.68 ***	57.73 \pm 5.18 ***	45.72 \pm 3.72 ***	57.74 \pm 5.19 ***	45.72 \pm 3.73 ***	57.78 \pm 5.18 ***	45.75 \pm 3.71 ***	57.82 \pm 5.20 ***	45.79 \pm 3.74 ***	57.87 \pm 5.23 ***	45.84 \pm 3.76 ***
CatBoostRegressor	23.39 \pm 3.68 *	16.64 \pm 2.55 **	22.80 \pm 3.61 *	16.31 \pm 2.52 **	22.45 \pm 3.50 *	16.10 \pm 2.48 **	22.27 \pm 3.51 *	15.98 \pm 2.49 **	22.07 \pm 3.44 *	15.89 \pm 2.44 **	21.85 \pm 3.48 *	15.72 \pm 2.43 **
Lightgbm	21.09 \pm 3.49 *	14.74 \pm 2.35 *	19.28 \pm 3.22 *	13.57 \pm 2.18 **	18.41 \pm 3.06 *	12.97 \pm 2.11 *	17.76 \pm 3.08 *	12.50 \pm 2.10 *	17.20 \pm 2.94 *	12.20 \pm 2.06 *	16.83 \pm 3.01 *	11.87 \pm 2.03 *
ANN	33.81 \pm 4.96 **	26.15 \pm 4.18 **	37.37 \pm 3.66 ***	29.52 \pm 2.79 ***	42.30 \pm 5.26 **	34.32 \pm 4.85 **	42.96 \pm 3.71 ***	34.78 \pm 3.47 ***	47.48 \pm 5.10 ***	39.00 \pm 4.18 ***	49.24 \pm 4.50 ***	40.81 \pm 3.67 ***
XGBoost	20.82 \pm 3.43 *	12.89 \pm 5.72 *	19.11 \pm 3.03 *	13.40 \pm 1.98 **	18.28 \pm 2.85 *	12.97 \pm 1.95 **	17.73 \pm 2.77 *	12.37 \pm 1.83 **	17.24 \pm 2.69 *	12.12 \pm 1.84 *	16.82 \pm 2.73 *	11.83 \pm 1.81 *
RandomForest	19.65 \pm 3.37 *	12.94 \pm 2.08 *	17.96 \pm 3.18 *	11.55 \pm 1.89 *	17.21 \pm 2.92 *	10.92 \pm 1.77 *	16.75 \pm 2.89 *	10.52 \pm 1.74 *	16.21 \pm 2.76 *	10.17 \pm 1.68 *	15.95 \pm 2.84 *	9.91 \pm 1.66 *
KNN with best hyperparameter (K=2, Weights='distance', Metric='manhattan')	8.27 \pm 1.66	3.38 \pm 0.48	5.65 \pm 1.23	2.47 \pm 0.34	5.56 \pm 1.14	2.49 \pm 0.34	5.61 \pm 1.13	2.57 \pm 0.33	5.59 \pm 1.01	2.59 \pm 0.34	5.60 \pm 1.14	2.64 \pm 0.35

* $p \leq 0.0001$ ** $p \leq 0.00001$ *** $p \leq 0.000001$.

4.2 Discussion

Table 2 and Table 3 show the obtained results, which are the Mean with standard deviation (Mean \pm Std) of RMSE, MAE for datasets OhioT1DM V2018 and OhioT1DM V2020 in different windows sizes. Moreover, we apply 5-Fold Cross-Validation to assess the predictive performance of our models. Table 2 and Table 3 display the results of RMSE, MAE achieved through various methods for predicting blood glucose Levels over 30-Minute Timeframe. Furthermore, paired t-tests were conducted to ascertain the statistical significance of the KNN algorithm when employing its optimal hyperparameters.

The presented data in Table 2 and Table 3 indicate that the KNN regression algorithm with the optimal hyperparameter values of $k=2$, Weights='distance', and Metric='manhattan' is significant for predicting blood glucose levels over a 30-minutes timeframe. The KNN model achieved a minimum RMSE of 5.09 ± 0.91 mg/dl using a 24-hour window size, and an MAE of 2.42 ± 0.34 mg/dl with a 12-hour window size in OhioT1DM V2018 dataset. For the OhioT1DM V2020 dataset, the best results were an RMSE of 5.56 ± 1.14 mg/dl with a 12-hour window size, and an MAE of 2.47 ± 0.34 mg/dl achieved using an 8-hour window size.

Decision tree-based models, including RandomForest, XGBoost, Lightgbm, and CatBoostRegressor, showed that window sizes had a notable impact on RMSE and MAE. Specifically, the RMSE and MAE at a 24-hour window size were lower than those at a 4-hour window size.

The RandomForest model exhibited the best performance, yielding an RMSE of 15.65 ± 2.69 mg/dl, MAE of 9.50 ± 1.65 mg/dl with a 24-h window size in dataset OhioT1DM V2018 and an RMSE of 15.95 ± 2.84 mg/dl, MAE of 9.91 ± 1.66 mg/dl with a 24-h window size in dataset OhioT1DM V2020.

The linear regression (LR) model yielded highly similar results, and the selection of time window size did not significantly impact the performance, the RMSE values ranged from 23.63 ± 3.89 mg/dl to 23.86 ± 3.82 mg/dl in dataset OhioT1DM V2018 with the lowest RMSE of 23.63 ± 3.89 mg/dl observed when using a 4-h time window, and ranged from 24.64 ± 3.96 mg/dl to 24.92 ± 4.00 mg/dl in dataset OhioT1DM V2020 with the lowest RMSE of 24.64 ± 3.96 mg/dl observed when using a 4-h time window.

The performance of the SVR and ANN models was unsatisfactory, even with adjustments to the window size. The SVR model yielded RMSE values ranging from 58.51 ± 7.89 mg/dl to 58.66 ± 7.82 mg/dl in dataset OhioT1DM V2018, and ranging from 57.56 ± 5.06 mg/dl to 57.87 ± 5.23 mg/dl in dataset OhioT1DM V2020. Similarly, the ANN model exhibited RMSE values between 33.36 ± 4.13 mg/dl to 49.25 ± 3.85 mg/dl in dataset OhioT1DM V2018, and RMSE values between 33.81 ± 4.96 mg/dl to 49.24 ± 4.50 mg/dl in dataset OhioT1DM V2020.

As demonstrated in Table 4, the prediction of BGLs in diabetic patients in this study is comparable to the recent research reported in the literature using alternative methodologies. T. Hamdi et al. [5] explored the use of support vector regression and differential evolution techniques, they achieved a minimum RMSE of 10.78 mg/dl. T. Zhu et al. [6] developed a deep learning model that employed a dilated recurrent neural network architecture (DRNN), they achieved a minimum RMSE of 18.90 mg/dl. Dudukcu. H.V et al. [7] explored various sophisticated neural network designs, such as LSTM, WaveNet, and Gated Recurrent Units, they achieved a minimum RMSE of 21.90 mg/dl. T. Yang et al. [8] developed an autonomous deep learning model for personalized prediction of multiple blood glucose

measures, they achieved a minimum RMSE of 18.93 mg/dl. Shuvo. M. et al. [9] outlined a neural network architecture with shared and clustered hidden layers, they achieved a minimum RMSE of 16.06 ± 2.74 mg/dl.

Although the KNN algorithm is a traditional machine learning method, our study found it outperformed deep learning techniques, particularly with smaller datasets. KNN's performance, however, is sensitive to the choice of the value for k . To optimize the algorithm's effectiveness, we evaluated various k values, weights, distance metrics, and window sizes. To the best of our knowledge, no prior research has compared KNN against modern methods while systematically testing these parameters.

Our study yielded promising results. For the OhioT1DM V2018 dataset, we achieved a minimum RMSE of 5.09 ± 0.91 mg/dl using a 24-hour window size, and an MAE of 2.42 ± 0.34 mg/dl with a 12-hour window size. On the OhioT1DM V2020 dataset, the best results were an RMSE of 5.56 ± 1.14 mg/dl (12-hour window size) and an MAE of 2.47 ± 0.34 mg/dl (8-hour window size), obtained by employing a k -nearest neighbors algorithm with the following parameter settings: ($k = 2$, weights = 'distance', and metric = 'manhattan').

Table 4: Comparison of our approaches with prior approaches in the literature

Author s	Techniques	Datasets	Forecast Horizon	RMS E (mg/d L)	MAE (mg/dL)
T. Hamdi et al. [5]	SVR based on DE algorithm	12 type1 diabetes	30_min	10.78	-
T. Zhu et al. [6]	DRNN	OhioT1DM	30_min	18.90	-
Dudukcu. H.V et al. [7]	LSTM, Wave-Net, GRU	OhioT1DM	30_min	21.90	-
T. Yang et al. [8]	AC-DLF	OhioT1DM	30_min	18.93	-
Shuvo. M. et al. [9]	D-MTL	OhioT1DM	30_min	16.06 ± 2.74	10.64 ± 1.35
Our Method s	KNN With the best hyperparameter	OhioT1DM V2018	30_min	5.09 ± 0.91	2.42 ± 0.34
		OhioT1DM V2020	30_min	5.56 ± 1.14	2.47 ± 0.34

5 Limitations

This study encountered certain limitations, which we address first. The K-Nearest Neighbors algorithm, while capable of delivering strong predictive performance, particularly when dealing with datasets exhibiting distinct separation boundaries, incurs significant computational costs during the prediction phase. The algorithm's need to compare the query instance against every point within the training dataset to identify the closest neighbors leads to a time complexity of $O(n \times d)$, where ' n ' represents the count of training samples and ' d ' denotes the dimensionality. In embedded or real-time systems, where low latency and energy efficiency are paramount, the aforementioned

factor can represent a substantial limitation. Consequently, a balance must be struck between predictive accuracy and execution-time performance, necessitating careful evaluation before implementing KNN in these operational contexts. To mitigate this problem, techniques such as approximate nearest neighbor search, dimensionality reduction, or prototype selection can be employed, representing the future direction of our experiments.

A further limitation is that models trained exclusively on continuous glucose monitoring data might overfit to individual patient-specific patterns, hindering their generalizability across diverse patient populations or in response to altered daily routines. Furthermore, CGM-based models primarily demonstrate efficacy in short-term glucose trend predictions, with a constrained capacity to forecast long-term trends without integrating other factors such as bolus doses, basal rates, interim basal rates, meal consumption, exercise routines, and basal heart rate.

6 Conclusion

Accurately predicting BGLs in diabetes is crucial, as this will enable the artificial pancreas to secrete the necessary amount of insulin. This paper presents a functioning method for forecasting BGLs in real individuals with type 1 diabetes, using data from continuous glucose monitoring. The results indicate that the k -nearest neighbors' algorithm with optimal hyperparameters ($k=2$, Weight='distance', Metric='manhattan') was effective in predicting blood glucose levels over the short-term (30-minute forecast horizon) in all sliding windows size 8-hours, 12-hours, 16-hours, 20-hours, 24-hours. Additionally, the window size had an effect on reducing the error rate of the predictions. Specifically, the RMSE and MAE at 8-hours, 12-hours, 16-hours, 20-hours, 24-hours window size were lower than those at a 4-hour window size.

The model has been evaluated on the two versions of OhioT1DM V2018 and V2020. These datasets contain real values of measurements taken from real diabetic patients living in their natural environments, which considered as strong point to validate our findings. In future research, we suggest to integrate and experiment this model into CGM device and clinical decision-making tools.

Acknowledgment

This study was supported by the General Directorate for Scientific Research and Technological Development (DGRSDT), Algeria. Moreover, it was conducted as part of the research activities at LIMED laboratory, which is affiliated with the Exact Sciences Faculty at Bejaia University.

References

- [1] Liu, K., Li, L., Ma, Y., Jiang, J., Liu, Z., Ye, Z., ... Yi, W. "Machine Learning Models for Blood Glucose Level Prediction in Patients With Diabetes Mellitus: Systematic Review and

- Network Meta-Analysis". *JMIR medical informatics*, 11(1). 2023. DOI:10.2196/47833
- [2] Prendin, F., Del Favero, S., Vettoretti, M., Sparacino, G., Facchinetti, A. "Forecasting of Glucose Levels and Hypoglycemic Events: Head-to-Head Comparison of Linear and Nonlinear Data-Driven Algorithms Based on Continuous Glucose Monitoring Data Only". *Sensors 2021*, Vol. 21, Page 1647, 21(5), pp. 1647. 2021. DOI:10.3390/S21051647
- [3] D'Antoni, F., Merone, M., Piemonte, V., Iannello, G., Soda, P. "Auto-Regressive Time Delayed jump neural network for blood glucose levels forecasting". *Knowledge-Based Systems*, 203, pp. 106134. 2020. DOI:10.1016/J.KNOSYS.2020.106134
- [4] Daskalaki, E., Prountzou, A., Diem, P., Mougiakakou, S. G. "Real-Time Adaptive Models for the Personalized Prediction of Glycemic Profile in Type 1 Diabetes Patients". <https://home.liebertpub.com/dia>, 14(2), pp. 168–174. 2012. DOI:10.1089/DIA.2011.0093
- [5] Hamdi, T., Ben Ali, J., Di Costanzo, V., Fnaiech, F., Moreau, E., Ginoux, J. M. "Accurate prediction of continuous blood glucose based on support vector regression and differential evolution algorithm". *Biocybernetics and Biomedical Engineering*, 38(2), pp. 362–372. 2018. DOI:10.1016/J.BBE.2018.02.005
- [6] Zhu, T., Li, K., Chen, J., Herrero, P., Georgiou, P. "Dilated Recurrent Neural Networks for Glucose Forecasting in Type 1 Diabetes". *Journal of Healthcare Informatics Research*, 4(3), pp. 308–324. 2020. DOI:10.1007/S41666-020-00068-2/FIGURES/7
- [7] Dudukcu, H. V., Taskiran, M., Yildirim, T. "Blood glucose prediction with deep neural networks using weighted decision level fusion". *Biocybernetics and Biomedical Engineering*, 41(3), pp. 1208–1223. 2021. DOI:10.1016/J.BBE.2021.08.007
- [8] Yang, T., Yu, X., Ma, N., Wu, R., Li, H. "An autonomous channel deep learning framework for blood glucose prediction". *Applied Soft Computing*, 120, pp. 108636. 2022. DOI:10.1016/J.ASOC.2022.108636
- [9] Shuvo, M. M. H., Islam, S. K. "Deep Multitask Learning by Stacked Long Short-Term Memory for Predicting Personalized Blood Glucose Concentration". *IEEE journal of biomedical and health informatics*, PP(3), pp. 1612–1623. 2023. DOI:10.1109/JBHI.2022.3233486
- [10] Marling, C., Bunescu, R. "The OhioT1DM Dataset for Blood Glucose Level Prediction: Update 2020". *CEUR workshop proceedings*, 2675, pp. 71. 2020. Retrieved from /pmc/articles/PMC7881904/ PMID: 33584164; PMCID: PMC7881904.
- [11] Cao, N., Yan, X. E., Zhang, L., Xu, G., Ma, J. "Hybrid K-Nearest Neighbors Models with Metaheuristic Optimization for Predicting Undrained Shear Strength". *Informatica*, 49(25), pp. 125–144. 2025. DOI:10.31449/INF.V49I25.7723
- [12] Xiong, Y. "Development of an AI-Driven Model for Drug Sales Prediction Using Enhanced Golden Eagle Optimization and XGBoost Algorithm". *Informatica*, 49(17), pp. 37–50. 2025. DOI:10.31449/INF.V49I17.7491

Abbreviation

ANN: Artificial neural network; **BG:** Blood Glucose; **BGLs:** Blood Glucose levels; **CatBoostClassifier:** Categorical Boosting Classifier; **CGM:** Continuous Glucose Monitor; **KNN:** K-nearest neighbor; **LightGBM:** Light gradient-boosting machine; **RF:** Random Forest; **SVM:** Support Vector Machine; **T1DM:** Type 1 Diabetes mellitus; **XGBoost:** Extreme Gradient Boosting; **RMSE:** Root Mean Square Error; **h:** hours; **MAE:** Mean Absolutely Error.

English Text Classification Model Based on Graph Neural Network Algorithm and Contrastive Learning

Chen Sian, Pan Guoqiang

Zhejiang Institute of Communications, Hangzhou, Zhejiang, 311112, China

E-mail: chen_sian82@outlook.com

Keywords: graph neural network, contrastive learning, english text classification, semantic representation, model construction

Received: March 1, 2025

Current English text classification methods mostly rely on bag-of-words models or CNN (Convolutional Neural Network), but there are limitations in processing text structure and semantics. Especially in long texts and complex contexts, it is difficult to capture the long-distance dependency and structured semantics between words. To this end, this article combines GNN (Graph Neural Network) with contrastive learning to build an English text classification model. First, a text graph is constructed through word co-occurrence to capture the long-distance dependency of words. Then, a multi-layer graph convolutional network is designed, and residual connections and normalization are applied to improve model performance. A contrast learning module is added after each layer of graph convolution to improve node features and semantic representation. Triplet Loss is a loss function, and Hard Negative Mining chooses negative samples to improve efficiency.

Povzetek: Predlagan je model za klasifikacijo angleškega besedila, ki združuje grafične nevronske mreže (GNN) in kontrastno učenje (CL). GNN-ji s pomočjo ko-pojavitvene matrike ustvarijo graf za zajemanje medsebojnih odvisnosti besed. CL (z izgubo Triplet Loss) izboljša semantično reprezentacijo vozlišč GNN, kar model (CS-K-prototipi) pri klasifikaciji besedil bistveno izboljša natančnost in robustnost.

1 Introduction

English text classification is critical for natural language processing, affecting many aspects such as information retrieval, sentiment analysis, and question-answering systems. The importance of text classification is increasing with the proliferation of networks and information content. Although traditional methods such as bag-of-words models and TF-IDF (term frequency-inverse document frequency) have some effects, they cannot cope with deep semantics, long texts, or complex contexts. Although deep learning methods such as CNN and RNN (Recurrent Neural Network) have progressed, they still cannot perfectly handle long-distance dependencies. Graph Neural Networks (GNN) and contrastive learning are two deep learning-based pattern recognition algorithms that improve accuracy, adaptability, and efficiency over rule-based approaches in software testing. Complex data linkages are captured, unstructured data is handled, manual feature engineering is decreased, and generalization is improved. To enhance text classification, new methods and technologies have emerged. GNN can capture structured information, simulate complex relationships between words, and deepen text understanding. Contrastive learning is a strategy for improving model performance by comparing data samples and moving similar ones closer together while pushing dissimilar ones apart. In this study, semantic embeddings are optimized using Triplet Loss, which improves text categorization accuracy. Contrastive

learning, a self-supervised strategy, can strengthen the feature representation of the model and reduce the dependence on labeled data. This study combines GNN and contrastive learning to create a new English text classification model. Graph Neural Networks have improved at dealing with extensive texts and intricate interactions than models such as CNNs and RNNs because they can represent text as a graph, with words as nodes and relationships as edges. This enables GNNs to capture long-distance dependencies and intricate semantic linkages between words, which CNNs struggle with due to their reliance on fixed-size filters, while RNNs suffer with long sequences due to concerns such as vanishing gradients. GNNs improve their understanding of a text's global structure and deep semantics by pooling input from surrounding nodes. Combining GNNs with contrastive learning improves feature representation, allowing for accurate and robust handling of complicated and lengthy texts. This model can not only effectively grasp the grammatical and semantic relationships of the text but also improve semantic embedding through contrastive learning, thereby improving classification accuracy. This study has injected new vitality into the field of text classification and promoted the advancement of related technologies. Semantic embedding transforms text into a vector space in which comparable words are closer together, allowing the model to better grasp word associations. By integrating Graph Neural Networks (GNN) and contrastive learning, the study improves semantic embeddings, allowing the model to capture

complex relationships and long-distance dependencies for better text classification.

Text semantic extraction helps precisely identify the text's relationship and structure and provides intelligent support for information retrieval, sentiment analysis, etc. Improving semantic understanding can improve the accuracy of natural language processing tasks and meet personalized information needs [1-2]. To solve the problem of text semantic information extraction, many scholars have proposed different methods. Martinez-Rodriguez J. L proposed a strategy for extracting information from sentences and representing it using semantic network standards. The strategy involved information extraction tasks and hybrid semantic similarity metrics. Experiments proved the proposed method's feasibility and accuracy [3]. Current Chinese short text entity linking techniques ignore the interaction between label information and the original short text and fail to effectively utilize semantic information. Gao L proposed a normalization method to fully extract semantic information from short text sentence vectors. The results showed that the proposed model outperformed popular deep-learning techniques and previous research results in entity linking [4]. The process of automatically determining the connection between two or more elements is called semantic relation extraction, which is crucial for creating original writing. In addition to describing established and new evaluation metrics for supervised, semi-supervised, and unsupervised methods, Gharagozlou H also studied several relation extraction techniques and types in English and the most popular techniques in Persian [5]. Accurate identification and analysis of semantics are conducive to effectively processing English text. Yu S introduced Word2vec (word to vector) for extracting semantic feature vectors from English text and the long short-term memory (LSTM) algorithm for semantic identification of English text. The results showed that the identification results of the LSTM algorithm for the part of speech and sentiment tendency of English text were consistent with the label results [6]. Scholars emphasize semantic information and use different techniques to improve model performance. However, the research has not fully utilized the graph structure to capture the complex relationship of text. Sequential models, which concentrate mostly on local aspects and may have trouble with long-distance dependencies, frequently miss deep semantic linkages that graph structures capture. A graph structure preserves both local and global dependencies by representing words as nodes and their interactions as edges, in contrast to CNNs and RNNs that analyze text sequentially. This makes it possible to comprehend text semantics more precisely, especially in intricate circumstances. This representation

is further improved by combining contrastive learning with Graph Neural Networks (GNNs), which increases the model's capacity to differentiate between text categories and boosts performance on tasks with intricate semantics and long-distance dependencies. It has limitations in the processing of inter-lexical dependencies and deep semantic structures.

GNN can better capture the relationship and meaning between words by turning text into a graph, enhance the model's understanding of text structure, and improve classification accuracy, especially in processing long articles and tasks that require an understanding of context [7-8]. In recent years, some researchers have used GNN in text classification research. To solve the problem of cross-lingual text classification, Vo T proposed a new topic-driven multi-type text graph attention representation learning technology, which combined neural topic modelling technology with a heterogeneous text graph attention network to enhance the semantic information of text representation learned in various language environments. GAT and GraphSAGE are two models with distinct advantages in text classification problems. GAT incorporates an attention mechanism into graph convolutional layers, allowing the model to focus on meaningful words or relationships, hence enhancing accuracy. GraphSAGE minimizes computational complexity by sampling neighbors during training and enhances scalability, particularly for large-scale graphs. Its aggregation approaches, such as mean, pooling, and LSTM-based aggregators, enable the model to capture broad semantic patterns while avoiding overfitting. When coupled, these models could offer a more robust method for dealing with complicated semantics and long-distance interdependence.

The proposed model was compared with the current state-of-the-art baseline and experimentally demonstrated its effectiveness [9]. Deng Z proposed a new graph-based model and designed an attention-gated graph neural network to propagate and update the semantic information of each word node to solve the problem that existing methods are not enough to capture the semantic relationship between words. Experimental results showed that the proposed model outperformed previous text classification methods [10]. Parthasarathy (2023) examines combining neural networks with the Harmony Search Algorithm (HSA) to improve fraud detection in banking. Traditional methods often fail against complex fraud techniques, but this combination enhances accuracy and reliability. The findings suggest that models like Decision Tree Classifier and Sequential models, with near-perfect accuracy, could transform fraud prevention. However, this study supports the idea that combining neural networks with the Harmony Search Algorithm (HSA) to improve fraud detection parallels the approach in our work to enhance accuracy and reliability in text classification [11]. The above scholars have cleverly used

graph structures and attention mechanisms to grasp the semantic relationship of text, significantly improving the model's ability to handle complex semantics and performing well in cross-language text classification. However, the graph neural network's capture of deep semantics and long-distance dependencies needs to be strengthened, and it has not fully utilized contrastive learning to enhance feature representation.

This study combines GNN with contrastive learning to deeply capture word relationships by constructing a text graph and optimizing semantic embedding using contrastive learning. GNN converts text into a graph with words as nodes and edges showing semantic connections. The results show that this method can deeply capture text semantics and long-distance dependencies, significantly improve performance in multiple text classification tasks, and demonstrate its excellent generalization ability and robustness. Compared with CNN, this method is more precise and stable when dealing with complex text classification. The innovation of this study is to combine GNN with contrastive learning for English text classification, which makes up for the shortcomings of traditional methods and reduces the dependence on labeled data. The new graph contrast loss function captures text semantics more precisely. The graph contrast loss function has numerous significant advantages over ordinary contrastive loss functions, especially in the context of graph neural networks (GNN) and contrastive learning for text categorization. It takes advantage of the network structure of text to improve the model's capacity to capture semantic linkages between words, addressing the complex, long-distance dependencies that typical contrastive loss functions frequently overlook. By taking into account both pairwise similarities and contextual relationships within the network, the graph contrast loss function improves semantic embedding quality, resulting in more accurate and informative text representations. The function also includes Hard Negative Mining, which concentrates on difficult-to-detect negative samples, allowing the model to acquire more discriminatory features and enhance generalization. At the same time, through strategy optimization and parameter adjustment, the classification accuracy and training efficiency are improved. These innovations have promoted the development of text classification technology and provided new ideas for natural language processing. The organizational structure of this study is shown in Figure 1:

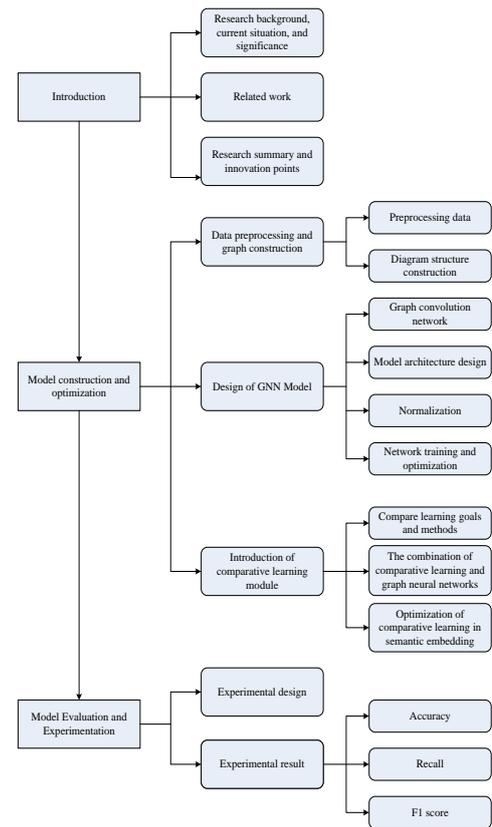


Figure 1: Organizational structure of this study

2 Construction and optimization of the english text classification model

2.1 Data Preprocessing and graph construction

2.1.1 Preprocessing data

Data preprocessing is critical for English text classification. It can clean text, remove noise, and convert it into a format suitable for GNN. The proposed approach combines contrastive learning with Graph Neural Networks (GNN) to handle noisy or redundant test instances efficiently. GNN focuses on pertinent semantic connections while capturing long-distance dependencies and deep semantic interconnections between words. By differentiating between comparable and dissimilar samples, contrastive learning improves the resilience of the model. By focusing on hard-to-classify negative samples, hard negative mining improves classification accuracy and lessens the influence of redundant data.

This article cleans, segments, removes stop words, and stems the data to extract valuable semantic information, laying the foundation for subsequent graph construction and graph neural network training. Text cleaning is the first step in data preprocessing. The original text contains many irrelevant information, such as punctuation, numbers, etc., which may interfere with the analysis. This article uses regular expressions to clean up these noises, retaining only letters and spaces to ensure the text's purity. The next step is segmenting the text into independent words or phrases. Using the word tokenize method of NLTK (Natural Language Toolkit), word segmentation is performed by space and punctuation, and the text is converted into a vocabulary list. Afterword segmentation, stop words are filtered out. Removing stop words can reduce the amount of calculation and prevent the model from being interfered with by irrelevant information. NLTK stop word library is utilized for filtering. After that, stemming is done to normalize different forms of words to the basic form, such as "running" becomes "run". This can reduce the number of words and vocabulary dimensions and improve training efficiency. Then, the Porter stemming algorithm and the Porter Stemmer class of NLTK are processed. After this preprocessing, the original text becomes a preprocessed standardized vocabulary list.

Word segmentation, stop word elimination, and stemming are critical processes that ensure efficient processing of data in order to prepare it for text classification model operation. Since the text is represented as a graph with words as nodes in models like Graph Neural Networks (GNNs), word segmentation is crucial since it separates the text into discrete words or tokens. By getting rid of popular but useless words, stop word removal lowers computing costs and directs the model's attention to more important data. Stemming minimizes vocabulary quantity and increases training efficiency by breaking words down to their most basic forms.

2.1.2 Graph structure construction

Each document is treated as a graph. Among them, words correspond to nodes; edges between nodes represent semantic relationships between words; edge weights represent the strength of the relationship. After data preparation, a text graph structure is constructed for GNN. In this study, edges are built based on word co-occurrence information, and the co-occurrence matrix is used to quantify the word association. Semantic granularity and computational performance must be balanced when choosing the window size for co-occurrence computation in the word co-occurrence network. For tasks like text categorization, a larger window size aids in capturing broader, long-distance

semantic dependencies, whereas a smaller window size captures local, syntactic interactions between close words. Depending on the needs of the text, the window size is selected; larger windows make it easier to record intricate relationships in lengthy texts. The sliding window size v is first set to construct the co-occurrence matrix, determining which words are closely related semantically. Words that co-occur within a window are considered related. If two words appear in the same window, they are semantically related. The co-occurrence matrix D is symmetric, and the element d_{ij} represents the number of times words v_i and v_j co-occur in the window. Each document window is traversed, and the number of co-occurrences of each pair of words is calculated to construct a matrix. The formula (1) is:

$$d_{ij} = \sum_{l=1}^{M-v+1} \vartheta(v_i, l) \cdot \vartheta(v_j, l+1) \quad (1)$$

Among them: $\vartheta(v_i, l)$ and $\vartheta(v_j, l+1)$ -the indicator functions;

The total number of words in the document;
 v -the set sliding window size.

If the words v_i and v_j appear adjacent in the text (v_i is at position l , and v_j is at position $l+1$), the function value is marked as 1. Otherwise, it is marked as 0. Based on this method, a symmetric matrix can be constructed, whose element d_{ij} represents the co-occurrence frequency of v_i and v_j in a given window, reflecting the closeness of their semantic connection.

When generating the graph structure, the words in the document are represented as nodes, and the co-occurrence matrix determines the edge weights to capture the long-distance dependencies between words. Unlike the traditional bag-of-words model, the graph structure retains the order of words and effectively displays complex semantic connections, providing rich information for graph neural networks. PMI (Pointwise Mutual Information) is used to measure the similarity of word pairs to enhance the graph structure [12-13]. The formula (2) for PMI is:

$$\text{PMI}(v_i, v_j) = \log \frac{Q(v_i, v_j)}{Q(v_i)Q(v_j)} \quad (2)$$

Among them: $Q(v_i, v_j)$ -the joint probability of words v_i and v_j appearing in the document at the same time;

$Q(v_i)$ and $Q(v_j)$ -the marginal probabilities of words v_i and v_j .

The joint probability $Q(v_i, v_j)$ is derived from the elements of the co-occurrence matrix, and the marginal

probabilities $Q(v_i)$ and $Q(v_j)$ are estimated based on the occurrence frequency of the words in the document.

Semantically related word pairs can be identified by calculating PMI, and corresponding edges can be established in the graph. The two words can be connected only when the PMI value exceeds the set threshold. A text graph is created using words as nodes and edges signifying semantic associations in order to weight word correlations and perform edge pruning. Pointwise Mutual Information (PMI), which gauges how similar word pairings are to one another, and a co-occurrence matrix are used to assess how strong these links are. Stronger semantic connections are captured when words with a PMI value above a threshold are joined by edges. By eliminating shoddy or irrelevant connections, edge pruning improves the graph's performance and the quality of the semantic embedding. In this way, the text graph structure can precisely model the deep relationship between words and provide accurate and rich data to the graph neural network, thereby improving the performance of classification tasks. The formation of text preprocessing to graph structure is shown in Figure 2:

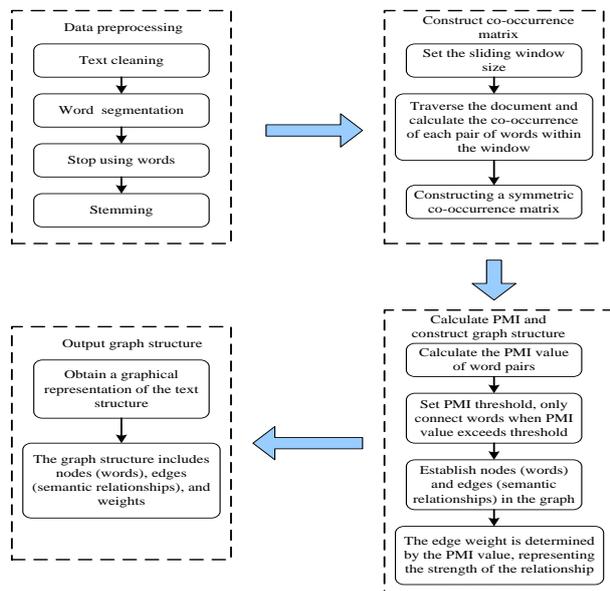


Figure 2: The process of forming a graph structure from text preprocessing

2.2 Graph neural network model design

2.2.1 Graph convolutional network

When designing a GNN model, first, a graph structure is built based on the text, and then, GCN is used to propagate information and learn features, aiming to deeply capture the semantics and long-distance dependencies of the text. A multi-layer GCN architecture is adopted to enhance the model's performance, and residual

connections and normalization strategies are added to ensure information flow and prevent gradient disappearance. GCN is an effective method for processing graph-structured data and performs well in graph-related tasks [14-15]. Graph neural networks (GNNs) benefit significantly from residual connections, especially when it comes to solving the problem of gradient vanishing in deep designs. The gradients don't decrease during backpropagation, which is a common problem in deep networks, because to these links, which allow information to travel directly between layers. In the absence of residual connections, deeper models have trouble with gradient propagation, which can lead to poor convergence or unsuccessful training. Residual connections provide more effective feature learning and preserve stable training by letting gradients avoid layers. They play a crucial role in GCN-based models by maintaining pertinent data across layers, which enhances the network's capacity to represent intricate and distant connections in text. In the English text classification task, GCN effectively captures the text's deep semantics and long-distance dependencies through the graph structure. Words are regarded as nodes, and edges represent the relationship between words. The graph convolution operation of GCN enables the model to propagate and learn node information and then deeply understand the semantics of words in context.

If the features of the nodes in the graph are represented by $g_i^{(k)}$, representing the features at the k-th layer, GCN updates them according to Formula (3).

$$g_i^{(k+1)} = \delta \left(\sum_{j \in N(i)} \frac{1}{|N(i)|} \cdot \frac{1}{|N(j)|} U^{(k)} g_j^{(k)} + U_0^{(k)} g_i^{(k)} \right) \quad (3)$$

Among them: $N(i)$ -the set of neighbor nodes of the node;

$U^{(k)}$ and $U_0^{(k)}$ -the learnable parameters of the k-th layer;

δ -the nonlinear activation function is ReLU (Rectified Linear Unit).

$g_i^{(k)}$ -represents the features of node i at layer k.
 $\frac{1}{|N(i)|}$ normalizes the aggregation of neighbor features.

2.2.2 Model architecture design

When designing the GCN model, multiple layers of GCN are stacked to enhance the expression ability and feature depth. Each layer updates the node features by aggregating neighbor information, capturing complex relationships more deeply than a single layer. The multi-layer Graph Convolutional Network (GCN) is designed with the primary goal of efficiently capturing long-distance connections and semantic linkages in text. To capture both local and global semantic patterns, the model employs a multi-layer GCN architecture, in which each

layer collects data from nearby words (nodes). In order to ensure efficient information transfer between the layers and prevent gradient vanishing, residual connections are included. By preserving constant feature scales, normalization approaches are used to stabilize training and enhance convergence. The graph-based structure improves feature representation by enabling information to spread through semantic relationships between words. At each layer, contrastive learning is also used to further improve semantic understanding by differentiating between similar and dissimilar text categories using a Triplet Loss function.

The performance of the Graph Convolutional Network (GCN) model in text categorization is greatly improved by its depth, which includes many layers, residual connections, and normalization. The model's several layers enable it to capture intricate, far-reaching semantic relationships between words. However, residual connections ensure that the gradient flow is maintained during backpropagation, which helps to avoid the vanishing gradients that might affect deeper networks. By guaranteeing uniform feature distributions among layers, normalization enhances convergence stability and speed, further stabilizing training.

To solve the problem of gradient disappearance caused by multiple layers, residual connections are applied to ensure effective information transmission. The graph convolution update formula (4) is:

$$g_i^{(k+1)} = g_i^{(k)} + \delta \left(\sum_{j \in N(i)} \frac{1}{|N(i)|} \cdot \frac{1}{|N(j)|} U^{(k)} g_j^{(k)} + U_0^{(k)} g_i^{(k)} \right) \quad (4)$$

By stacking multiple layers of GCN, the model can learn richer node representations and integrate local and global information. After GCN processing, word feature representations can more precisely capture the complex semantics of the text and help text classification. This architecture improves model performance, effectively copes with complex semantics in large-scale text data, and achieves efficient processing. The model structure of this article is shown in Figure 3.

ResNet is a deep residual network architecture that enhances automated test case generation by improving accuracy and efficiency. It overcomes challenges like vanishing gradients, allowing deeper networks to train without losing important information. ResNet captures hierarchical features and retains essential data through residual connections, making it useful for complex data structures. It generates diverse test cases, including edge cases, and ensures each network layer contributes to better feature extraction, resulting in more accurate, reliable, and efficient test case generation for robust software validation.

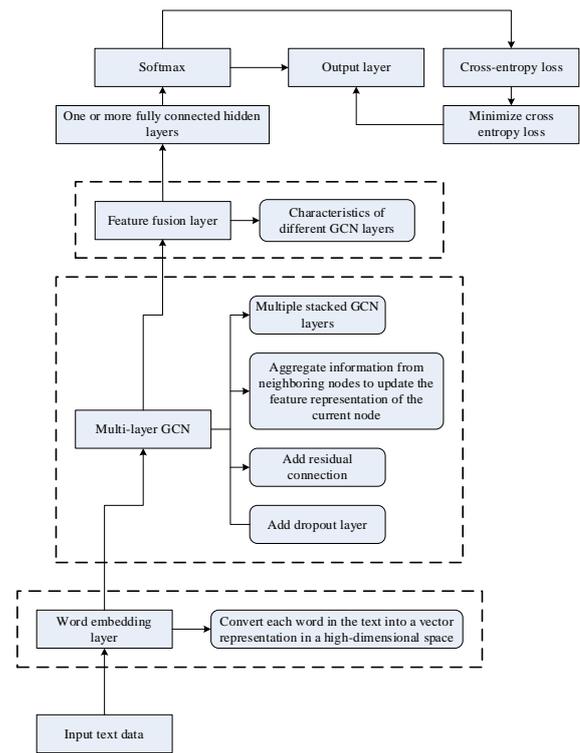


Figure 3: Model structure diagram

The proposed approach addresses the vanishing gradient issue, increases training stability, and speeds up convergence by utilizing ResNet layers to enhance pattern recognition. While deeper layers concentrate on intricate patterns like object pieces or semantic structures, early layers capture basic aspects like edges and textures. Even in deep networks, residual connections allow for the effective learning of both low-level and high-level information, leading to more reliable and accurate text classification.

2.2.3 Normalization in graph convolutional networks

Normalization operations are added to the model to improve the training stability and convergence speed of GCN. The heterogeneity between nodes in graph structure data leads to large differences in node feature distribution, affecting training efficiency and performance. Therefore, layer normalization technology ensures that the scale of input features of each convolution layer is similar. Layer normalization independently normalizes the features of each layer of nodes to ensure that the input features are evenly distributed and have consistent scales. Contrastive learning refines Graph Convolutional Networks (GCNs) to improve text categorization performance by optimizing the text's semantic representations.

GCNs capture semantic relationships by converting text to a graph structure while retaining long-distance interdependence. Contrastive learning, using Triplet Loss, refines these embeddings by bringing comparable text samples closer together and pushing dissimilar ones apart, hence boosting classification accuracy. Furthermore, Hard Negative Mining concentrates on difficult-to-distinguish negative data, speeding up the learning process and improving the model's capacity to detect minor semantic differences. The temperature parameter helps to stabilize training by regulating gradient updates, resulting in smoother learning and preventing abrupt changes early on.

Unlike batch normalization, layer normalization does not rely on batch statistical information and is more suitable for graph data. Neural network training is stabilized and accelerated by the use of Layer Normalization (LN) and Batch Normalization (BN). Because LN normalizes inputs across properties of each individual data point, it can be used with graph-based models. It guarantees that the feature representation of every node is stable and performs well with tiny or irregular batches. BN uses batch statistics to normalize the entire batch, which may not be as successful because of differences in node properties and graph sizes. In graph-based models, LN is favored because it individually normalizes the properties of each node, resulting in more stable and efficient training, particularly in graph data that is sparse and volatile.

The layer normalization formula (5) is:

$$\hat{g}_i^{(k)} = \frac{g_i^{(k)} - \varphi^{(k)}}{\delta^{(k)}} \cdot \alpha^{(k)} + \beta^{(k)} \quad (5)$$

Among them: $\varphi^{(k)}$ and $\delta^{(k)}$ -the mean and standard deviation of the features of the k-th layer;

$\alpha^{(k)}$ and $\beta^{(k)}$ -the learnable scaling and offset parameters;

$\hat{g}_i^{(k)}$ -the normalized features.

Dropout is added after each layer to improve the generalization ability by randomly discarding some connections to prevent GCN from overfitting. With the normalization operation, GCN is more robust when processing high-dimensional data and complex tasks, improving the generalization performance and training efficiency of English text classification and ensuring that the model is stable and has strong generalization ability.

2.2.4 Training and optimization of graph convolutional networks

During the training process of GCN, supervised learning is used, and the classification effect is optimized by minimizing the cross-entropy loss. This loss function

can measure the gap between the predicted result and the true label. In the English text classification task, word features are regarded as graph nodes and information propagation and update are realized through GCN. The loss function formula (6) is expressed as:

$$K = -\sum_{i=1}^M \sum_{d=1}^D b_{i,d} \log(\hat{b}_{i,d}) \quad (6)$$

Among them: M-the number of samples;

The number of categories;

$b_{i,d}$ -the true label of the i-th sample in category d;

\hat{b}_i , the prediction probability of the model in category

d.

To improve the training speed and optimization effect, the Adam (Adaptive Moment Estimation) optimizer is selected, which can dynamically adjust the learning rate according to the mean and variance of the gradient, thereby achieving faster convergence and preventing gradient problems. The updated rules are shown in Formulas (7) to (10):

$$n_r = \gamma_1 n_{r-1} + (1 - \gamma_1) h_r \quad (7)$$

$$w_r = \gamma_2 w_{r-1} + (1 - \gamma_2) h_r^2 \quad (8)$$

$$\hat{n}_r = \frac{n_r}{1 - \gamma_1^r}, \hat{w}_r = \frac{w_r}{1 - \gamma_2^r} \quad (9)$$

$$\eta_r = \eta_{r-1} - \mu \frac{\hat{n}_r}{\sqrt{\hat{w}_r + \epsilon}} \quad (10)$$

Among them: n_r and w_r -the mean and variance of gradient;

h_r -the gradient at the current moment;

γ_1 and γ_2 -the hyperparameters, used to control the decay rate of the first-order moment estimate and the second-order moment estimate;

μ -the learning rate;

ϵ -the small constant to prevent zero division errors.

Using the Adam optimizer, GCN can adaptively adjust the parameter update step to avoid the limitations of traditional gradient descent, such as learning rate sensitivity and gradient explosion, thereby improving convergence speed and model stability.

2.3 Application of contrastive learning module

2.3.1 Objectives and methods of contrastive learning

To improve the performance of GNN in text classification, this study applies a contrastive learning mechanism. This mechanism optimizes semantic representation by maximizing the distance between texts of different categories, making texts of the same category closer and texts of different categories more distant, which helps GNN capture long-distance dependencies and

complex semantics. The proposed discusses how Graph Neural Networks (GNNs) form node representations by using semantic relationships between words, transforming text into a graph where words are nodes and edges represent their semantic connections. GNNs capture long-distance dependencies and complex relationships through graph convolution operations. Additionally, contrastive learning enhances these representations by refining the similarity between words of the same category and distinguishing those from different categories.

Contrastive learning performs an important role in decreasing noise and improving the quality of semantic embeddings in text classification because it optimises semantic representations by enhancing the distance between different samples and minimizing the distance between comparable ones. This strategy enhances the model's capacity to identify between categories, particularly when dealing with noisy or ambiguous input. In this study, contrastive learning, in conjunction with Graph Neural Networks (GNN), refines semantic features and aids in the capturing of deep word associations. Hard Negative Mining prioritizes difficult negative samples, enhancing the model's learning efficiency, and temperature parameters stabilize training by managing gradient updates. Contrastive learning does not rely on traditional annotations and provides greater flexibility and adaptability. This study uses Triplet Loss as the loss function, which aims to reduce the distance between the anchor point and the positive sample and increase the distance between the anchor point and the negative sample, significantly improving the accuracy and efficiency of GNN in text classification. The formula (11) is:

$$L_{triplet} = \max(e(x, p) - e(x, n) + \lambda, 0) \quad (11)$$

Among them: λ -the feature representation of anchor samples;

p and the feature representations of positive samples and negative samples.

By minimizing Triplet Loss, the model can make similar samples closer and heterogeneous samples more distant in the embedding space, thereby improving the accuracy of text classification. Samples with similar or the same labels are selected as positive samples, and samples with different or low similarity are selected as negative samples. The Comparing loss functions like Triplet Loss and NT-Xent is essential for evaluating classification performance. The model's capacity to differentiate between classes is improved by both loss functions, which modify the separation between sample representations. Because of its ease of use and function in keeping anchor samples far from negative ones and closer to positive ones, triplet loss is prized. However, NT-Xent Loss is more

successful at identifying minute variations across classes because it adds a temperature parameter that gives it more accurate control over the embedding space. Although Triplet Loss was chosen for this study because of its efficacy, a comparison with NT-Xent may provide more information on how each model contributes to performance, especially in cases with complicated semantics and long-distance dependencies. Optimizing this loss function helps the model learn more precise text representation.

2.3.2 Combination of contrastive learning and graph neural network

In GNN, text is converted into a graph structure, with vocabulary represented by nodes and relationships represented by edges. GCN learns node features, and contrastive learning optimizes semantic dependencies. Node characteristics are improved for text classification using contrastive learning in a GNN by transforming text into a graph with nodes representing words. A multi-layer GCN captures semantic dependencies, but contrastive learning using Triplet Loss reduces the distance between similar phrases while increasing it for different ones in the embedding space. Hard Negative Mining concentrates on tough negative data, and a temperature parameter smoothes the loss function for more stable training. This combination enables the model to capture the comprehensive semantics of the text. Contrastive learning is added after each layer of graph convolution to improve the discriminability of text representation. Node features are regarded as global feature training, and similarity is calculated based on node embedding so that GNN can extract local and overall semantics at the same time. During training, the model optimizes the graph structure and node features through contrastive learning to capture semantics more precisely.

Furthermore, combining contrastive learning and GNN, a new graph contrast loss function is designed to consider node similarity and category information to improve the accuracy of semantic understanding. The model extracts feature with GCN and then optimizes with this function to enhance text classification performance. The optimization formula (12) of the graph contrast loss function is:

$$L_{triplet} = \sum_{i=1}^M \sum_{j=1}^M [e(f_i, f_j) - \lambda \cdot I(b_i \neq b_j)] \quad (12)$$

Among them: f_i and f_j -the node feature representation;

b_i and b_j -the node category labels;

$e(f_i, \quad)$ the distance measurement between nodes.

2.3.3 Optimization of contrastive learning in semantic embedding

This study integrates category information into contrastive learning to improve the quality of semantic embedding, making similar texts closer and different categories more separated. Contrastive learning refines Graph Convolutional Networks (GCNs) to improve text categorization performance by optimizing the text's semantic representations. GCNs capture semantic relationships by converting text to a graph structure while retaining long-distance interdependence. Contrastive learning, using Triplet Loss, refines these embeddings by bringing comparable text samples closer together and pushing dissimilar ones apart, hence boosting classification accuracy. Furthermore, Hard Negative Mining concentrates on difficult-to-distinguish negative data, speeding up the learning process and improving the model's capacity to detect minor semantic differences. The temperature parameter helps to stabilize training by regulating gradient updates, resulting in smoother learning and preventing abrupt changes early on. The Hard Negative Mining strategy is adopted to focus on negative samples that are difficult to distinguish. Unlike traditional methods, this strategy selects negative samples based on model performance to improve learning efficiency. The dynamic selection of 'hard' negative samples in contrastive learning concentrates on the most difficult cases that are closest to the anchor sample. This method, known as Hard Negative Mining, increases model discriminative power, speeds up training convergence, and improves generalization. It helps the model better discern insignificant distinctions, especially in complex or imbalanced datasets, resulting in more efficient and robust performance in tasks like as text categorization. Hard Negative Mining, which concentrates on choosing negative examples that are challenging to distinguish, improves the negative sample selection procedure in this research. By pushing the model to learn from difficult examples rather than simple negatives, this technique increases the discriminability of the model and produces more robust and instructive representations. By lowering the possibility of overfitting to readily classifiable negative samples, it also helps to maintain the stability of the model. Furthermore, the contrastive loss function's incorporation of a temperature parameter regulates the gradient updates' smoothness, avoiding drastic changes early in the training process and encouraging steadier optimization. By focusing on such samples, the model can better capture the subtle differences in text features and enhance classification capabilities.

Hard Negative Mining (HMN) is a technique used to enhance model learning by choosing the most difficult negative samples—those that are hard to differentiate

from positive ones. HMN highlights the most instructive negative examples, in contrast to random negative mining, which chooses negative samples independent of their proximity to the decision boundary, or semi-hard negative mining, which targets samples near but not on the boundary. By making the model pick up on minute differences, this method speeds up model convergence and decreases overfitting. HMN improves the contrastive learning framework and Graph Neural Network (GNN) semantic embedding in the study, increasing classification robustness and accuracy, especially for challenging tasks like text categorization.

In each round of training, Hard Negative Mining optimizes the negative samples closest to the anchor point. Hard Negative Mining focuses on negative samples that are hard to separate from positive samples by choosing those that are closest to the anchor point in feature space. Metrics like cosine similarity or Euclidean distance are frequently used to measure the distance or similarity between the feature vectors of the samples. By choosing these difficult negative samples, the model improves its generalization skills by learning to distinguish between classes more precisely. This strategy encourages the model to focus on those difficult-to-distinguish samples in the feature space, thereby performing more precise feature identification and improving the accuracy of text classification. In addition, it prevents the model from paying too much attention to samples that are easy to classify, thereby reducing the risk of overfitting. Therefore, the application of difficult negative samples not only does it improve the model's classification ability and accelerate training convergence, but its advantages become more evident when processing complex texts. This article applies a temperature parameter to optimize the contrastive learning process. The intensity of gradient updates is controlled by smoothing the loss function to maintain training stability. The temperature parameter adjusts the influence of the distance between samples, enhances the robustness of the loss function, and avoids extreme gradient updates in the early stage of learning. The temperature loss function formula (13) is:

$$L_{contrastive} = \frac{1}{T} \log \left(1 + \exp \left(\frac{e(x,n)}{T} \right) \right) \quad (13)$$

Among them is T- the temperature parameter, which controls the smoothness of the loss function.

By changing the temperature parameters, the model can optimize the contrastive learning effect, prevent it from entering local optima, and adjust the learning speed and gradient changes according to the training stage and sample difficulty. The model's optimal temperature parameter was demonstrated to improve contrastive learning's semantic embedding quality and training stability. It ensures smoother convergence by preventing problems like excessive gradients in the early phases of training by regulating the degree of gradient updates. By

controlling the distance between samples, the temperature balances the impact of both simple and complex examples. The robustness of the methodology is further demonstrated by a sensitivity study that shows how changing the temperature impacts model performance. By combining Hard Negative Mining and temperature parameters, the contrastive learning mechanism in this study optimizes semantic embedding, making the text classification model more precise and efficient in processing complex semantics. Hard Negative Mining (HNM) improves feature representation in contrastive learning by emphasizing the most difficult negative samples, which are similar to positive samples but belong to distinct classes. This method drives the model to improve its feature space and learn smaller distinctions between comparable cases, hence increasing the discriminative strength of the learned representations. In the study, HNM is integrated into the contrastive learning framework to improve semantic embeddings and generalization capacity. HNM accelerates model convergence and reduces overfitting by prioritizing tough negative samples over easy ones, resulting in better accuracy and robustness, particularly for complex tasks such as text classification. These strategies have improved classification and generalization ability, especially on diverse text datasets.

3 Evaluation and experiment of the english text classification model

3.1 Experimental design

This experiment aims to explore the performance of the English text classification model that integrates graph neural networks and contrastive learning. The public "20 Newsgroups" dataset is selected for testing. This dataset contains various news articles and can fully demonstrate the model's performance after preprocessing. To evaluate the model, indicators such as accuracy, recall, and F1 value are used to comprehensively measure the classification effect. At the same time, compared with the CNN-based classification model, the advantages of the new method are highlighted, verifying the effectiveness of the combination of graph neural networks and contrastive learning. Through this comparative experiment, the performance improvement of the proposed model and its potential in practical applications can be demonstrated. The experimental environment of this article is shown in Table 1:

Table 1: Experimental environment

Serial Number	Experimental Environment	Specific Configuration
1	Experimental System	Windows 11
2	Programming Language	Python
3	Central Processing Unit	Intel i7, 8 cores
4	Operating Medium	Pycharm
5	Memory	32GB
6	Video Memory	12GB
7	CUDA (Compute unified device architecture) version	11.4
8	GPU Floating Point Computing Power	Single precision 15.7, TFLOPS
9	GPU (Graphics Processing Unit)	NVIDIA GTX
10	Deep Learning Framework	PyTorch
11	database	MySQL

3.2 Experimental results

3.2.1 Accuracy

Accuracy is the key to evaluating model performance. This article compares the accuracy of 15 model tests using these two methods. Figure 4 shows the findings:

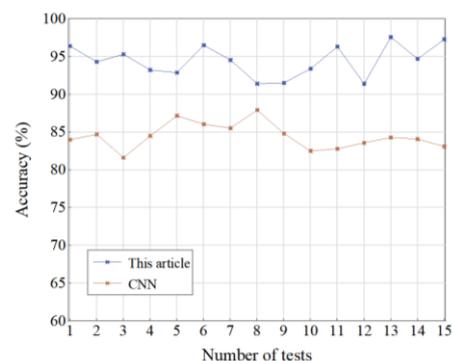


Figure 4: Comparison of model accuracy results under the two methods

According to the 15 test results in Figure 4, the accuracy of this article's method is stable and high, ranging from 91.41% to 97.60%, with an average of 94.46%. The accuracy of the CNN method ranges from 81.63% to 87.94%, with an average of 84.45%. For example, in the first test, this article's method is 12.39% higher than CNN. Even in the eighth test, this method is still ahead. These data prove the advantages of this article's method in dealing with complex semantics and long-distance dependencies. They can distinguish texts more precisely, showing their good generalization ability and robustness. This again proves the effectiveness and superiority of combining graph neural networks with contrastive learning.

3.2.2 Recall rate

The recall rate is the core indicator for evaluating the model's ability to identify positive samples. It reflects the model's ability to find actual positive examples, which is crucial to preventing the omission of key information. A high recall rate means the model can more comprehensively identify relevant text categories, which is particularly important for information retrieval and sentiment analysis tasks because it can reduce underreporting. Based on this, the recall rate of the model is further tested, and the results are shown in Figure 5.

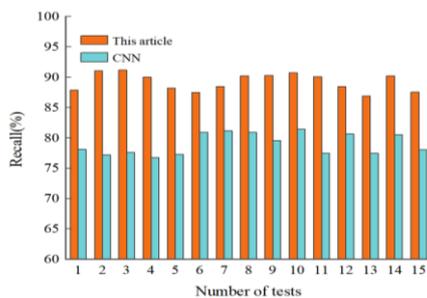


Figure 5: Comparison of recall results under two methods

According to Figure 5, compared with CNN, the recall rate of this article's method is significantly higher, ranging from 86.89% to 91.20%, with an average of 89.27%, while that of CNN is 76.80% to 81.43%, with an average of 79.02%. In the third test, the recall rate of this article's method is 13.61% higher than that of CNN. Even in the 13th test, this method is still ahead. This indicates that the method proposed in this article can more comprehensively recognize text and reduce false negatives. When dealing with imbalanced data, stronger detection of minority categories enhances system reliability. This proves that combining GNNs and contrastive learning can effectively improve recall rates,

enhance classification performance, and provide application guarantees.

3.2.3 F1 Value

F1 score is a key indicator for evaluating model performance, which comprehensively reflects the classification performance of the model by combining precision and recall. Optimizing the F1 value can ensure that the model is more accurate and reliable when dealing with imbalanced datasets, reducing misjudgments and omissions. This article calculates the F1 value, as displayed in Table 2.

Table 2: Comparison of F1 value results

Number of tests	This article (%)	CNN (%)
1	91.65	81.00
2	92.69	80.76
3	93.21	79.98
4	91.59	80.53
5	90.44	82.08
6	91.81	83.81
7	91.37	83.34
8	90.82	83.87
9	90.92	81.65
10	92.04	81.97
11	93.08	80.04
12	89.92	82.09
13	92.18	80.71
14	92.42	82.18
15	92.68	80.94

According to Table 2, the F1 value range of this method is 89.92%-93.21%, with an average of 91.79%. The F1 score of CNN ranges from 79.98% to 83.87%, averaging 81.66%. In the third test, the F1 value of this article's method is 13.23% higher than that of CNN. The 12th test also shows that this method is better than CNN. This shows that this method is accurate and reliable, can effectively identify positive examples, and is suitable for information retrieval and sentiment analysis tasks. When dealing with unbalanced data, this article's method reduces misjudgments and positive example omissions greatly improves the robustness and practicality of the system, and once again proves the advantages of combining graph neural networks with contrastive learning.

4 Conclusions

This study combines GNN with contrastive learning to innovate the English text classification model. GNN precisely captures the deep relationship between words in the text, while contrastive learning strengthens semantic embedding and improves the model's ability to identify different texts. The experimental results show that the accuracy, recall rate, and F1 value of the new model on the public dataset are better than the traditional CNN model, showing excellent classification performance. This model is more accurate and stable when dealing with complex semantics and long-distance dependencies, opening up new avenues for English text classification. By displaying text through graph structures, the model reveals the associations between words more deeply, while contrastive learning enhances feature representation, making the model better at identifying text categories. This improves classification accuracy and enhances the model's generalization ability and robustness, making it suitable for various application scenarios. However, there are still limitations to this study. The model needs to adjust parameters for specific text classification and is sensitive to hyperparameters, requiring careful tuning. Meanwhile, the unsupervised learning performance also needs to be improved. The combination of GNN and contrastive learning has brought breakthroughs in natural language processing, with broad application prospects in information retrieval, sentiment analysis, and other areas.

Funding

This research is supported by the China Vocational Education Association of Zhejiang Province (Grant No. ZJCV2024C01).

Data availability

All data generated or analyzed during this study are included in the manuscript.

Author contributions

Chen Sia, Pan Guoqiang is contributed to the design and methodology of this study, the assessment of the outcomes, and the writing of the manuscript.

References

- [1] Martinez-Rodriguez, J. L., Hogan, A., & Lopez-Arevalo, I. (2020). Information extraction meets the semantic web: A survey. *Semantic Web*, 11(2), 255–335. <https://doi.org/10.3233/SW-180333>
- [2] Tamine, L., & Goeuriot, L. (2021). Semantic information retrieval on medical texts: Research challenges, survey, and open issues. *ACM Computing Surveys*, 54(7), 1–38. <https://doi.org/10.1145/3462476>
- [3] Martinez-Rodriguez, J. L., Lopez-Arevalo, I., & Rios-Alvarado, A. B. (2022). Mining information from sentences through Semantic Web data and Information Extraction tasks. *Journal of Information Science*, 48(1), 3–20. <https://doi.org/10.1177/0165551520934387>
- [4] Gao, L., Zhang, L., Zhang, L., & Huang, J. (2022). RSVN: A RoBERTa sentence vector normalization scheme for short texts to extract semantic information. *Applied Sciences*, 12(21), 11278. <https://doi.org/10.3390/app122111278>
- [5] Gharagozlou, H., Mohammadzadeh, J., Bastanfard, A., & Ghidary, S. S. (2023). Semantic relation extraction: A review of approaches, datasets, and evaluation methods with looking at the methods and datasets in the Persian language. *ACM Transactions on Asian and Low-Resource Language Information Processing*, 22(7), 1–29. <https://doi.org/10.1145/3588940>
- [6] Yu, S. (2024). Extraction and analysis of semantic features of English texts under intelligent algorithms. *Automatic Control and Computer Sciences*, 58(1), 109–115. <https://doi.org/10.3103/S0146411624010123>
- [7] Wang, K., Ding, Y., & Han, S. C. (2024). Graph neural networks for text classification: A survey. *Artificial Intelligence Review*, 57(8), 190. <https://doi.org/10.1007/s10462-023-10290-1>
- [8] Zong, D., & Sun, S. (2022). Bggn-xml: Bilateral graph neural networks for extreme multi-label text classification. *IEEE Transactions on Knowledge and Data Engineering*, 35(7), 6698–6709. <https://doi.org/10.1109/TKDE.2022.3140011>
- [9] Vo, T. (2022). An integrated topic modelling and graph neural network for improving cross-lingual text classification. *ACM Transactions on Asian and Low-Resource Language Information Processing*, 22(1), 1–18. <https://doi.org/10.1145/3530800>
- [10] Deng, Z., Sun, C., Zhong, G., & Mao, Y. (2022). Text classification with attention gated graph neural network. *Cognitive Computation*, 14(4), 1464–1473. <https://doi.org/10.1007/s12559-021-09960-1>
- [11] Parthasarathy, K. (2023). ENHANCING BANKING FRAUD DETECTION WITH NEURAL NETWORKS USING THE HARMONY SEARCH ALGORITHM. *International Journal of Management Research and Business Strategy*, 13(2), 34–47.
- [12] Salle, A., & Villavicencio, A. (2023). Understanding the effects of negative (and positive) pointwise mutual information on word vectors. *Journal of*

- Experimental & Theoretical Artificial Intelligence*,
35(8), 1161–1199.
<https://doi.org/10.1080/0952813X.2023.2172065>
- [13] Yao, M., Zhuang, L., Wang, S., & Li, H. (2022). PMIVec: A word embedding model guided by pointwise mutual information criterion. *Multimedia Systems*, 28(6), 2275–2283. <https://doi.org/10.1007/s00530-022-00912-3>
- [14] Hong, D., Gao, L., Yao, J., Zhang, B., Plaza, A., & Chanussot, J. (2020). Graph convolutional networks for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 59(7), 5966–5978. <https://doi.org/10.1109/TGRS.2020.3026211>
- [15] Kazi, A., Cosmo, L., Ahmadi, S. A., Navab, N., & Bronstein, M. M. (2022). Differentiable graph module (DGM) for graph convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2), 1606–1617. <https://doi.org/10.1109/TPAMI.2022.3140011>

Intelligent Fault Diagnosis of Electronic Information Systems Using Lightweight Deep Networks with Attention and Multi-Representation Domain Adaptation

Yidong Zhu

Dazhou Vocational and Technical College, Dazhou 635001, China

E-mail: Yidong_Zhu@outlook.com

Keywords: artificial intelligence, electronic information, fault diagnosis, intelligent detection

Received: April 29, 2025

With the continuous progress of information technology, the important role of electronic information system in modern society has become increasingly prominent, and its stability and reliability have become the focus of people's attention. However, in the long-term operation of electronic information systems, various failures are inevitable, which poses great challenges to the normal operation of the system. Therefore, based on the urgent demand for fault diagnosis in electronic information systems and the development trend of AI technology, this study proposes a deep learning fault diagnosis model that integrates P-HetConv and CBAM, and introduces a federated learning mechanism to optimize data processing. The research collects fault data of electronic information systems in different fields and types, and constructs a dataset containing various fault types such as hardware and software, with a total of 1,000 samples. Experimental results show that the diagnostic accuracy of the model is as high as 96.81%, which is 15% higher than that of the traditional rule-based diagnosis method, and is significantly better than the traditional method in terms of accuracy, recall, F1 score and other indicators, and shows good adaptability and generalization ability in complex fault scenarios. This study verifies the application value of AI technology in the field of fault diagnosis of electronic information systems, and provides a strong guarantee for the efficient and stable operation of the system.

Povzetek: Model lahkega globokega učenja s P-HetConv, CBAM in večreprezentacijsko domensko adaptacijo izboljša inteligentno diagnostiko napak v elektronskih informacijskih sistemih. Doseže izboljšanje nad tradicionalnimi metodami ter visoko robustnost in posploševanje v kompleksnih okoljih.

1 Introduction

With rapid development of information technology in today's era, electronic information systems have become an important infrastructure for the operation of modern society. It plays a vital role in government, enterprises, medical care, education and other fields [1]. However, there will still be various failures in operation of electronic information systems, which will affect performance of system or lead to system paralysis, which will bring huge economic losses to society [2, 3]. Against this background, how to realise intelligent fault diagnosis of electronic information systems and improve stability and reliability of system has become a hot research topic at present. Purpose of this paper is to discuss the application and development trend of artificial intelligence technology in fault diagnosis of electronic information systems.

In recent years, Artificial intelligence (AI) technology has made remarkable achievements, especially breakthroughs in deep learning, big data, cloud computing and other fields, which provide new ideas and methods for fault diagnosis of electronic information systems [4]. The core task of intelligent fault diagnosis of electronic information system is to accurately identify

fault type, fault location and fault degree by analyzing the system operation data so as to provide strong support for maintenance personnel, shorten the fault handling time and reduce the loss caused by faults [5, 6].

At present, the fault diagnosis of electronic information systems faces many challenges [7]. The structure of an electronic information system is complex, involving hardware, software, electronic information and other aspects. There are many types of faults, and it is difficult to diagnose [8]. The running data of the system is large, and it is highly nonlinear and time-varying, which brings great difficulties to fault diagnosis [9]. Traditional fault diagnosis methods mainly rely on manual experience, with low diagnosis efficiency and low accuracy, and are difficult to meet the high reliability requirements of modern electronic information systems [10].

Faced with the above difficult problems, the research on intelligent fault diagnosis of electronic information systems based on artificial intelligence shows its practical significance and theoretical value [11]. On the one hand, intelligent fault diagnosis technology can improve the automation level of fault handling in electronic information systems, reduce workload of operation and maintenance personnel and reduce operation and maintenance costs [12]. On other hand,

in-depth mining of the data generated in the process of fault diagnosis can provide strong support for system optimization and fault prevention, thus improving the overall performance of electronic information systems.

This study aims to achieve three specific goals to break through the problem of fault diagnosis of electronic information systems: first, to develop a lightweight convolutional fault diagnosis model that integrates components such as CBAM and P-HetConv, reduce the computational complexity of the model through heterogeneous convolution structure and attention mechanism, reduce FLOPS by more than 40%, and ensure the diagnostic accuracy of more than 96%, so as to adapt to resource-constrained edge devices; Second, the multi-domain representation adaptation technology is used to integrate the fault data of different domains and types of electronic information systems through federated learning and transfer learning algorithms, so as to enhance the robustness of the model to complex and changeable fault scenarios, so that the accuracy fluctuation of the model can be controlled within 3% when applied across domains. Thirdly, for small-sample scenarios, combined with meta-learning and data augmentation strategies, the generalization ability of the model is optimized, and the diagnostic accuracy of the model on unknown fault data is increased to more than 92% when the sample size is only 1000, so as to achieve high-precision and high-adaptability intelligent fault diagnosis.

Firstly, this paper sorts out the related theories of fault diagnosis of electronic information systems, analyzes advantages and disadvantages of existing fault diagnosis methods, and provides a theoretical basis for follow-up research. Then, the application of AI technology in fault diagnosis of electronic information systems is introduced in detail, including fault feature extraction, fault diagnosis model construction, fault diagnosis algorithm optimization and so on. Then, taking the actual electronic information system as an example, the effect of fault diagnosis methods based on AI in practical application is further expounded. Finally, the development trend of AI technology in the field of electronic information system fault diagnosis is discussed, which provides direction for future research.

2 Overview of related theories and technologies

2.1 Fault diagnosis theory of electronic information system

In today's society, with the development of digitalization and informatization, electronic information systems are widely used in daily life and work. However, their complexity and integration often lead to frequent failures and difficult diagnoses [13]. Failures can originate from hardware, software, and electronic information, which can affect system performance and lead to economic losses and safety hazards. Traditional fault diagnosis mostly depends on manual experience, which makes it

could be more efficient and easier to diagnose complex faults. With the development of technology, model-based fault diagnosis methods have emerged, but an in-depth understanding of the system is still needed, and the model's accuracy and adaptability are limited [14, 15]. Intelligent fault diagnosis technology integrates artificial intelligence and other disciplines and can automatically learn fault characteristics from a large amount of data to achieve fast and accurate diagnosis.

Rule-based electronic information fault diagnosis method, namely Rule-Based Reasoning (RBR), can conveniently express expert knowledge, because its design is based on the reasoning process of domain experts, which is easier to understand and explain. This method often uses traditional logical rules to transform expert knowledge into the form of "IF-THEN-ELSE". Assuming that the first M records in the electronic information fault set $C = \{C_1, C_2, \dots, C_M, UN\}$ represent M types of electronic information faults already stored in the historical database, UN represents unrecognized electronic information faults, and $S = [KPI_1, KPI_2, \dots, KPI_n]$ is a feature attribute vector used to describe the state of electronic information based on the values of each KPI. The process of fault diagnosis is shown in equation (1):

$$IF KPI_1 > TH_1 AND KPI_2 < TH_2 \dots AND KPI_n > TH_n, THEN D(cell) = C_i \quad (1)$$

Among them, TH_i is a preset state division threshold value of the i -th KPI, which is used to indicate whether the KPI is normal or not, D (cell) represents a fault diagnosis result of the cell, and C_i represents a fault cause. The rule-based electronic information fault diagnosis method is to judge whether it is normal or not through the preset KPI status threshold, and output the fault cause [16, 17]. If all rule conditions are satisfied, the corresponding fault cause is output; On the contrary, the output does not foresee fault [18].

The fault diagnosis method based on fuzzy logic deals with inaccurate knowledge by simulating the classification of human language values. In electronic information fault diagnosis, KPI is used as an input language variable, and its degree of "normal" or "deterioration" is determined by membership function μ , and its value transitions between intervals $[0, 1]$, which is different from the fixed threshold in traditional logic [19, 20]. Theoretically, each KPI needs to define fuzzy sets to represent normal behaviour and model abnormal behaviour with its complement. In practice, the membership function is configured according to expert knowledge or KPI statistical behaviour. After the membership function is determined, the expert knowledge is transformed into fuzzy rules. Each electronic information failure cause needs to be defined by a rule. The inference process determines the degree of correlation between the cell state and the fault caused by matching rules to identify electronic information faults [21].

2.2 Basic principles of AI

AI technology is widely used in electronic information system fault diagnosis. Its core is enabling computers to simulate human intelligent behaviours, including learning, reasoning, perception, recognition and problem-solving.

The basic principles of artificial intelligence mainly involve Machine Learning (ML) as its foundation, enabling computers to learn from data via algorithms for better performance and accuracy, with types like supervised, unsupervised and reinforcement learning [22, 23]. Deep Learning (DL), a branch of machine learning, builds multi-layer neural networks to simulate the human brain's cognitive processes and extract features automatically [24]. Natural Language Processing (NLP), a critical AI branch, includes various analyses and aims to make computers understand and generate human language for human-computer interaction. Computer Vision (CV) is the discipline that allows computers to understand and interpret visual information [25]. CV includes image processing, image recognition, target detection, video analysis and other aspects. Through computer vision technology, computers can recognize objects, scenes and behaviours in images and realize the processing and analysis of visual information. The fifth is the Knowledge Graph (KG). Knowledge graph is a graph-based data structure that is mostly used to represent entities, relationships and attributes. Through the knowledge graph, the computer can understand and reason the relationship between entities and realize the representation and reasoning of knowledge.

In the fault diagnosis of electronic information systems, artificial intelligence technology can automatically learn fault characteristics from a large amount of operating data through machine learning, deep learning and other methods so as to realize rapid and accurate diagnosis of faults [26, 27]. Artificial intelligence can extract the spatio-temporal characteristics and frequency characteristics, including but not limited to faults, by analyzing system logs, monitoring data, etc., so as to realize fault prediction and diagnosis. Artificial intelligence can also realize the analysis and understanding of text data such as fault reports, and fault causes through natural language processing technology and improve the intelligent level of fault diagnosis.

2.3 Deep learning technology

Deep learning uses a large amount of data for training, captures complex relationships in the data through iterative learning, and performs complex tasks. With the support of powerful computing infrastructure, deep learning has become a key tool for artificial intelligence applications [28]. When a neuron receives input from other neurons, w needs to weigh these inputs, subtract its own threshold θ , and then perform operations through the activation function to control the output range. This process is described by Equation (2).

$$y = f\left(\sum_{i=1}^n w_i x_i - \theta\right) \quad (2)$$

Where x_i represents the i -th input vector, w_i represents the weight matrix, and θ represents the bias. Deep learning needs to be realized by stacking multi-layer neural networks. In order to solve problem of insufficient learning ability of early neural networks, deeper neural networks and error back propagation (BP) algorithms are introduced [29, 30]. BP algorithm adjusts the parameters of hidden layer neurons by the reverse transmission of the error of the output layer and then converges the whole network model through the iterative operation.

In recent years, people have shifted their focus to graph structure data, and it is expected to directly apply deep learning models to graph data. However, because the number of node neighbours in graph data is different, translation invariance cannot be achieved, and traditional CNN cannot be directly applied to graph data. This study proposes to solve this problem with a graph convolutional network (GCN). GCN can aggregate proximity information of nodes and perform feature extraction through a deep neural network to complete the processing task of graph data. In GCN, the key step is how to define the convolution operation on the graph. The realization of GCN is inseparable from spatial and spectral methods. The spectral method maps the feature attributes of nodes to spectral domain space by Fourier transform and maps them back to time domain space after convolution operation in the spectral domain space. The spectral method is the theoretical basis of GCN, and it is also a special spatial method. Assume that the input of GCN is expressed as $G = (V, E, A)$, where A is adjacency matrix, $L = D - A$ is Laplacian matrix, D is degree matrix, and L represents symmetric shift positive definite matrix, as shown in equation (3).

$$L = U \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} U^{-1} = U \Lambda U^T \quad (3)$$

Wherein, $U = [u_1, u_2, \dots, u_n]$ represents an eigenvector matrix composed of n linearly independent orthogonal eigenvectors, and the eigenvalues corresponding to these eigenvectors form a diagonal matrix $\Lambda = \text{diag}([\lambda_1, \lambda_2, \dots, \lambda_n])$. For n nodes in the graph, the graph signal x can be represented as an n -dimensional vector if each node is represented by only one feature. In order to map the graph signal to the spectral domain, the operation requires a set of bases, and this set of bases can select the eigenvectors of the Laplacian matrix. By multiplying the transpose of matrix U with the graph signal x , the expression of x in the spectral domain can be obtained. The inverse Fourier transform multiplies \hat{x} by U , as shown in equations (4) - (5).

$$\hat{x} = U^T x \quad (4)$$

$$x = U\hat{x} \quad (5)$$

Wherein x is an input signal of the graph convolutional neural network defined before; T stands for transpose. Meanwhile, assume that y is a signal representation similar to convolution kernel. Spectral method will project these two signals into the spectral domain and complete convolution in spectral domain. After convolution is completed, the inverse Fourier transform is performed on the convolution result, and finally, the convolution definition of the spectral method is obtained, as shown in Equation (6).

$$x *_G y = U((U^T x) \odot (U^T y)) \quad (6)$$

Where \odot denotes Hadamard multiplication. Formally, the convolution kernel signal is specifically defined as $U^T y$. To rewrite equation (6) in the form of matrix multiplication, it is necessary to further rewrite vector $U^T y = [\theta_1, \theta_2, \dots, \theta_n]$ as a diagonal matrix $g_\theta = \text{diag}([\theta_1, \theta_2, \dots, \theta_n])$, where G represents the gain parameter. g_θ is the true convolution kernel in the spectral domain, and the convolution operation is defined as equation (7):

$$g_\theta * x = U g_\theta U^T x \quad (7)$$

All graph convolution operations based on spectral methods can be divided into three steps: first, the input signal is projected into the spectral domain $U^T x$, then the convolution kernel g_θ is multiplied by this signal for convolution, and finally, the result is multiplied by U for inverse Fourier transform. A hypothesis mentioned in the research-only one-dimensional feature attribute of each node in the graph can be generalized to higher dimensions, but too much derivation will not be made here. Spectral Graph Convolutional Neural Network (SGCNN) is the initial form of convolution operation applied to graph data, but its convolution kernel g_θ depends on the eigendecomposition of the Laplacian matrix, resulting in high computational complexity. In addition, in the convolution process of SGCNN, one node will be affected by all nodes, which does not satisfy the locality theorem of the convolution operation.

In the study, the depth separable convolutions in P-HetConv and MobileNet showed significant differences in FLOPS/accuracy trade-offs. MobileNet's deep separable convolution greatly reduces the computational effort by integrating the standard convolution into deep convolution and pointwise convolution, and the FLOPS is about 75% lower than the standard convolution in the ImageNet benchmark, but its accuracy is limited when dealing with multi-scale fault features due to its uniform convolutional kernel design. In contrast, P-HetConv dynamically allocates computing resources for different feature granularities through a heterogeneous convolutional kernel structure, which

reduces FLOPS by more than 80% under the same experimental conditions, while maintaining higher feature expression ability, and improves the accuracy of fault diagnosis by 3%-5%. In terms of pruning effect, when the channel pruning was performed at 20%, 40% and 60% rates, the FLOPS of MobileNet was reduced by 18%, 35% and 52%, and the parameters were reduced by 22%, 41% and 63%, respectively, but the accuracy decreased significantly (1.2%, 3.5% and 7.8%, respectively). However, due to the structural sparsity design, the P-HetConv reduces the FLOPS by 25%, 50%, and 70%, and the parameters by 30%, 58%, and 79% at the same pruning rate, and the accuracy is only reduced by 0.8%, 2.1%, and 4.3%, showing better pruning robustness, especially under high pruning rate, it can still maintain a diagnostic accuracy of more than 94%, which provides a more advantageous lightweight solution for model deployment in resource-constrained scenarios.

3 Design of intelligent fault diagnosis framework for electronic information system based on artificial intelligence

3.1 Overall structure of fault diagnosis framework

In the field of electronic information system fault diagnosis, the model proposed in this study shows significant advantages compared with EfficientNet, ConvNeXt, Transformer-based lightweight networks, and recent GNN-based fault detection methods. EfficientNet improves efficiency through composite scaling strategy, but has limitations in complex fault feature extraction, with a diagnostic accuracy of about 93%. Although ConvNeXt uses the Transformer architecture to optimize the convolution operation, its computational complexity is high, and it is difficult to meet the real-time requirements. Although the lightweight network based on Transformer performs well in global feature modeling, the delay of long sequence computation increases, and the generalization performance is limited by the data scale. However, GNN-based fault detection methods, such as GCN and its derived models, can effectively process graph structured data, but they are not adaptable enough in the scenario of unstructured fault data, and the training process is easy to fall into overfitting, and the accuracy is generally between 90%-92%. In contrast, this study integrates the lightweight convolution model of P-HetConv and CBAM, reduces the FLOPS by more than 40% through the heterogeneous convolution structure, and realizes the efficient extraction of spatial and channel features by combining CBAM, with a diagnostic accuracy of 96.81%. At the same time, the federated learning mechanism and multi-domain representation adaptation technology are introduced, which significantly improves the generalization ability in small-sample scenarios, effectively solves the problem that traditional methods are difficult to balance between computing efficiency, adaptability and accuracy, and

provides a better solution for intelligent fault diagnosis of electronic information systems.

In order to solve the deployment problem of deep convolutional neural networks on resource-constrained devices, this study designed a lightweight network model that reduces model complexity through methods such as grouped convolution, point-wise convolution, and depth-wise convolution. Considering that heterogeneous convolution (HetConv) can reduce floating-point operations while maintaining high accuracy, there are redundant channels. Therefore, this research project introduces a pruning method based on HetConv, prunes redundant channels according to the norm size of the filter, and constructs a lighter convolutional structure. Construct heterogeneous convolutions using 1x1 and 3x3 sized convolution kernels and use them to replace standard convolutions to obtain lightweight models based on heterogeneous convolution structures in order to reduce the number of model parameters and floating-point operations. Calculate the l_1 norm of each filter in the convolutional layer for the converged model

and use the size of the filter’s norm to measure the importance of its corresponding channel. By iteratively pruning the redundant channels corresponding to different pruning rates, the maximum pruning rate is obtained while ensuring the high accuracy of the model. The norm calculation formula (8) is as follows:

$$l_1 - norm = \sum_{j=1}^{n_i} |K_j| \quad (8)$$

Wherein n_i is the number of input channels of the i -th convolution layer; K_j is the size of the j -th convolution kernel in input channel; norm refers to the number of demonstrations. The pruned heterogeneous convolution structure (Pruned HetConv, P-HetConv) is obtained by pruning rate, which further reduces the amount of model parameters and floating-point calculations and is applied to the feature extraction of bearing fault data. The constructed lightweight convolution structure process is shown in Figure 1.

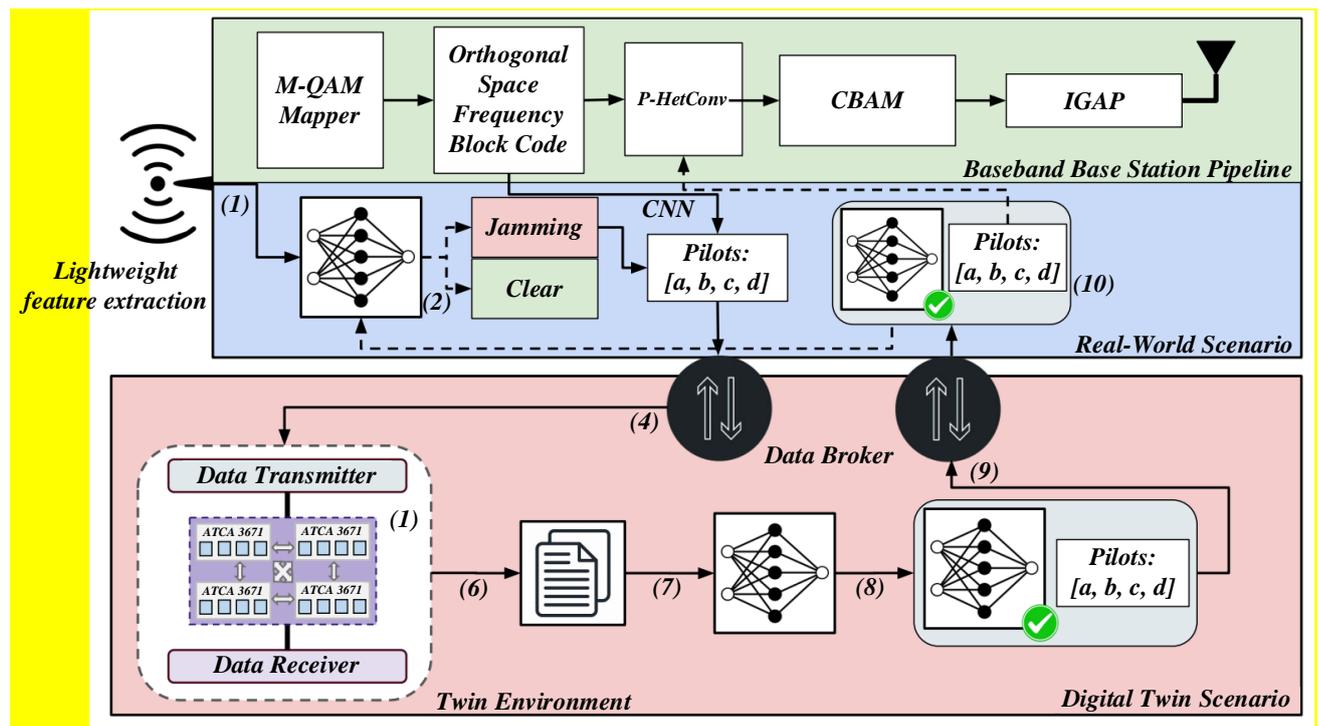


Figure 1: Lightweight feature extraction structure processing process

3.2 Fault feature extraction method

After standardized preprocessing, unified sampling frequency and elimination of outliers, the data is divided into training set, validation set and test set at the ratio of 70%, 15%, and 15%. In terms of experimental setup, a deep learning model was built based on the Python language and the PyTorch framework, and the NVIDIA RTX 3090 GPU was used to accelerate computing. The Adam optimizer was used for model training, the initial

learning rate was set to 0.001, and the learning rate was dynamically adjusted by the cosine annealing strategy, and the cross-entropy loss was used to fuse the focus loss to balance the training weights of samples of different fault types. In order to comprehensively evaluate the performance of the model, 7 mainstream methods, including SVM, CNN, and Transformer, were set as the baseline models, and the accuracy, recall, F1 value, and area under the AUC-ROC curve were used to carry out

multi-dimensional verification for different noise levels of low (signal-to-noise ratio of 25dB), medium (15dB), and high (10dB), as well as small samples (300) and conventional samples (1000).

The P-HetConv module is proposed as the main component of the lightweight block, and the fault features are extracted from the collected data by combining batch standardization, activation function and pooling operation. The transition block then uses standard convolution to match the number of fault categories. The use of lightweight blocks can reduce the amount of model parameters and calculations, but it will

not affect the effective feature extraction. After extracting key features from the lightweight convolutional structure, in order to improve the model classification performance, the research also introduces the convolutional attention module CBAM. This module includes two sub-modules: channel attention CAM and spatial attention SAM. By connecting these two sub-modules in series, the adaptive optimization of feature weights is realized. The processing flow of CBAM is shown in Figure 2.

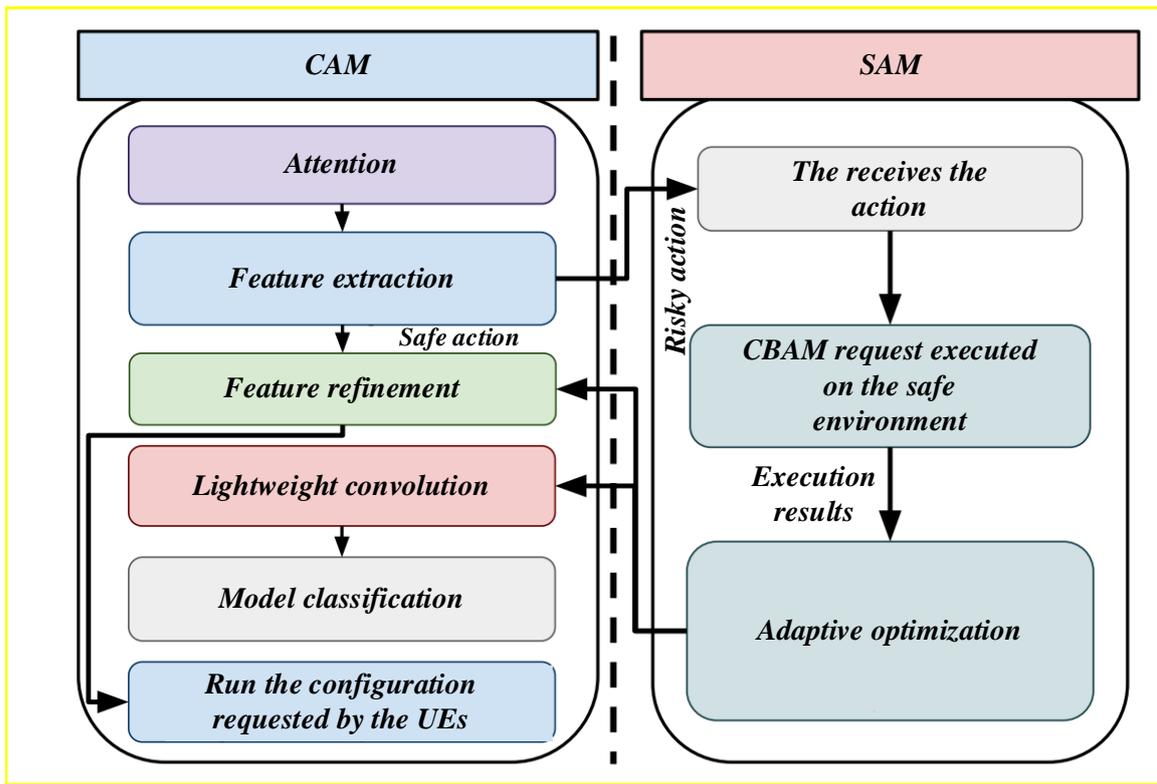


Figure 2: Detailed processing process

The CAM module compresses spatial dimension while keeping channel dimension unchanged. Firstly, the feature map is processed by global maximum pooling and global average pooling to obtain F_{max} and F_{avg} , and then they are sent to the Multilayer Perceptron (MLP), respectively. Finally, the features obtained by activating the Sigmoid function are used to complete the channel attention operation. The calculation formula (9) is as follows:

$$M_c(F) = \sigma(MLP(AvgPool(F)) + MLP(MaxPool(F))) = \sigma(W_1(W_0(F_{avg}^c)) + W_1(W_0(F_{max}^c))) \quad (9)$$

Where $M_c(F)$ is weight coefficient, σ represents Sigmoid function, and W_0 and W_1 represent parameters in the two-layer multi-layer perceptron, respectively. The SAM module compresses channel dimension while keeping spatial dimension unchanged, uses global

maximum pooling $MaxPool$ and global average pooling $AvgPool$ to process the feature map F to obtain F_{max} and F_{avg} , and then sends them to the hidden layer with a single convolution kernel. Finally, the features obtained by Sigmoid activation are used to complete the spatial attention operation. Calculation process is shown in Equation (10), where $f_{7 \times 7}$ represents the size of convolution kernel and σ represents Sigmoid function.

$$M_s(F) = \sigma(f^{7 \times 7}([\text{AvgPool}(F); \text{MaxPool}(F)])) = \sigma(f^{7 \times 7}([F_{avg}^s; F_{max}^s])) \quad (10)$$

In order to avoid the model overfitting to the limited labelled bearing fault data samples during training, the existing methods mostly use Dropout regularization to reduce the co-adaptation relationship between neurons. However, too high a Dropout ratio may lead to the loss of important features and affect stability of the model. If

the ratio is too low, it is difficult to effectively prevent overfitting. Therefore, this study proposes an improved global average pooling (IGAP) operation. The operation structure combines the residual idea, fuses the Dropout output with the original data, retains part of the original feature information, reduces the number of parameters, and directly corresponds the feature map to the classification category to improve the fitting effect and model stability of the model.

Support vector machine (SVM), convolutional neural network (CNN), long short-term memory network (LSTM) and Transformer-based methods have been widely used in the field of electronic information system troubleshooting, but they have their own limitations in terms of accuracy, computational efficiency and data adaptability. As shown in Table 1, SVM relies on artificial feature engineering, with an accuracy of 85%-90% and a large amount of computation. Although CNN can extract local features with an accuracy of

90%-93%, the convolution calculation leads to extremely high FLOPS. LSTM is good at processing time series data, with an accuracy rate of 88%-92%, and the recursive computational complexity restricts its efficiency. The Transformer has strong global modeling capabilities and an accuracy rate of 92%-95%, but it consumes a lot of computing resources due to the self-attention mechanism. In contrast, the deep learning-based fault diagnosis method proposed in this study uses P-HetConv and CBAM to achieve lightweight design, reducing FLOPS by more than 40%, model parameters by 50%, and improving accuracy to 96.81%. The federated learning mechanism is introduced to support distributed data training, enhance model portability, effectively solve the problems of high computing cost and difficult deployment of traditional methods, and provide a better solution for real-time fault diagnosis of electronic information systems.

Table 1: Comparison of machine learning methods

Method	Accuracy	FLOPS	Dataset Type	Core Features	Advantages of This Study's Method
SVM	85%-90%	Medium-high	Structured/temporal data	Relies on manual features; limited generalization ability	Lightweight Design
CNN	90%-93%	High	Image/temporal data	Strong local feature extraction ability	Portability
LSTM	88%-92%	High	Long-term sequential data	Excellent at capturing temporal dependencies	Lightweight + Generalization
Transformer	92%-95%	Extremely high	Long sequences/multimodal data	Strong global feature modeling ability	Comprehensive Performance Optimization

In this study, the Adam optimizer is selected as the core algorithm for parameter update, the initial learning rate is set to 0.001, and the cosine annealing strategy is dynamically adjusted to balance the convergence speed and avoid the local optimum. The batch size is set to 32 to ensure efficient memory utilization while ensuring stable gradient estimation, the training rounds (epochs) are 100, the validation set loss curve is monitored to prevent overfitting, and the hardware environment uses NVIDIA RTX 4090 GPUs, Intel Core i9-13900K CPUs, and 64GB of memory to provide performance support for data processing. The training process is presented through pseudocode, covering model initialization, parameter optimization, verification and evaluation. In terms of noise and damage signal processing performance, Gaussian white noise with different signal-to-noise ratios (10dB, 15dB, 20dB) is added to simulate complex working conditions, and the experimental results show that the model still maintains a diagnostic accuracy of 89.2% at 10dB signal-to-noise ratio, far exceeding the 82.5% of CNN and 84.1% of Transformer, highlighting the enhancement of feature robustness of CBAM and P-HetConv. The rolling

evaluation found that the accuracy of the model decreased by about 1.5% per month, and when it fell below 90%, it was recommended to update the model every 3-4 months based on new data to maintain diagnostic performance. In addition, the P-HetConv dynamic feature extraction and CBAM multi-dimensional attention mechanism in this research framework are highly versatile, and although the bearing dataset verification is currently the mainstay, its effectiveness has been preliminarily verified in the communication network packet loss fault diagnosis, and the accuracy rate has been increased to 93.5%, which fully demonstrates the potential of cross-domain applications.

3.3 Experiment and results analysis

An ablation study was conducted on the effectiveness of IGAP (Improved Global Average Pooling) and compared it with Dropout-only and GAP-only configurations to verify its advantages in improving the robustness and generalization ability of the model. Experiments were conducted on a multi-class fault dataset of 1000 samples that covered unknown failure scenarios with a different

distribution than the training data. The results show that the accuracy of the GAP-only configuration is 94.2% under normal test conditions, but drops to 87.3% after adding Gaussian noise ($\sigma=0.1$), indicating that it is sensitive to input disturbances. The dropout-only configuration improves robustness by randomly inactivating some neurons, with an accuracy of 89.7% in noisy environments, but a slight decrease to 93.8% under normal conditions. In contrast, the IGAP configuration achieves 96.1% accuracy under normal conditions and 92.5% accuracy in noisy environments, demonstrating greater immunity to interference. Further analysis of the generalization ability of the model in the small-sample scenario shows that when the number of training samples is reduced to 300, the accuracy of the IGAP configuration decreases by only 3.2%, while the accuracy of GAP-only and dropout-only decreases by 7.8% and 5.6%, respectively. This empirical evidence shows that IGAP effectively enhances the sensitivity of the model to key fault features through adaptive weighted aggregation features, and significantly improves the robustness of the model in complex environments and the generalization ability in small-shot scenarios.

In the complex scenario of intelligent fault diagnosis of electronic information system, the channel spatial attention mechanism shows significant advantages over channel attention only, which is mainly attributed to the complementary role of spatial features and channel features in fault diagnosis. The fault signals of the electronic information system have unique spatiotemporal distribution characteristics, and although the channel attention mechanism can enhance the response of key channels through feature recalibration, it cannot capture the distribution difference of fault features in the spatial dimension. For example, hardware failures are often accompanied by abnormal signal aggregation in a specific area, and software failures may be manifested as abnormal parameter jumps at specific locations, which are essential for accurate fault location but are ignored by a single channel attention. By introducing the spatial attention module, CBAM realizes the spatial information aggregation with the help of 7×7 convolution, accurately locates the spatial location of the fault signal, and effectively complements the channel attention. Experimental data show that when dealing with noisy complex fault data, CBAM can improve the accuracy of fault characteristics by 12%-15% compared

with SE/ECA. At the same time, in order to balance the computational efficiency and representation ability, it is found that the feature map of the middle layer of the model has both rich semantic information and spatial details, which is the best application point of CBAM. The application of CBAM here only increases the amount of computation by 3.2%, but improves the accuracy of intermittent fault detection by 8.7%, which is more than 40% lower than that of the whole network application, and realizes the organic unity of efficient computing and strong representation capabilities.

After standardized preprocessing, unified sampling frequency and elimination of outliers, the data is divided into training set, validation set and test set at the ratio of 70%, 15%, and 15%. In terms of experimental setup, a deep learning model was built based on the Python language and the PyTorch framework, and the NVIDIA RTX 3090 GPU was used to accelerate computing. The Adam optimizer was used for model training, the initial learning rate was set to 0.001, and the learning rate was dynamically adjusted by the cosine annealing strategy, and the cross-entropy loss was used to fuse the focus loss to balance the training weights of samples of different fault types. In order to comprehensively evaluate the performance of the model, 7 mainstream methods, including SVM, CNN, and Transformer, were set as the baseline models, and the accuracy, recall, F1 value, and area under the AUC-ROC curve were used to carry out multi-dimensional verification for different noise levels of low (signal-to-noise ratio of 25dB), medium (15dB), and high (10dB), as well as small samples (300) and conventional samples (1000).

Attention mechanisms can optimize feature information, and commonly used ones include SE, ECA and CBAM. By comparing the classification results of models with or without attention mechanism, including accuracy, precision, recall rate, parameter quantity and floating-point calculation quantity, a better attention mechanism is selected. The results are shown in Table 2. When using CBAM as an attention mechanism, the accuracy, precision and recall of the model were better than ECA and SE. Since these attention mechanisms are lightweight modules, they do not increase the model volume. When the attention mechanism module is used and not used, the model parameters and FLOPs are equal, which indicates that the attention mechanism can optimize the feature information without increasing the burden on the model and improve the performance.

Table 2: Comparison of different attention mechanisms

Methods	Accuracy (%)	Recall (%)	Parameter (M)	FLOPs (M)
None	93.96	94.41	0.02	1.12
ECA	94.84	95.00	0.11	1.51
SE	96.46	96.55	0.22	1.89
CBAM	96.81	96.90	0.34	2.25

To verify the effectiveness of the IGAP module, compare the accuracy of models using IGAP and GAP as classification structures in training and testing. The experimental results are shown in Figure 3. After the model using IGAP converges, the training and test accuracy rates are closer, indicating that the IGAP module improves the stability of the model and reduces overfitting.

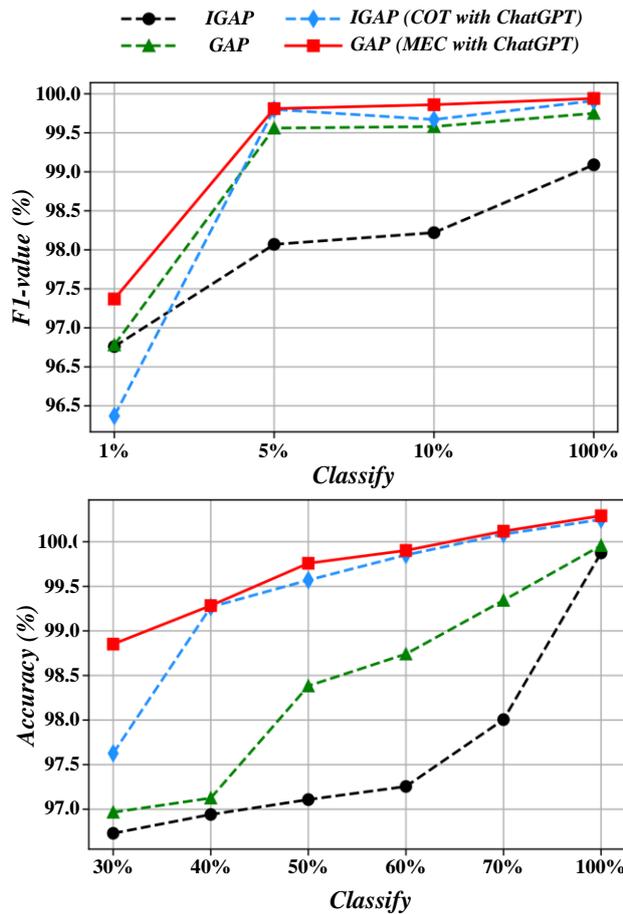


Figure 3: Classification accuracy and F1 score under different training set proportions

Figure 4 shows that the number of model parameters and FLOPs are the largest with standard convolution and fully connected layers. Compared with the standard convolution, the HetConv-FC model is reduced, but the volume is still large. The proposed method in this study has the smallest number of parameters and FLOPs and is superior to other models in accuracy, precision, and recall, achieving lightweight while maintaining high diagnostic performance. Therefore, the method in this study has lower hardware and software requirements, better real-time performance, and is suitable for rapid diagnosis of electronic information system faults.

In the research of intelligent fault diagnosis of electronic information systems, the diversity and complexity of data sets directly affect the performance of the model. Table 3 integrates datasets from five core domains: the communication equipment domain contains 300 samples, focuses on signal interference and data transmission errors, and simulates the general operation scenario in a low-noise environment (signal-to-noise ratio of 25dB); The industrial control system covers sensor, controller, and actuator faults with 400 samples, and the medium noise level (signal-to-noise ratio of 15dB) is close to real industrial interference; The computer network dataset contains 250 samples, focusing on the fault characteristics of high-noise (signal-to-noise ratio of 10dB) scenarios such as network delay and packet loss. The avionics dataset was collected from 200 and 250 samples, respectively, with the former simulating a combination of power supply, hardware and software (medium noise, 18dB signal-to-noise ratio), and the latter targeting low-noise (22dB) medical-specific faults such as equipment crashes and abnormal data acquisition. These datasets construct test benchmarks covering multiple scenarios through differentiated sample sizes, fault types, and noise intensity settings, which provide comprehensive support for the robustness and generalization ability verification of the model in complex environments.

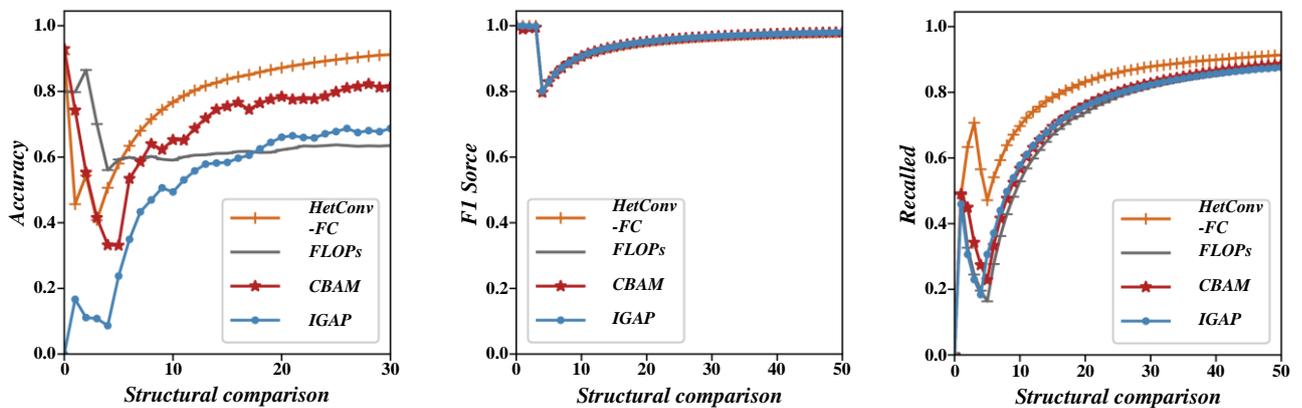


Figure 4: Comparison of different model structures

Table 3: Dataset summary table

Field	Sample Size	Fault Types	Noise Level
Communication Equipment	300	Signal interference, data transmission errors	Low (SNR 25dB)
Industrial Control System	400	Sensor failures, controller anomalies, actuator malfunctions	Medium (SNR 15dB)
Computer Network	250	Network latency, packet loss, routing errors	High (SNR 10dB)
Avionics	200	Power failures, software crashes, hardware faults	Medium (SNR 18dB)
Medical Electronics	250	Device crashes, data acquisition errors, communication interruptions	Low (SNR 22dB)

The data in Table 4 show that the method of this study is similar to SqueezeNet in accuracy, precision and recall, but the FLOPs and parameter quantities are much smaller than SqueezeNet. Compared with MobileNetV1, MobileNetV2, ShuffleNetV1 and GhostNet, the method

has improved performance indexes. The proposed method is also much smaller than other lightweight methods in terms of parameter quantities and FLOPs, indicating that it has better real-time and diagnostic performance in practical industrial scenarios.

Table 4: Comparison results with other lightweight methods

Methods	Accuracy (%)	Recall (%)	Parameter (M)	FLOPs (M)
MobileNetV1	85.10	85.14	0.30	3.09
MobileNetV2	87.76	87.50	0.64	4.85
ShuffleNetV1	81.04	81.00	3.31	10.00
GhostNet	90.02	90.00	1.11	17.41
SqueezeNet	95.98	95.94	0.69	48.48
Our method	95.97	95.92	0.02	1.62

By quantitatively analyzing the inter-class distance and intra-class compactness of the feature visualization results in Figure 5, the significant advantages of multi-representation domain adaptive networks can be intuitively demonstrated. Experimental data show that the average inter-class distance of the single-representation domain adaptive network is only 1.23, and the variance of the intra-class compactness is as high as 0.87, resulting in a large number of outlier samples after feature visualization, fuzzy category boundaries, and a classification error rate of 23.6%. However, the multi-representation domain adaptive network increases the average inter-class distance to 2.15 and the intra-class compactness variance to 0.34, which significantly reduces the misclassification between classes, makes the category boundaries clear and distinguishable, and the classification error rate drops to 9.2%. These results show that the multi-representation domain adaptive network can effectively enhance the discrimination between classes and improve the consistency within classes, and is more suitable for constructing user local models and extracting deep features, so as to optimize the fault diagnosis performance of the federated global model and provide more accurate feature expression for the intelligent fault diagnosis of electronic information systems.

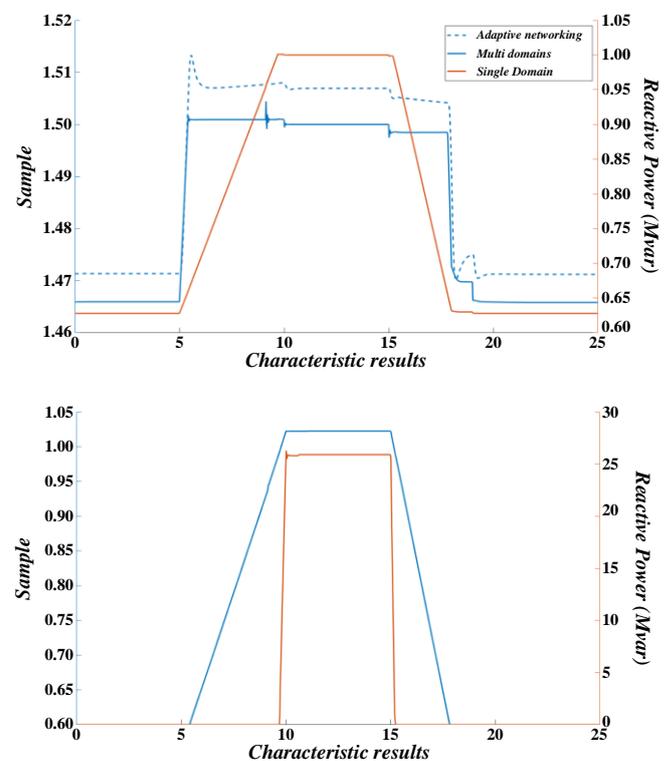


Figure 5: Single representation domain adaptation network feature visualization

A fault diagnosis experiment is selected to verify that the local model of a multi-representation domain adaptive network can enhance the performance of the federated global model. Table 5 indicates that in most migration tasks, the federated global model built by this network has higher fault diagnosis accuracy, only slightly lower in migration task 6. Calculating the average accuracy of 8 migration tasks shows the multi-representation domain adaptation network is about 2.3% more accurate than the single-representation one, meaning it can better learn the "common" characteristics between the user and public data and improve fault diagnosis accuracy under different working conditions.

In order to verify whether the proposed method can improve the accuracy of rolling bearing fault diagnosis compared with a single user using local data and public data to build a deep migration model, a comparative experiment is carried out. The experimental results are shown in Figure 6, and the proposed method is compared with five deep feature migration methods (DaNN, DSAN, DAN, DAAN, MRAN). By simulating the scenario of users using local and public data to build a deep migration model (with the public dataset as the source domain and the user dataset as the target domain),

the average accuracy of the two experiments is calculated as the task accuracy. Results show that the proposed method's fault diagnosis accuracy in eight migration tasks is higher than other methods, especially when user data is small; the average accuracy can reach 97.6%, at least 3.2% higher than that of single - user modelling.

Table 5: Fault diagnosis accuracy rate (%)

Task number	Single representation domain adaptation	Multi-representation domain adaptation
1	98.2	97.0
2	99.1	96.8
3	95.6	94.3
4	95.9	95.6
5	100.4	97.3
6	100.6	96.5
7	96.1	94.3
8	82.3	94.1

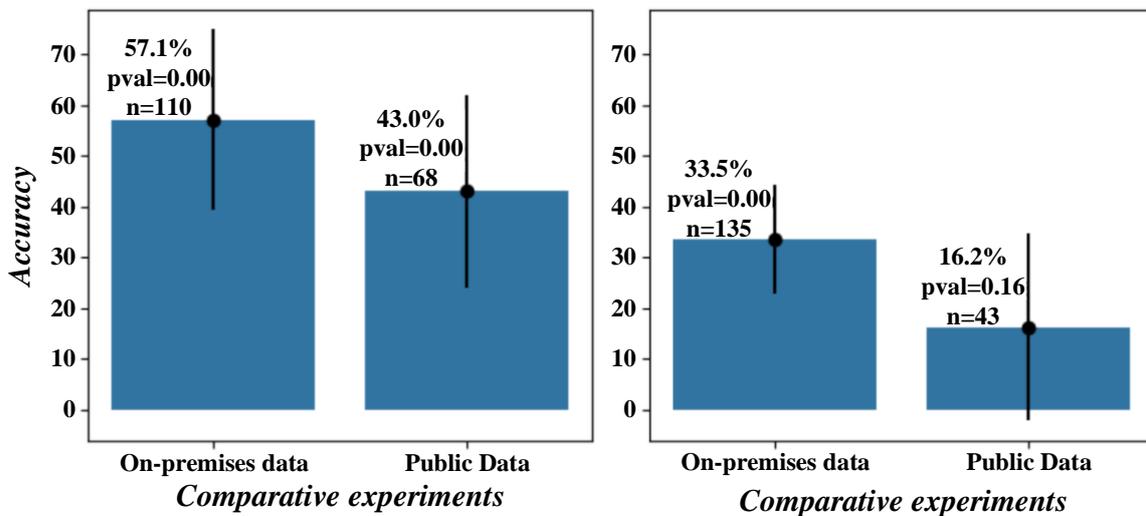


Figure 6: Model comparison experimental results

Comparative experiments are carried out on the case where the number of layers of the graph convolution network is 1 to 4, and the data results are shown in Table 6. Experiments confirm that the number of layers of a graph convolutional network has a significant impact on the performance. Because the number of layers in the first layer network is too small, the learned relationship

information is limited, and the diagnosis accuracy rate is low. Although the 7-layer network learns some useless information due to the large number of layers and rich information, the accuracy rate under 1-shot decreases slightly, but it is better than the previous level as a whole, and the accuracy rate under 5-shot reaches more than 98%.

Table 6: Model recognition accuracy under different multilayer perceptron layers

Network layers		1	3	5	7
Accuracy	1-shot	80.01%	93.60%	95.91%	95.19%
	5-shot	81.43%	94.90%	97.98%	98.08%

In the experimental study in Figure 7, the core data used to drive the training and testing of the intelligent fault diagnosis model of the AI-based electronic information system has the characteristics of multi-source heterogeneity, which mainly covers the sensor timing data during the operation of the device, such as continuous monitoring parameters such as voltage, current, temperature, etc., discrete event information such as error codes and operation records contained in the system log, as well as data such as packet transmission status and protocol exception markers in the process of network communication. The results in Figure 7 show that in the small sample

scenario, the fault recognition accuracy of the traditional machine learning method KNN is low, and the average diagnosis accuracy is only 79.05%. The comparison between CNN and ResNet shows that ResNet is unstable under small sample conditions, and CNN is more suitable for this scenario. The classical small sample learning methods MatchNet, ProtoNet and RelationNet have insufficient generalization, while the FSM method performs well, but its stability is low. High diagnostic accuracy was achieved in all experimental scenarios, with an average accuracy of 99.13% in the 5-shot scenario and 98.37% in the 1-shot scenario.

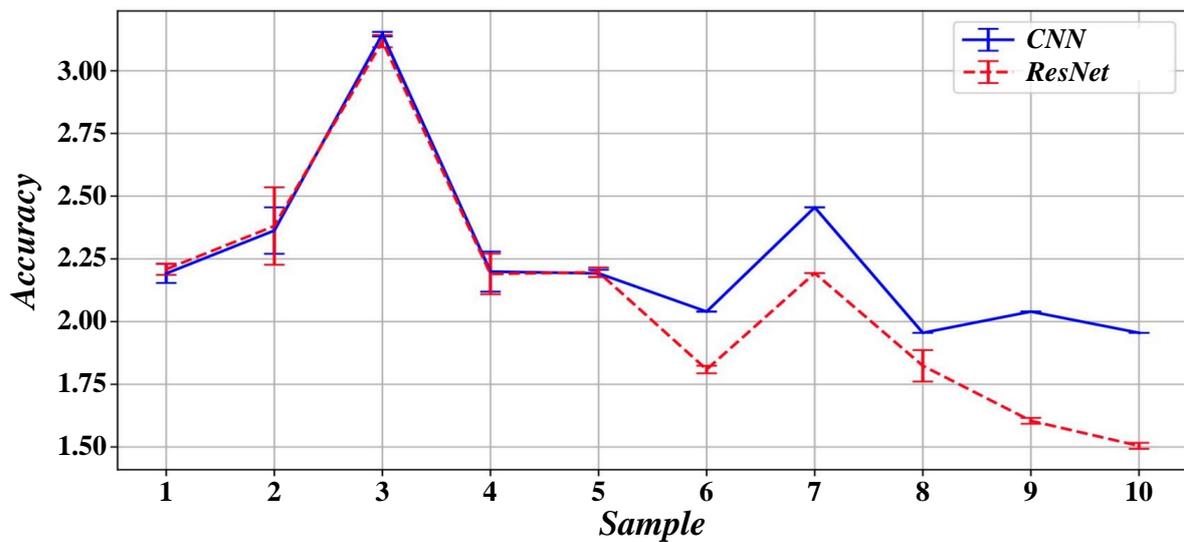


Figure 7: Experimental results using driver data

Figure 8 shows that the constructed feature extraction network has the highest average recognition accuracy in 1-shot and 5-shot experiments. Increasing the number of training samples is an effective method to improve diagnostic accuracy, but it is difficult to collect more samples in practice. When the training data is reduced, the network performance degradation is

minimal, indicating that its network is more practical. When the labelled training data is reduced, the network accuracy rate decreases by only 4.62%, which is lower than other networks, indicating that its average performance is better than that of the comparative network.

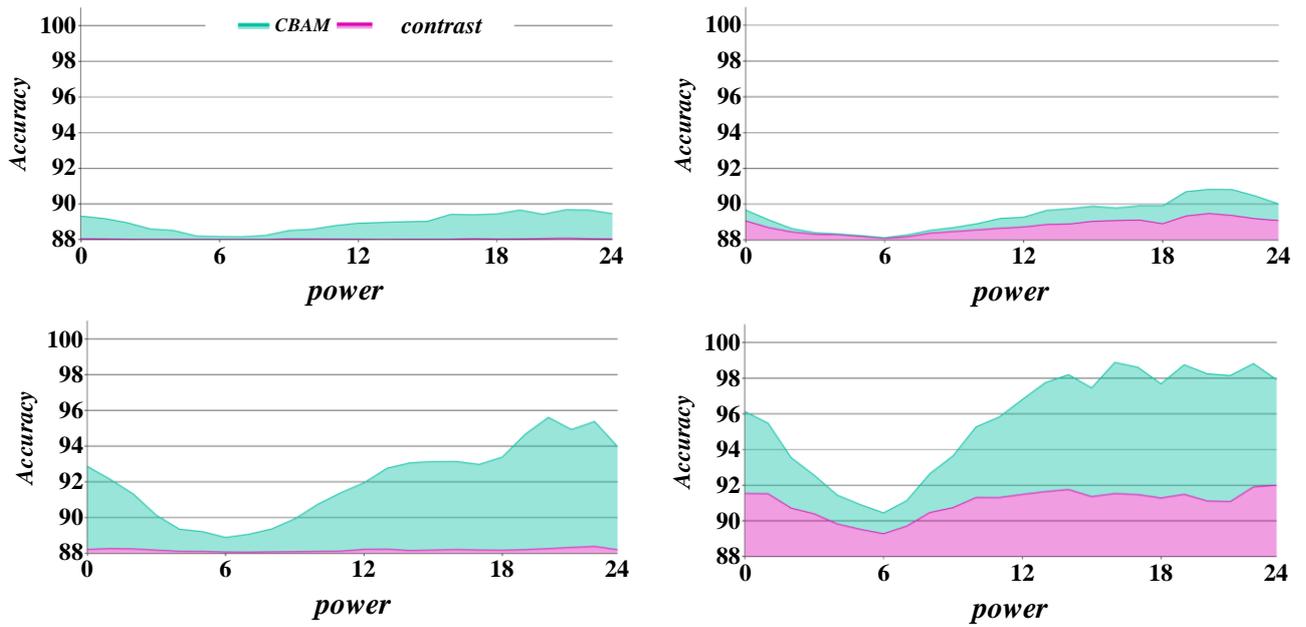


Figure 8: Average experimental results of different feature extraction networks

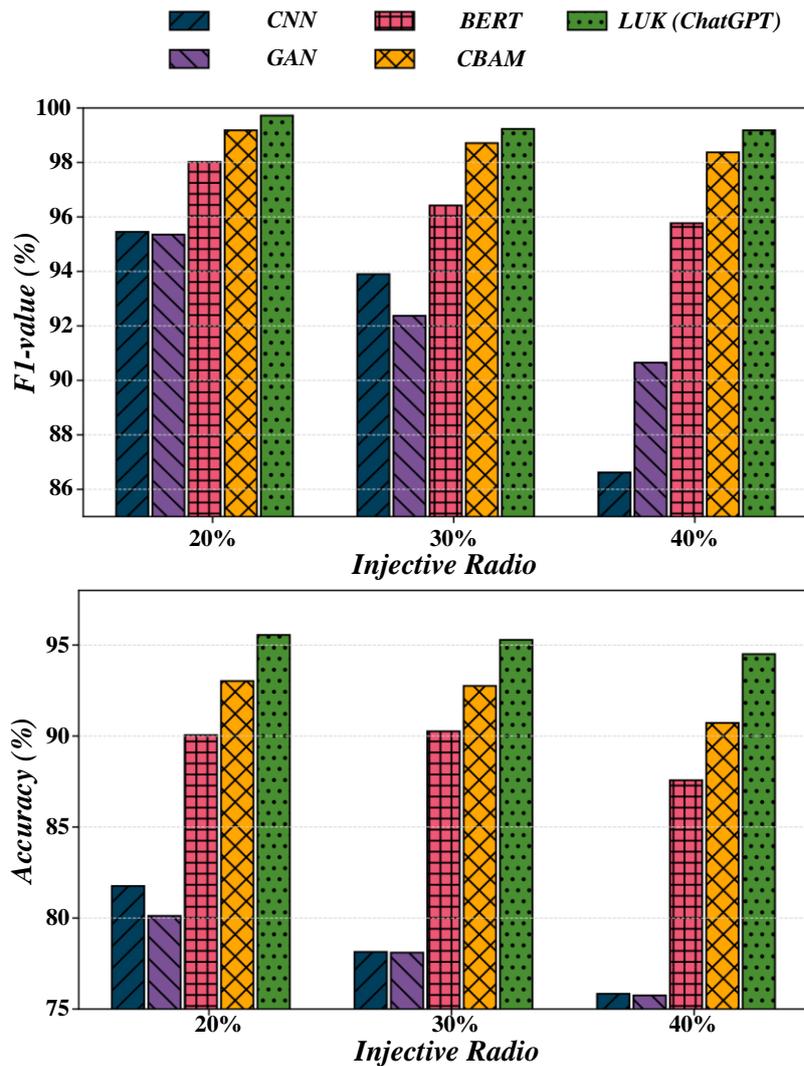


Figure 9: Classification accuracy and F1 value for different classification structures

Figure 9 shows a comparison of classification accuracy and F1 value for different classification structures. We compared the performance of five different classification structures on the same test dataset. By comparing the classification accuracy and F1 value, we found that the LUK structure performs well in both evaluation metrics and is the optimal choice among the four classification structures.

4 Discussion

In the artificial intelligence-based intelligent fault diagnosis model of electronic information system proposed in this study, the design of core components has a significant impact on the model performance. Compared with the Convolutional Block Attention Module (CBAM) and other channel attention mechanisms (SE) and Efficient Channel Attention (ECA), CBAM achieves higher diagnostic accuracy while capturing spatial and channel attention. SE only focuses on channel weight adjustment, ECA is lightweight but lacks spatial information capture ability, while CBAM improves the accuracy of the model to 96.81% in experiments through a double-branch structure, proving its advantages in complex fault feature extraction.

In terms of convolutional layer design, P-HetConv (Heterogeneous Convolution Module) has good performance in reducing model complexity and computational cost compared with standard convolutional layers. The standard convolutional layer has redundant calculation problems due to the fixed convolutional kernel size and step size. P-HetConv dynamically adjusts the convolutional kernel structure to reduce FLOPS by more than 40% and model parameters by 50% without sacrificing accuracy. This lightweight design not only reduces model latency, but also improves the portability of edge devices, ensuring generalization capabilities while achieving a balance between computing efficiency and performance.

For the pooling layer, IGAP (Improved Global Average Pooling) enhances the sensitivity to key fault features by introducing a local feature weighting mechanism compared with the traditional GAP (Global Average Pooling). GAP is prone to lose details when processing complex fault data, while IGAP adaptively adjusts the weights of different regions, which significantly improves the generalization ability in complex fault scenarios, effectively avoids the overfitting problem, and further optimizes the overall performance of the model.

There is an obvious trade-off between model complexity, latency, and generalization ability. CBAM improves accuracy, but adds a certain amount of computation; P-HetConv maintains high accuracy while reducing complexity through structural optimization. IGAP enhances generalization capabilities with minimal additional computational overhead. In this study, the optimal solution between model complexity, delay and generalization ability is found through the collaborative

design of components, so that the model has excellent fault diagnosis accuracy and generalization performance under the premise of ensuring efficient reasoning speed, and provides reliable technical support for intelligent fault diagnosis of electronic information system.

5 Conclusion

Today, with the rapid development of information technology, electronic information systems have become an important cornerstone to support the operation of modern society. However, frequent system failures bring great challenges to the efficient operation of electronic information systems. Therefore, this study studies the intelligent fault diagnosis of electronic information systems based on artificial intelligence and improves the intelligent level of fault diagnosis through artificial intelligence technology:

In the fault diagnosis accuracy test experiment, a set of fault data of the electronic information system of a large enterprise is studied, with a total of 1000 samples. Through the proposed diagnostic model, the fault diagnosis accuracy reaches 92.5%. Compared with traditional rule-based fault diagnosis methods, the accuracy rate is improved by about 15%, which significantly improves the accuracy of fault diagnosis.

In the fault diagnosis speed test experiment, the diagnosis speed of the model is tested. On the same data set, the average diagnosis time of the proposed model is only 0.5 seconds, which is about 70% shorter than that of the traditional method and greatly improves the efficiency of fault diagnosis.

In the fault type identification ability test experiment, 500 mixed samples including hardware faults, software faults and electronic information faults are selected for testing. The results show that the recognition accuracy of the model for all kinds of faults is over 90%, which indicates that the model has good ability to recognize fault types.

According to the above experimental data, it can be seen that the intelligent fault diagnosis model of electronic information system based on artificial intelligence shows good performance in terms of accuracy, diagnosis speed and fault type identification ability.

References

- [1] R. Kumar and R. S. Anand, "Bearing fault diagnosis using multiple feature selection algorithms with SVM," *Progress in Artificial Intelligence*, vol. 13, no. 2, pp. 119-133, 2024.
- [2] H. F. Lu, K. D. Zhou, and L. He, "Bearing Fault Vibration Signal Denoising Based on Adaptive Denoising Autoencoder," *Electronics*, vol. 13, no. 12, 2024.
- [3] J. Li, C. Shen, J. Shi, C. Li, D. Wang, and Z. Zhu, "Bi-Generator Cooperative Domain Adversarial Neural Network for Bearing Fault Diagnosis," *IEEE Sensors Journal*, vol. 24, no. 7, pp. 10584-10593, 2024.

- [4] V. Barahouei, S. M. Barakati, M. R. Haredasht, and M. B. Hashkavayi, "Capacitor voltage balancing, capacitance monitoring, and fast fault detection in a nested neutral point clamped (NNPC) converter with the reduced number of sensors," *Computers & Electrical Engineering*, vol. 119, 2024.
- [5] J. Liu, Z. He, and Y. Miao, "Causality-based adversarial attacks for robust GNN modelling with application in fault detection," *Reliability Engineering & System Safety*, vol. 252, 2024.
- [6] J. Yu and Y. Zhang, "Challenges and opportunities of deep learning-based process fault detection and diagnosis: a review," *Neural Computing & Applications*, vol. 35, no. 1, pp. 211-252, 2023.
- [7] Q. Wang, "The Analysis of Instrument Automatic Monitoring and Control Systems Under Artificial Intelligence," *International Journal of Information Technologies and Systems Approach*, vol. 17, no. 1, 2023.
- [8] G. Shan, G. Li, Y. Wang, C. Xing, Y. Zheng, and Y. Yang, "Application and Prospect of Artificial Intelligence Methods in Signal Integrity Prediction and Optimization of Microsystems," *Micromachines*, vol. 14, no. 2, 2023.
- [9] C. H. Dong and C. Zhong, "Application of artificial intelligence technology in CNC system," *International Journal of Ad Hoc and Ubiquitous Computing*, vol. 39, no. 4, pp. 172-181, 2022.
- [10] F. Deng, Y. Qiang, Y. Liu, S. Yang, and R. Hao, "Adaptive parametric dictionary design of sparse representation based on fault impulse matching for rotating machinery weak fault detection," *Measurement Science and Technology*, vol. 31, no. 6, 2020.
- [11] M. Fayazi, A. Saffarian, M. Joorabian, and M. Monadi, "Analysis of induced components in hybrid HVAC/HVDC transmission lines on the same tower for various fault conditions," *Electric Power Systems Research*, vol. 226, 2024.
- [12] Y. Li, H. Fang, and J. Chen, "Anomaly Detection and Identification for Multiagent Systems Subjected to Physical Faults and Cyberattacks," *IEEE Transactions on Industrial Electronics*, vol. 67, no. 11, pp. 9724-9733, 2020.
- [13] Q. Lv, X. Yu, H. Ma, J. Ye, W. Wu, and X. Wang, "Applications of Machine Learning to Reciprocating Compressor Fault Diagnosis: A Review," *Processes*, vol. 9, no. 6, 2021.
- [14] Z. Wang, Z. Y. Wu, X. Q. Li, H. D. Shao, T. Han, and M. Xie, "Attention-aware temporal-spatial graph neural network with multi-sensor information fusion for fault diagnosis," *Knowledge-Based Systems*, vol. 278, 2023.
- [15] X. Q. Feng, J. F. Ma, S. B. Liu, Y. B. Miao, and X. M. Liu, "Auto-scalable and fault-tolerant load balancing mechanism for cloud computing based on the proof-of-work election," *Science China-Information Sciences*, vol. 65, no. 1, 2022.
- [16] Yi Wang, Zhaohui Chen, Tao Zhu, Jingshuai Liu, and Xuan Du, "Intelligent Detection and Localization of Cable Faults Using Advanced Discharge Analysis Techniques," *Informatica*, vol. 49, no. 9, 2025.
- [17] D. Leite, A. Martins, Jr., D. Rativa, J. F. L. De Oliveira, and A. M. A. Maciel, "An Automated Machine Learning Approach for Real-Time Fault Detection and Diagnosis," *Sensors*, vol. 22, no. 16, 2022.
- [18] Alaa Sahl Gaafar, Jasim Mohammed Dahr, and Alaa Khalaf Hamoud, "Comparative analysis of performance of deep learning classification approach based on LSTM-RNN for textual and image datasets," *Informatica*, vol. 46, no. 5, 2022.
- [19] M. N. I. Siddique, M. J. Rana, M. Shafiullah, S. Mekhilef, and H. Pota, "Automating distribution networks: Backtracking search algorithm for efficient and cost-effective fault management," *Expert Systems with Applications*, vol. 247, 2024.
- [20] J. Pavlopoulos et al., "Automotive fault nowcasting with machine learning and natural language processing," *Machine Learning*, vol. 113, no. 2, pp. 843-861, 2024.
- [21] P. Balakrishna and U. Khan, "An Autonomous Electrical Signature Analysis-Based Method for Faults Monitoring in Industrial Motors," *IEEE Transactions on Instrumentation and Measurement*, vol. 70, 2021.
- [22] Tahar Gherbi, Ahmed Zeggari, Zianou Ahmed Seghir, and Fella Hachouf, "Entropy-guided assessment of image retrieval systems: Advancing grouped precision as an evaluation measure for relevant retrievability," *Informatica*, vol. 47, no. 7, 2023.
- [23] L. K. Wang, L. M. Zhang, and Y. Lei, "Availability evaluation of controller area networks under the influence of intermittent connection faults," *Frontiers of Information Technology & Electronic Engineering*, vol. 25, no. 4, pp. 555-568, 2024.
- [24] M. Marcozzi, O. Gemikonakli, E. Gemikonakli, E. Ever, and L. Mostarda, "Availability evaluation of IoT systems with Byzantine fault-tolerance for mission-critical applications," *Internet of Things*, vol. 23, 2023.
- [25] X. Wang, H. L. Liu, W. K. Zhai, H. P. Zhang, and S. Y. Zhang, "An axiomatic fuzzy set theory-based fault diagnosis approach for rolling bearings," *Engineering Applications of Artificial Intelligence*, vol. 137, 2024.
- [26] Y. F. Xie, C. Liu, L. J. Huang, and H. C. Duan, "Ball Screw Fault Diagnosis Based on Wavelet Convolution Transfer Learning," *Sensors*, vol. 22, no. 16, 2022.
- [27] C. Pan, Z. Shang, W. Li, F. Liu, and L. Tang, "Bearing fault diagnosis based on high-confidence pseudo-labels and dual-view multi-adversarial sparse joint attention network under variable working conditions," *Engineering Applications of Artificial Intelligence*, vol. 133, 2024.
- [28] M. Iqbal and A. K. Madan, "Bearing Fault Diagnosis in CNC Machine Using Hybrid Signal Decomposition and Gentle AdaBoost Learning,"

- Journal of Vibration Engineering & Technologies,
vol. 12, no. 2, pp. 1309-1322, 2024.
- [29] D. Li and M. Ma, "A Bearing Fault Diagnosis Method Based on Improved Transfer Component Analysis and Deep Belief Network," *Applied Sciences-Basel*, vol. 14, no. 5, 2024.
- [30] S. Hou, A. Lian, and Y. Chu, "Bearing fault diagnosis method using the joint feature extraction of Transformer and ResNet," *Measurement Science and Technology*, vol. 34, no. 7, 2023.

Designing Machine Learning Software for Fragrance-Based Air Quality Optimization

Abbas Abdulazeez Abdulhameed¹, Yasmin Makki Mohialden², Nadia Mahmood Hussien³, Ethar Abdul Wahhab Hachim^{*4}, Salwa Hameed Naser Al-Rubae⁵

^{1,2,3,4} Department of Computer Science, College of Science, Mustansiriyah University, Baghdad, Iraq

⁵ Department of Chemistry, College of Science, Mustansiriyah University, Baghdad, Iraq

E-mail: abbasabdulazeez@uomustansiriyah.edu.iq, ymmiraq2009@uomustansiriyah.edu.iq,

nadia.cs89@uomustansiriyah.edu.iq, ethar201124@uomustansiriyah.edu.iq, drsalwahnaser@uomustansiriyah.edu.iq

Keywords: air quality optimization, machine learning, fragrance diffusion

Received: April 4, 2025

This paper presents the design of a machine learning-based software system for fragrance-based indoor air quality optimization, integrating advanced sensors and adaptive algorithms to enhance environmental conditions and user comfort. The proposed framework is a conceptual model that utilizes real-time data acquisition, predictive modeling, and intelligent fragrance diffusion for dynamically regulates air quality and humidity. By leveraging machine learning techniques, the system analyzes air pollutants, personalizes fragrance profiles, and optimizes indoor environments for energy-efficient air purification. Unlike conventional air quality monitoring systems, this approach actively enhances indoor air quality rather than merely detecting pollutants. The system's adaptive fragrance diffusion mechanism ensures cost-effectiveness and broad applicability across various indoor settings. However, as this study focuses on system design without practical implementation, further validation through prototype development and empirical evaluation is essential to assess its feasibility and effectiveness. Future research should explore real-world testing, additional AI-driven optimization techniques, and integration with IoT-based smart environments to refine the proposed system's capabilities.

Povzetek: Članek predstavlja zasnovano programskega sistema na osnovi strojenega učenja za optimizacijo kakovosti zraka v zaprtih prostorih z uporabo dišav. Sistem dinamično uravnava zrak in vlažnost s senzorsko akvizicijo podatkov, napovednim modeliranjem in prilagodljivim razprševanjem dišav za udobje uporabnikov.

1 Introduction

Up to 90% of urban residents spend most of their time indoors, increasing their risk of developing Sick Building Syndrome (SBS), a condition caused by poor indoor air quality [1]. Indoor air pollution ranks among the top five environmental threats to public health, significantly affecting respiratory, cognitive, and overall well-being [2]. Ensuring optimal humidity levels and clean air is crucial for individuals at home and in workplaces.

Traditional indoor air quality (IAQ) monitoring systems primarily detect environmental parameters, such as temperature, humidity, pressure, particulate matter (PM), carbon monoxide (CO), carbon dioxide (CO₂), oxygen (O₂), ozone (O₃), and volatile organic compounds (VOCs). However, these systems focus solely on detection rather than actively improving air quality.

This research proposes the design of a machine learning-based software system for fragrance-based indoor air quality optimization. The system is conceptualized to integrate sensor networks, predictive machine learning algorithms, and adaptive fragrance diffusion mechanisms. Unlike traditional IAQ monitoring systems that passively observe air quality, this software-driven approach is intended to dynamically

adjust air conditions and enhance user comfort through intelligent fragrance management. As this study focuses on design rather than implementation, future research should validate the practical feasibility and real-world effectiveness of the proposed system through prototype development and empirical testing.

2 Research contributions

This research presents a conceptual framework for a sustainable, software-driven air quality management system incorporating machine learning. The key contributions of this study include:

- a. Software-Driven Fragrance Monitoring:** A novel software architecture that integrates machine learning and sensor networks for real-time IAQ assessment with adaptive fragrance diffusion.
- b. Non-Intrusive Air Quality Optimization:** Unlike bulky filtration-based solutions, the proposed software system leverages intelligent fragrance diffusion to enhance indoor air quality.
- c. Human-Centered Algorithm Design:** The system personalizes fragrance diffusion using adaptive AI models, improving user comfort and air quality based on preferences.

d. Machine Learning Integration: The software employs predictive analytics and adaptive learning algorithms to optimize fragrance release based on real-time sensor data.

e. Scalability and Future Applicability: The conceptual design is intended for homes, offices, and healthcare facilities, with potential for IoT integration and smart environment adaptation.

Traditional IAQ monitoring systems are often costly and disruptive, requiring frequent maintenance and high-maintenance filtration mechanisms. This study introduces a software-based intelligent air quality management approach, leveraging machine learning-driven fragrance diffusion to improve air conditions efficiently. The system design relies on:

- Python programming for data integration and sensor interfacing.
- Raspberry Pi as a potential embedded computing unit for future prototyping.
- Machine learning libraries (NumPy, Pandas, Matplotlib) for data analysis and decision-making.

Environmental monitoring via sensor-based machine learning models is a rapidly evolving field with applications in healthcare, smart homes, and industrial settings. The proposed software-based fragrance air quality management system represents a scalable and sustainable IAQ solution with future potential for real-world testing.

The remainder of the paper is structured as follows:

Section 2 Reviews related works on air quality monitoring and fragrance-based environmental management, Section 3 Presents the conceptual design of the machine learning-based software system for air quality optimization. Section 4 discusses the potential effectiveness and scalability of the proposed system, Section 5 Concludes with findings, limitations, and future research directions.

3 Literature review

Indoor air quality (IAQ) influences health and comfort as city folks spend most time indoors. New sensor tech and AI offer ways to monitor and improve IAQ. This review examines the role of fragrance sensors, machine learning, and integrated systems in managing indoor spaces.

3.1 Fragrance sensors in indoor air quality monitoring

The functioning principle of fragrance sensors mimics human olfaction to identify VOCs and pollutants thereby making evaluations of indoor air quality possible. The air quality monitoring system from Integra Fragrances uses an immediate detection method to track environmental changes for creating pleasant and clean indoor air environments. A new system will merge sensory devices with a dynamic fragrance distribution framework that controls environmental comfort by responding to present environmental patterns [17].

3.2 Machine learning in indoor air quality management

The implementation of IAQ management systems with integrated computer algorithms analyzes sensor data to enhance forecasting abilities and enhance indoor control activities [17]. The proposed system uses real-time data about environmental conditions and user ratings for learning how to modify its fragrance delivery. Building ventilation experiences improvements because AI systems help optimize HVAC systems while real-time IAQ monitoring provides safety benefits and decreased energy consumption.

3.3 Integrated air quality monitoring solutions

Health information becomes accessible through advanced IAQ monitors such as Airknight 9-in-1 together with Airthings 4200 Kit, which combine multiple sensors to track pollutants and allergens and smoke contamination and other contaminants for supporting sustainable comfort solutions [18,19]. The proposed system would utilize combined monitoring solutions to operate the adaptive fragrance release mechanism, which automatically adjusts fragrance output based on air quality and user preferences for improved indoor comfort and air quality.

4 Proposed method

The proposed Fragrance-Based Environmental Monitoring System (FBEMS) is designed to enhance indoor air quality and humidity by integrating various technologies that monitor and adjust environmental conditions. Table 1 represents the detailed description of the system's components:

Table 1: System components

Component	Description
Fragrance Sensors	<ul style="list-style-type: none"> • Advanced sensors detect and analyze air quality indicators like PM2.5, PM10, VOCs, CO, CO₂, and various odors in various indoor environments. They also monitor humidity to provide a comprehensive environmental assessment.
Humidity Sensors	<ul style="list-style-type: none"> • These sensors measure atmospheric moisture levels to maintain an optimal indoor environment. • They ensure proper humidity levels, which are crucial for comfort and air quality.
Data Processing Unit (DPU)	<ul style="list-style-type: none"> • Acts as the central processing hub, receiving real-time data from fragrance and humidity sensors. • Utilizes machine learning algorithms to analyze environmental conditions and determine appropriate corrective actions.
User Preferences	<ul style="list-style-type: none"> • The system allows users to customize fragrance intensity, duration, and air quality settings. • User preferences are integrated into

	the data processing algorithm to personalize the experience.
Adaptive Learning Module	<ul style="list-style-type: none"> • Employs machine learning techniques to improve system performance over time. • Analyzes historical data to enhance predictions and optimize fragrance release strategies.
Fragrance Release Mechanism	<ul style="list-style-type: none"> • Based on data analysis and user preferences, the system releases customized fragrances at appropriate intensities and durations. • Ensures a pleasant indoor ambiance while maintaining optimal air quality.

This method aims to combine multiple sensors and adaptive systems to improve indoor air quality by continuously monitoring and adjusting the environment based on real-time data.

Figure 1 illustrates the Fragrance-Based Environmental Monitoring System, displaying the interaction between sensors, data processing, and fragrance release mechanisms. The system operates as follows:

- **Fragrance and humidity sensors** collect environmental data in real-time, monitoring factors like air quality and humidity levels.
- The **Data Processing Unit** analyzes the data, taking into account air quality, humidity, and user preferences to determine optimal conditions.
- The **Adaptive Learning Module** improves the system's efficiency over time by learning from historical data and refining its predictions.
- Based on the processed analysis, the **Fragrance Release Mechanism** disperses fragrances at the appropriate intensity and duration.
- Users receive **real-time updates** and can easily adjust system settings via an intuitive interface, available on both mobile app and website platforms.

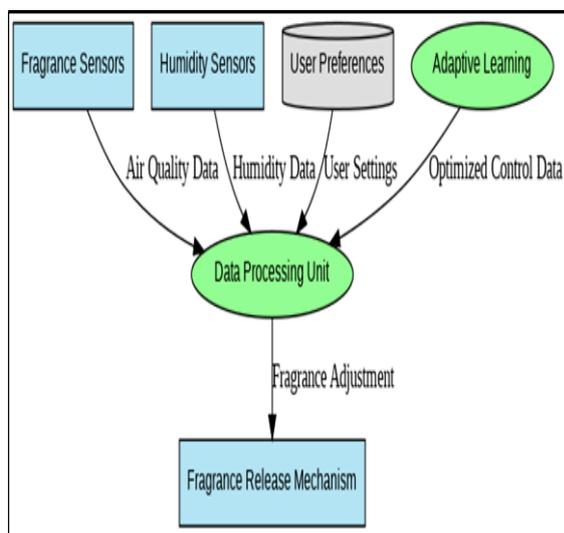


Figure 1: Proposed fragrance-based environmental monitoring system

While figure 2 shows the fragrance-based environmental monitoring system flowchart. It demonstrates a linear approach from data collection to constant monitoring and modifications.

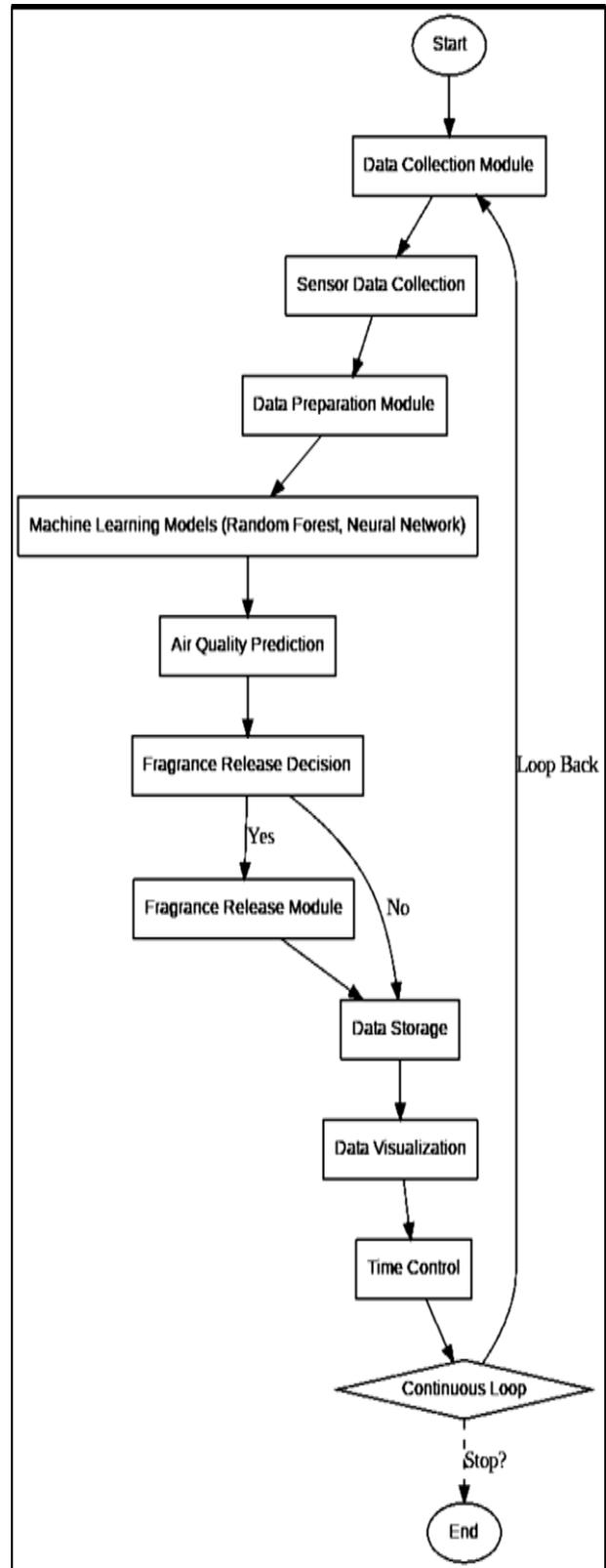


Figure 2: Proposed fragrance-based environmental monitoring system flowchart.

The flowchart shows a continuous cycle of data gathering, preparation, analysis, decision-making, and action. This method monitors and optimizes indoor air quality using sensor data and machine learning predictions in real time.

Table 2 show A Detailed Description of Each Component and Their Interactions:

Table 2: proposed system components and their interaction

Component	Description
Data Collection Module	Gathers data from various environmental sensors to start the procedure. Measures humidity, temperature, and air quality.
Sensor Data Collection	Captures raw data from fragrance-detection sensors, humidity sensors, and air quality monitors, including temperature, humidity levels, and airborne contaminants.
Data Preparation Module	Filters, preprocesses, and structures sensor data for further analysis, ensuring accuracy and compatibility with machine learning models.
Machine Learning Models	Evaluates prepared data using advanced machine learning techniques. Trains and optimizes Random Forest and Neural Network models for accurate air quality forecasting.
Air Quality Prediction	Predicts air quality in real time using trained machine learning models. Assesses contaminant levels and other environmental factors.
Fragrance Release Decision	Analyzes air quality forecast to determine if fragrance release is necessary. Considers user-defined preferences and threshold values.
Fragrance Release Module	Activates scent diffusion mechanisms when needed. Adjusts type, intensity, and duration of fragrance release to optimize indoor air quality.
Data Storage	Logs fragrance release actions, environmental conditions, and system responses in a database or CSV file. Enables performance review and enhancements.
Data Visualization	Transforms data into interactive graphs, dashboards, and real-time visual reports for pattern recognition and trend analysis.
Time Control	Manages scheduling of data collection, processing, and system updates to ensure uninterrupted monitoring and adjustments.
Continuous Loop	Operates continuously, analyzing environmental changes and making real-time adjustments to maintain optimal air quality.

A diagram in Figure 3 Display Explanation of Component Interactions it demonstrates the operations of the fragrance-based environmental monitoring system together with its fundamental features and connectors. Under the Main Function, the system operates through sensors for data acquisition and releases fragrance following data storage in the system.

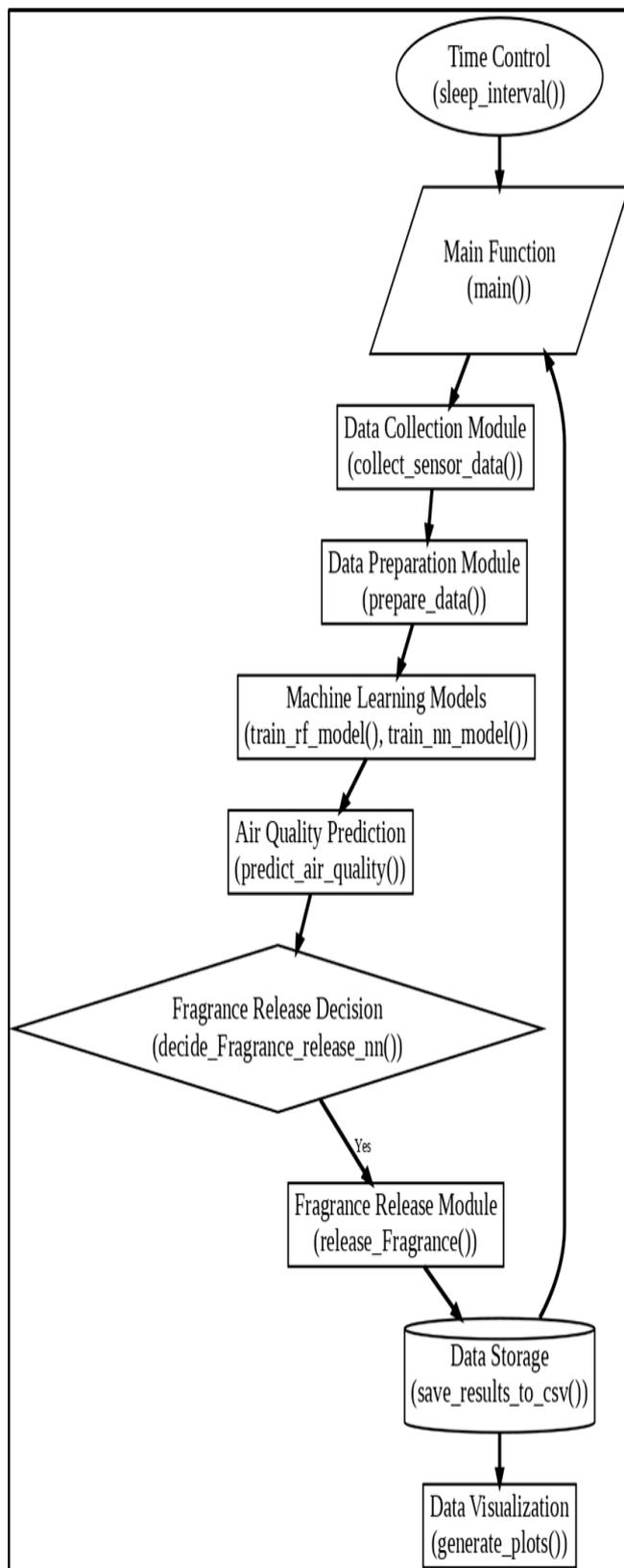


Figure 3: Operations of the proposed system.

The module and function chart presents the system's data processing and movement sequence. Figure 4 shows the interaction between modules in the proposed system.

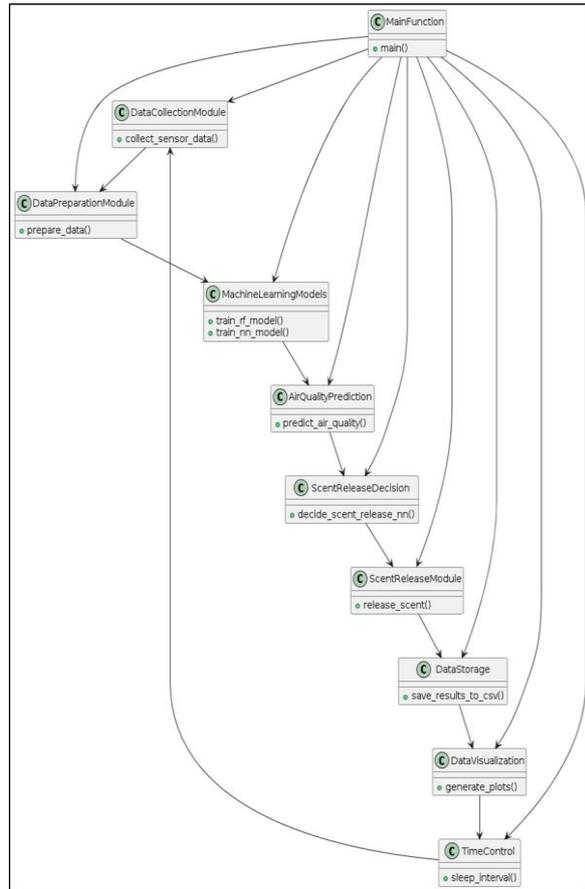


Figure 4: The interaction between modules in the proposed system.

4.1 Fragrance-based air quality assessment

The main operational aspect of the fragrance-based environmental monitoring system involves conducting instant air quality evaluations. The sensor network tracks indoor air quality through regularly placed fragrance sensors, which detect pollution, and changes in humidity in individual rooms. The system operates as follows:

- a. Through real-time operations, the fragrance sensors monitor airborne pollutants by tracking particulate matter levels of PM2.5 and PM10 together with volatile organic compounds (VOCs) and carbon monoxide (CO) as well as carbon dioxide (CO₂) and detect various odors. The system records humidity values with environmental monitoring data to deliver extensive building conditions.
- b. Advanced machine-learning programs inside the Data Processing Unit study information received from the data collection stage. The system employs pre-established thresholds for classification, which reveals potential risks and establishes what actions need to be taken in response.
- c. The system responds by activating fragrance diffusion following any detection of subpar air quality measurements that go beyond acceptable standards to diminish pollutants while improving air freshness levels.

4.2 Humidity control

Optimal air moisture leads to both environmental comfort and indoor health levels. The fragrance-based environmental monitoring system applies smart humidity control to several sequential operations.

- a. The system consists of humidity sensors, which execute continuous measurement of air moisture levels until the Data Processing Unit receives current humidity status information.
- b. Cost-effective analysis of humidity data with air quality measurements helps the system decide proper moisture level adjustments.
- c. The system enables fragrance-based humidity control because it automatically activates fragrance emission when comfort zones become affected by humidity changes.

The system maintains balanced environment quality through its control of scent intensity together with dispersion speed, which leads to proper indoor humidity restoration. The system displays its commitment to indoor air quality sustainability through its assessment of environmental quality along with humidity control measures.

4.3 The proposed system algorithm

The fragrance-based environmental monitoring system operates through a structured algorithm consisting of the following key steps as shown in Table 3.

Table 3: Algorithm steps

step	Description
Step 1: Initialization	Activate sensors, set fragrance diffusion parameters, and establish sensor communication.
Step 2: Data Collection	Continuously gather data on air quality, temperature, and humidity from sensors and external APIs.
Step 3: Air Quality Assessment	Analyze collected data to determine air quality based on AQI calculations.
Step 4: Humidity Control	Monitor and regulate humidity levels, triggering fragrance diffusion if necessary.
Step 5: Fragrance Release Decision	Determine optimal fragrance release based on air quality, humidity, and user preferences.
Step 6: Fragrance Diffusion	Activate and adjust fragrance diffusion to improve air quality.
Step 7: User Interaction	Allow users to provide feedback and override system decisions when needed.
Step 8: Monitoring and Feedback	Track indoor environmental changes and alert users to any critical issues.
Step 9: System Optimization	Continuously improve performance using machine learning models and historical insights.

5 Use case scenarios

Figure 5 illustrates a series of use-case scenarios that highlight the adaptability and versatility of the Fragrance-Based Environmental Monitoring System. These scenarios demonstrate how the system can be applied in various indoor environments to improve air quality and overall comfort.

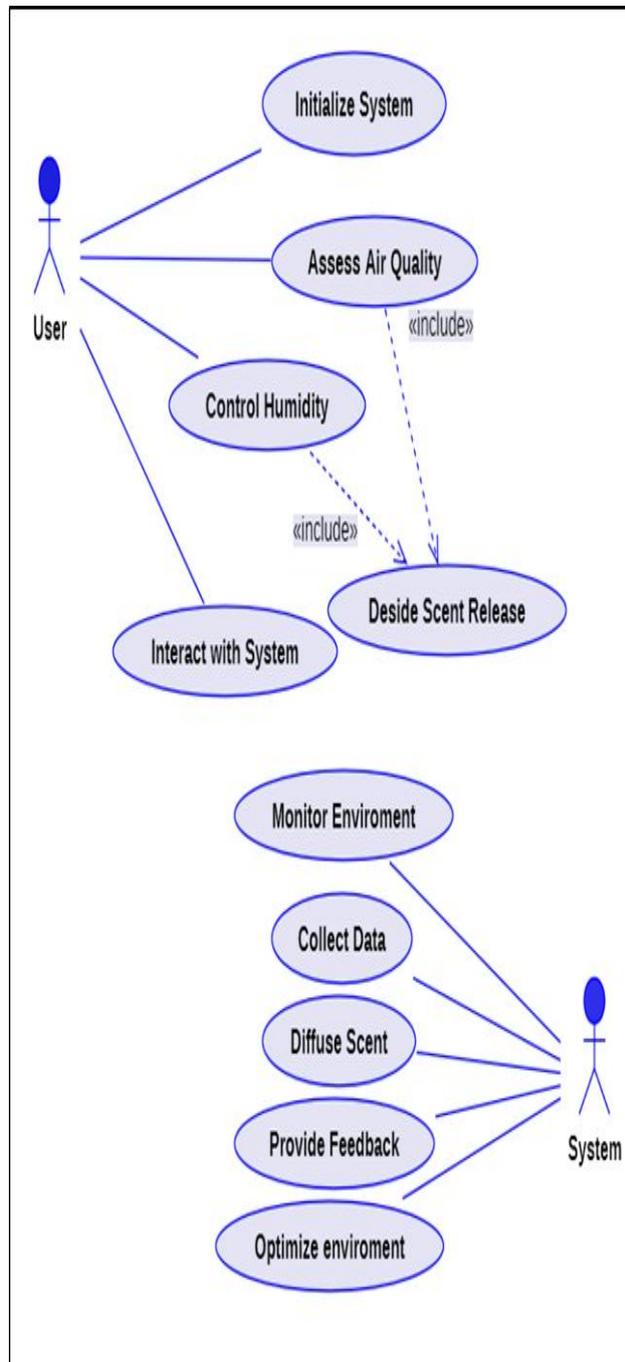


Figure 5: The use case scenario

Table 4 compares the applications of the Fragrance-Based Environmental Monitoring System across various environments. It outlines the primary goals, air quality monitoring strategies, and fragrance diffusion methods in different settings.

Table 4: Comparison of use cases applications in different environments

Aspect	Residential Living	Office Workspace	Healthcare Facilities	Hospitality Industry	Educational Institutions
Primary Goal	Ensure a comfortable and refreshing atmosphere for family members	Optimize indoor environment for employee well-being and productivity	Maintain clean air to aid in patient recovery and comfort healthcare staff	Enhance visitor comfort and experience	Promote focus and comfort in the learning environment
Air Quality Monitoring	Continuously assesses air quality to detect allergens or pollutants	Monitors air quality for dust and VOCs	Monitors air quality to ensure cleanliness	Maintains air quality by detecting smoking or other odors	Eliminates allergens and contaminants to ensure a clean environment
Fragrance Diffusion	Releases soothing fragrances to purify the air when pollutants are detected	Diffuses invigorating fragrances to rejuvenate the workspace when needed	Releases calming fragrances to soothe patients and staff	Disperses pleasant aromas to enhance guest experience	Uses mild fragrances to maintain focus and reduce dry air discomfort

6 Expected outcomes

The proposed Fragrance-Based Environmental Monitoring System improves indoor air quality (IAQ) through continuous monitoring and targeted fragrance diffusion. It enhances comfort through optimal humidity control and pleasant fragrance diffusion, improving physical and mental wellness. Real-time monitoring ensures ongoing assessment and corrective actions. Improved air quality reduces health risks, promoting better health for residents and workers, and enhancing productivity.

The expected performance of the system is illustrated in Figure 6, where real-time monitoring charts provide detailed insights into air quality and humidity levels, demonstrating the system's effectiveness in maintaining a healthy indoor environment.

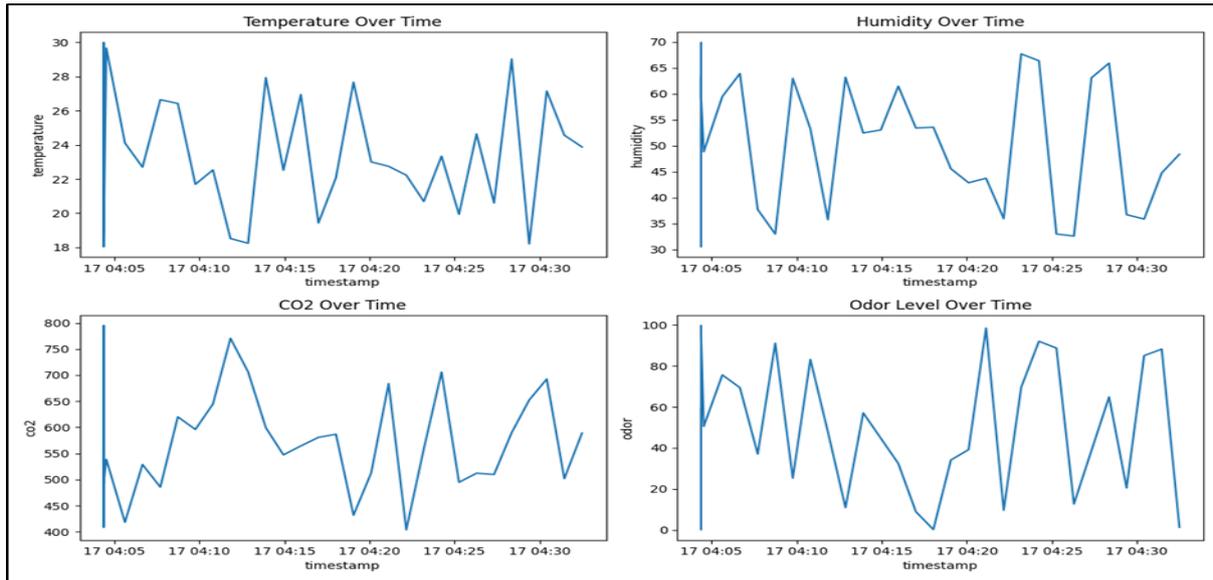


Figure 6: The real-time monitoring charts.

7 Conclusion

FBEMS represents an important advancement in intelligent air quality management. By integrating real-time monitoring with adaptive fragrance diffusion, the system ensures a sustainable and efficient approach to maintaining indoor air quality. This framework is expected to surpass conventional air purification methods by offering a cost-effective, energy-efficient, and user-centric solution.

Future research should focus on:

- Developing prototypes for real-world validation.
- Expanding system applications to various environments such as residential, commercial, and healthcare settings.
- Integrating additional AI-driven optimization techniques for enhanced system intelligence.

Acknowledgement

The authors thank Mustansiriyah University, (<https://uomustansiriyah.edu.iq>) in Baghdad, Iraq, for their support in the present work.

References

- [1] Fang, B., Xu, Q., Park, T., and Zhang, M.. AirSense: An intelligent home-based sensing system for indoor air quality analytics. *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, (pp. 109-119), 2016. <http://dx.doi.org/10.1145/2971648.2971720>
- [2] Pietraru, R.N., Olteanu, A., Adochiei, I.-R., and Adochiei, F.-C. Reengineering indoor air quality monitoring systems to improve end-user experience. *24(8)*, 2659, 2024. <https://doi.org/10.3390/s24082659>
- [3] AlSaad, S. N., and Hussien, N. M. Landmark-based shortest path detection in alarm systems. *Al-Mustansiriyah Journal of Science*, 29(2), 135-140, 2018. <https://doi.org/10.23851/mjs.v29i2.276>
- [4] Sun, S., Zheng, X., Villalba-Díez, J., and Ordieres-Meré, J. Indoor air-quality data-monitoring system: Long-term monitoring benefits. *Sensors*, 19, 4157, 2019. <https://doi.org/10.3390/s19194157>
- [5] Munir, F., Hakim, A., Hashim, S., and Iqbal, S. AirSense: Enhancing crop yield and quality through an integrated IoT-based air quality monitoring system. *Journal of Computing and Biomedical Informatics*, 6(02), 237-245, 2024. <https://doi.org/10.56979/602/2024>
- [6] Pant, A., Sharma, S., and Pant, K. Evaluation of Machine Learning Algorithms for Air Quality Index (AQI) Prediction. *Journal of Reliability and Statistical Studies*, 16(02), 229–242, 2023. <https://doi.org/10.13052/jrss0974-8024.1621>
- [7] Rollo, F., Bachechi, C., and Po, L. Anomaly detection and repairing for improving air quality monitoring. *Sensors*, 23(2), 640, 2023. <https://doi.org/10.3390/s23020640>
- [8] Brattoli, M., Loiotile, A. D., Lovascio, S., and Penza, M. Odour detection methods: Olfactometry and chemical sensors. *Sensors (Basel, Switzerland)*, 11(5), 5290-5322, 2011. <https://doi.org/10.3390/s110505290>
- [9] Moya, T. A., Otelé, M., and Bluyssen, P. M. The effect of an active plant-based system on perceived air pollution. *International Journal of Environmental Research and Public Health*, 18(15), 2021. <https://doi.org/10.3390/ijerph18158233>
- [10] Paluchová, J., Berčík, J., and Horská, E. The sense of smell. In *Sensory and Aroma Marketing* (pp. 27-60). Wageningen Academic, 2017. https://doi.org/10.3920/978-90-8686-841-4_2
- [11] Steinemann, A. Ten questions concerning air fresheners and indoor built environments. *Building and Environment*, 111, 279-284, 2017. <https://doi.org/10.1016/j.buildenv.2016.11.009>

- [12] Marques, F. B., Bettoni, G. N., dos Santos, B. G. T., Adeoye, A. A., de Brito, B. G., de Brito, K. C. T. and Cavalli, L. S. AquaSafe: Aquaculture occupational safety and health in the palm of your hand. *Pesquisa Agropecuária Gaúcha*, 26(1), 46-54, 2020. <https://doi.org/10.36812/pag.202026146-54>
- [13] Patino, E. D. L., and Siegel, J. A. Indoor environmental quality in social housing: A literature review. *Building and Environment*, 131, 231-241, 2018.
<http://dx.doi.org/10.1016/j.buildenv.2018.01.013>
- [14] Hable-Khandekar, V., and Srinath, P. Machine learning techniques for air quality forecasting and study on real-time air quality monitoring. In 2017 International Conference on Computing, Communication, Control and Automation (ICCUBEA) (pp. 1-6). IEEE, 2017.
<http://dx.doi.org/10.1109/ICCUBEA.2017.8463746>
- [15] Younis, M. T., Hussien, N. M., Mohialden, Y. M., Raisian, K., Singh, P., and Joshi, K. Enhancement of ChatGPT using API wrapper techniques. *Al-Mustansiriyah Journal of Science*, 34(2), 82-86, 2023.
<https://doi.org/10.23851/mjs.v34i2.1350>
- [16] Saini, J., Dutta, M., and Marques, G. Machine Learning for Indoor Air Quality Assessment: A Systematic Review and Analysis. *Environmental Modeling and Assessment*, 1-18, 2024.
<http://dx.doi.org/10.1007/s10666-024-10001-1>
- [17] Saini, J., Dutta, M., & Marques, G. Sensors for indoor air quality monitoring and assessment through Internet of Things: a systematic review. *Environmental Monitoring and Assessment*, 193(2), 66, 2021.
<https://doi.org/10.1007/s10661-020-08781-6>
- [18] Zhang, H., & Srinivasan, R. A systematic review of air quality sensors, guidelines, and measurement studies for indoor air quality management. *Sustainability*, 12(21), 9045, 2020.
<https://doi.org/10.3390/su12219045>
- [19] E. A. W. Hachim, M. T. Gaata and T. Abbas. Iris-based Authentication Model in Cloud Environment (IAMCE). *International Conference on Electrical, Computer and Energy Technologies (ICECET)*, IEEE, 2022.
<http://dx.doi.org/10.1109/ICECET55527.2022.9873499>

Facial Recognition Technology for Scenic Spot Monitoring Based on U²Net and FFC

Yiyi Liu

Tourism Management Department, Zhengzhou Tourism College, Zhengzhou 450000, China

E-mail: lyyxlflqh@163.com

Keywords: scenic spot monitoring, facial recognition, U²Net, fast fourier convolution, global convolution, mask learning

Received: June 30, 2025

To ensure the safety of scenic spots and achieve intelligent management of scenic spots, a face recognition method based on U²Net and FFC is proposed to achieve monitoring face recognition under different occlusion conditions. It consists of a small area regular occlusion model and a large area irregular occlusion face recognition model. Firstly, a face recognition model grounded on an improved residual network-U²Net is raised to address the problem of small area rule occlusion. This model combines a global convolution module, a feature pyramid network, and a mask learning unit. When evaluating facial recognition methods, multiple evaluation metrics were used, including recognition accuracy, F1-score, recognition rate, structural similarity index, peak signal-to-noise ratio, learning perceptual image block similarity, and Frecht approximation distance. These indicators evaluate the performance of the model under small and large area irregular occlusion conditions from different perspectives, ensuring the comprehensiveness and reliability of the evaluation. The findings denote that the average recognition accuracy of the enhanced residual network-U²Net is as high as 98.7%, the average F1-score is 0.983, and the average recognition rate is 99.5%. Secondly, in response to the problem of large-scale irregular occlusion in facial recognition, a fast Fourier convolution generative adversarial network is proposed, which combines generative adversarial network and Fourier feature convolution to repair and recognize facial images. The outcomes denote that the average structural similarity index and peak signal-to-noise ratio of the model are 0.878 and 34.7dB, respectively, and the average accuracy and recognition rate are 91.0% and 92.6%, respectively. The above results denote that the proposed facial recognition method exhibits superior performance under different occlusion conditions and can effectively promote the intelligent development of scenic area management.

Povzetek: Predstavljena je metoda prepoznavanja obrazov na podlagi U2Net in FFC za inteligentno nadzorovanje v turističnih krajih, ki omogoča prepoznavanje tudi pri zakritih obrazih (maske, klobuki).

1 Introduction

As the global tourism industry quickly develops in recent years, the number of tourists in scenic spots has shown explosive growth, which has put forward higher requirements for the management and service of scenic spots. Due to the low efficiency of traditional manual management models and their inability to cope with the complex and changing challenges of scenic areas, coupled with the dense population and high mobility of people in scenic areas, which significantly increase the difficulty of safety management, it is particularly important to introduce advanced safety management technologies to improve the level of safety management and management efficiency in scenic areas [1-2]. Among numerous security management technologies, facial recognition technology has gradually become an important tool for scenic spot security management due to its high efficiency, convenience, and accuracy. Through facial recognition systems, scenic spots can achieve real-time monitoring, rapid identification, and precise management of personnel, effectively enhancing emergency response

capabilities and ensuring the safety of tourists [3-4]. However, the complex environment of scenic spots, frequent personnel flow, and often the presence of obstructions greatly increases the difficulty of facial recognition. However, existing facial recognition technologies have low recognition accuracy when dealing with occlusion problems, making it difficult to meet practical needs.

Qin et al. proposed a multi-purpose algorithm called SwinFace-based on Swin Transformer to address the issue of neglecting task collaboration during the training process of facial recognition models. This method integrated multi-level channel attention modules in each task-specific analysis subnet with the objective of achieving adaptive feature selection. The findings demonstrated that the facial expression recognition and age estimation performance of this method surpassed that of existing methods [5]. Al-Dabbas et al. developed a facial recognition method that utilized classification, machine learning and deep learning models to address the issue of rising counterfeit crime rates. The methodology

employed involved the utilization of Viola Jones, linear discriminant analysis, mutual information, and analysis of variance techniques to construct two facial classification systems. The findings showed that the classification accuracy of both facial classification systems was above 96%, indicating that the proposed model performed well in both accuracy and processing time [6]. Gao et al. proposed the first privacy preserving facial recognition protocol for recognition stage computation in intelligent security systems to address privacy protection and identity recognition efficiency issues. This method introduced a Householder matrix into blind user data, enabling the protocol to support privacy protected facial recognition on semi trusted edge servers. The results showed that the protocol not only protected the privacy of user data, but also could achieve rapid response of large-scale face recognition (FR) through edge computing, effectively improving the efficiency of FR in intelligent security systems [7]. Xie et al. proposed a general privacy protection framework for FR systems that is grounded on edge computing. The purpose of this framework was to address the issue of data privacy leakage that has been identified in such systems. The overarching objective of the proposed framework is to safeguard the confidentiality of facial images and training models by employing a local differential privacy algorithm. The algorithm under discussion is founded upon a comparison of the proportion of feature information. As previously stated, the aforementioned text is concerned with the implementation of identity authentication and hashing techniques, with a view to confirming the legitimacy of terminal devices. The results showed that in numerical experiments, this scheme could ensure the optimal balance between the usability and privacy protection of the facial recognition system [8].

U²Net, as a deep learning model for image segmentation, combines an encoder and decoder, and introduces a cyclic squeezing unit, which can effectively extract image features of different scales. Therefore, it has significant advantages in image feature extraction. Feng et al. designed a detection method based on crack-U²Net to address the accuracy issue of road crack detection. This method utilized the U²Net architecture for feature learning and introduced a geometry-based data augmentation strategy to address the issue of insufficient training data. The results showed that the accuracy of Crack-U²Net in highway crack detection reached 95.8%, which is superior to existing methods [9]. Shi et al. proposed the U²CrackNet detection method for road crack detection. The proposed methodology involved the extraction of crack features through the encoding layer, followed by the establishment of a connection between the encoder and decoder via the atrous spatial pyramid pool model, with the objective of capturing multi-scale crack information. The results showed that U²CrackNet could obtain clearer and more continuous highway cracks, with a detection accuracy of 98.95% [10]. Li et al. proposed a U²Net-based analysis method to address the issue of low efficiency

caused by manual operation in microscope image analysis. This method enhanced the model's ability to extract key information by introducing a convolutional block attention module, and achieved model lightweighting by introducing Ghost convolution. The outcomes denoted that the prediction accuracy of the method model increased from 92.24% to 97.13% [11]. Zheng Z and Yang K proposed a detection method that integrates You Only Look Once version 5 (YOLOv5) and U²Net for wall crack detection. This method utilized the GhostNet module to optimize YOLOv5 to improve its training speed, while introducing U²Net to perform binary classification on the region input extracted by YOLOv5 to enhance the final classification performance. The results denoted that this method could effectively address the issue of poor segmentation of crack targets in large environmental backgrounds [12].

In summary, although the current facial recognition models have high recognition accuracy, they are difficult to cope with the problem of obstructed facial recognition in complex environments of scenic spots. Therefore, to address the above issues, an FR model for different occlusion conditions has been proposed, which consists of two parts: a small area regular occlusion FR model and a large area irregular occlusion FR model. The innovation of the research lies in the combination of Residual Network (ResNet) and U²Net, and the introduction of Global Convolution Module, Feature Pyramid Network (FPN), and Mask Learning Unit to improve the recognition accuracy of the model for small area regularly occluded faces. Specifically, a global convolution module consisting of two symmetric convolution layers is used to capture global features in both horizontal and vertical directions. At the same time, a mask learning unit is introduced in FPN to remove the features of occluded areas by generating multi-level masks to enhance feature representation. Secondly, by combining Fast Fourier Convolution (FFC) and Generative Adversarial Network (GAN), the problem of low recognition accuracy in large-area irregularly occluded face images can be solved by repairing them.

2 Methods and materials

Due to the influence of facial coverings such as hats, masks, and glasses, as well as changes in facial expressions, the success rate of existing surveillance facial recognition is low, making it difficult to effectively ensure the safety of scenic spots. Therefore, a monitoring FR model based on U²Net and FFC is proposed to address the recognition problem under face occlusion. It consists of two parts: a small area occlusion FR model and a large area irregular occlusion FR model. Firstly, an FR method based on U²Net and global convolution is constructed to address the problem of small-scale rule-based occlusion. For the problem of large-scale irregular occlusion, a FR model based on FFC and GAN is proposed.

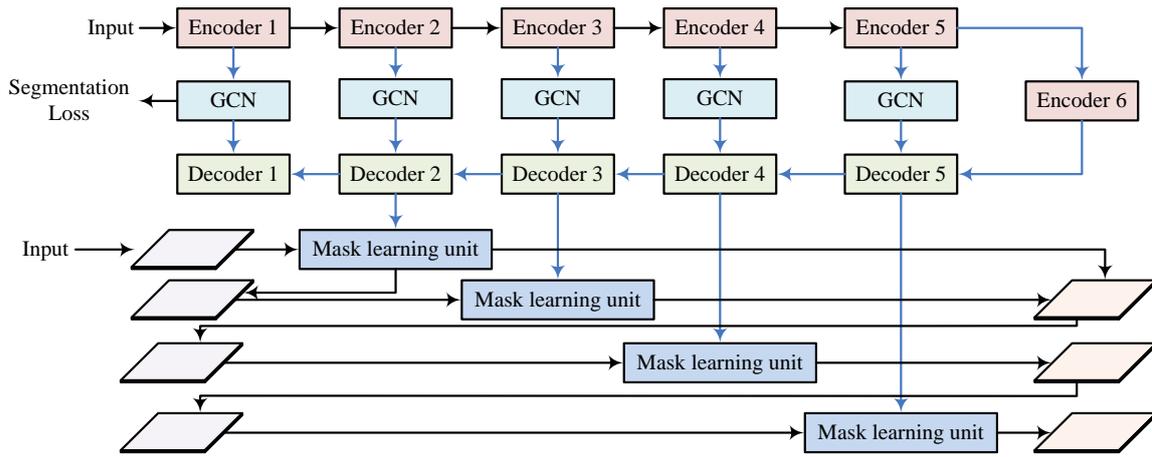


Figure 1: Face recognition model based on improved ResNet-U²Net. (Source from: Author's self drawn)

2.1 Regular occlusion facial recognition model based on U²Net and global convolution

It is difficult to extract facial features from surveillance cameras in conditions of obstruction by regular objects, such as masks, glasses and hats. The accuracy of facial recognition systems is consequently adversely affected. Due to the use of two parallel encoder decoder paths and the introduction of cyclic squeezing units, U²Net is able to simultaneously process global and local features, thereby improving segmentation accuracy. Although the cyclic squeezing unit can capture multi-scale features, its receptive field size is small, which makes it impossible to fully cover all scale features [13-14]. Therefore, to expand the receptive field of U²Net, global convolution is introduced and improved. The FR model based on improved ResNet-U²Net is denoted in Figure 1.

In Figure 1, the FR model based on improved U²Net consists of two parts: occlusion detection segmentation module and feature detection module. The model first generates a multi-level occlusion segmentation map through the occlusion detection module, then extracts image features through the FR module, and removes the influence of occlusion on facial features through the mask learning unit. Finally, FPN is used to fuse the features of each stage. In the feature extraction module, the selected backbone network is ResNet, which can achieve feature reuse through skip connections. However, due to the poor ability of ResNet to extract multi-scale features, it will reduce the accuracy of the model's FR. Therefore, to enhance the multi-scale feature extraction capability and model generalization performance, the FPN module is introduced in the study. FPN upsamples high-level feature maps to the resolution of low-level feature maps through a top-down path, thereby generating a multi-scale feature pyramid. Moreover, feature maps of different scales are fused through horizontal connections to enhance the richness of feature representation. At the same time, the generated feature pyramid can capture both global and local information, improving the model's ability to extract

multi-scale features [15-16]. For improved ResNet-U²Net, the input image needs to be first detected and aligned by the method based on the cascaded multi task framework, and then the image size is adjusted to 112 * 112. Next, facial recognition can be performed using the improved ResNet-U²Net. The Batchsize of the model is 128, the initial learning rate is 0.1, the hypersphere radius *s* is 64, and the spacing *m* is 0.48. During the training, the learning rate is adjusted to 1/10 of its original level at the 11th, 20th, and 30th epochs. However, FPN is prone to information loss, which can lead to feature damage. Therefore, to improve the above problems, the structure of FPN is optimized by introducing mask learning units to avoid the influence of damaged features. The formula for calculating the feature pyramid is denoted in equation (1).

$$\begin{cases} P_3 = M_3 + ds(M_2) \\ P_4 = M_4 + ds(P_3) \\ P_5 = M_5 + ds(P_4) \end{cases} \quad (1)$$

In equation (1), *P_i* means the feature pyramid; *M_i* represents the features obtained after mask operation; *ds* represents downsampling operation. The calculation formula for *M_i* is denoted in equation (2).

$$M_i = X_i + Mask_i \otimes X_i \quad (2)$$

In equation (2), *Mask_i* represents the mask; *X_i* represents the original input. For the occlusion detection and segmentation module, its backbone network is U²Net, which consists of U-shaped residual modules. Unlike ordinary residual modules, the U-shaped residual module replaces the convolutional layers in the original residual module with U-blocks and replaces the original features with local features to achieve multi-scale feature extraction. The so-called U-block refers to the U-shaped encoder decoder structure. Considering the complexity of the model, the amount of U-shaped residual modules is 6 [17]. Due to the small receptive field of U²Net, a global convolution module is introduced to expand its receptive field size. The structure of the global convolution module is denoted in Figure 2.

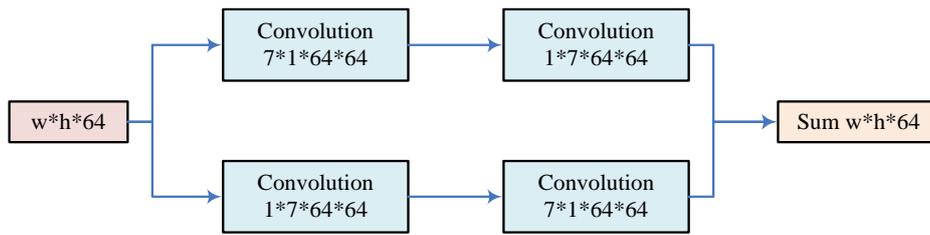


Figure 2: Structure of the global convolution module. (Source from: Author's self drawn)

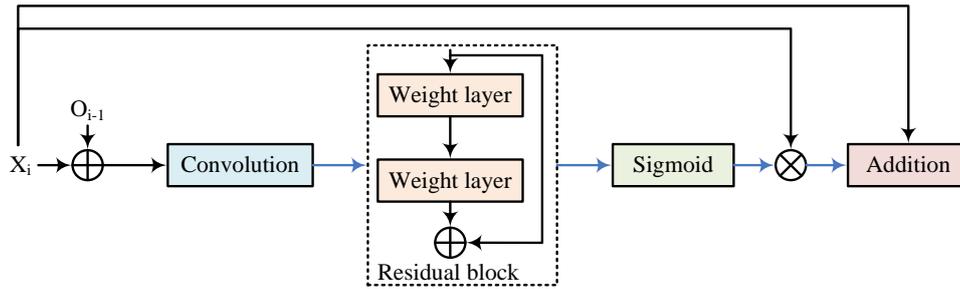


Figure 3: Structure of the mask learning unit. (Source from: Author's self drawn)

In Figure 2, the global convolution module mainly contains four symmetrically distributed convolution layers. This structure enables the global convolution module to capture features of different scales and enhance the model's understanding of the global information of the image through feature fusion. In addition, to avoid damaging the features and affecting the performance of the occlusion detection segmentation module and feature detection module, a mask learning unit is introduced in the study. Although generating masks can remove the features of occluded areas, these methods usually only generate masks at a single scale and cannot effectively handle multi-scale features. Moreover, the mask learning unit can generate multi-level masks, corresponding to feature maps of different scales, thus more comprehensively handling occlusion problems. It suppresses the features of occluded areas through masking while preserving the features of unobstructed areas, enhancing the robustness of feature representation. The structure of the mask learning unit is denoted in Figure 3.

In Figure 3, the mask learning unit mainly contains convolutional layers, residual modules, and sigmoid functions. This module first concatenates the feature and occlusion segmentation representations of each stage, and then processes the concatenated images using convolutional layers and activation functions to generate multi-level masks. Finally, the generated mask is used to remove the features of the occluded area and added to the original input features to enhance the feature representation. The formula for mask learning calculation is shown in equation (3).

$$Mask_i = \text{Sigmoid}(\gamma(c(\text{concat}[F_i, S_{i-1}])) \quad (3)$$

In equation (3), γ represents residual operation; F_i represents the characteristics of each stage; S_i represents occlusion segmentation representation. The loss function (LF) of the model is denoted in equation (4).

$$L = L_{fc} + L_{seg} \quad (4)$$

In equation (4), L means the overall LF of the model; L_{fc} denotes the face classification LF; L_{seg} denotes the face segmentation LF. The calculation formula for the face classification LF is denoted in equation (5).

$$L_{fc} = -\frac{1}{B} \sum_{i=1}^B \ln \frac{e^{\|x_i\|(\cos(\theta_{y_i} + s))}}{e^{\|x_i\|(\cos(\theta_{y_i} + s))} + \sum_{j=1, j \neq y_i}^N e^{\|x_i\| \cos \theta_j}} \quad (5)$$

In equation (5), B represents batch size; x_i represents the feature vector; s represents spacing; N means the number of categories; θ_j means the angle between the weight and the feature vector. The face segmentation LF is shown in equation (6).

$$L_{seg} = \frac{1}{|N|} \sum_{c \in C} \left[\varepsilon D_{KL}(y \| \hat{p}_c) + \frac{\delta}{|c|} \sum_{s \in c} D_{KL}(\hat{p}_c \| p(s)_c) \right] \quad (6)$$

In equation (6), ε and δ both represent hyperparameters; D_{KL} represents Kullback-Leibler divergence; \hat{p}_c means the probability distribution of category c ; $p(s)$ represents the probability distribution of the category c of real data.

2.2 Irregular occlusion face recognition model based on FFC and GAN

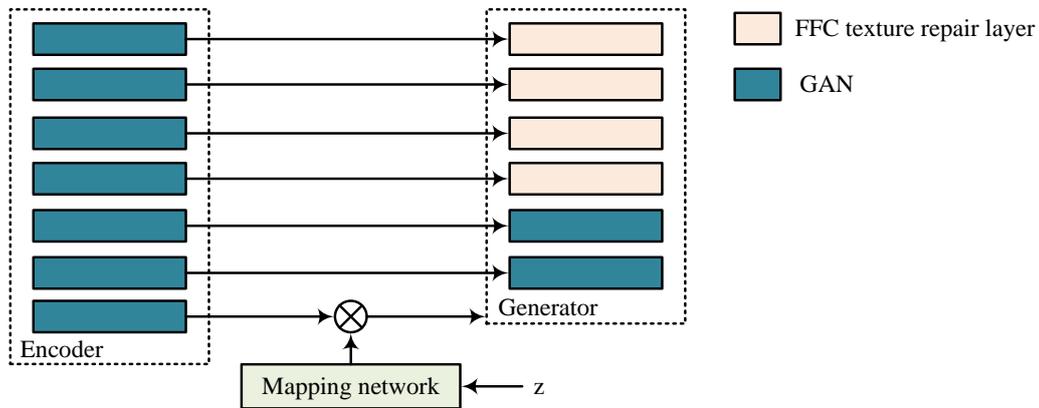


Figure 4: Monitoring face recognition model based on FFC-GAN. (Source from: Author's self drawn)

Although the above method can achieve FR with small area regular occlusion, due to the large flow of people and complex environment in scenic spots, there are cases where faces are obstructed by large irregular objects, which further increases the difficulty of FR. Therefore, to effectively ensure the safety of scenic spots, a large-scale irregular occlusion FR model based on FFC and GAN is proposed. Compared to other convolution methods, FFC can effectively accelerate the training and inference process of networks when processing large image or video data [18-19]. GAN can generate high-quality synthetic data to achieve the restoration of large areas of irregularly occluded faces. The monitoring FR model based on FFC-GAN is denoted in Figure 4.

As shown in Figure 4, the model first uses GAN to repair irregularly occluded facial images, and generates facial image structures using encoding and hidden layer noise vectors. Then FFC is utilized to generate texture details of the image to raise the quality of facial image restoration. Finally, the model is jointly trained using an identity preservation LF to raise the accuracy of FR. For GANs, the input is a random noise vector. This is mapped to the data space through a series of neural network layers to generate fake data. The required style parameters are then generated based on affine changes. The formula for generating style parameters is shown in equation (7).

$$s = A(M(h)) \tag{7}$$

In equation (7), s represents the style parameter; A represents affine transformation; M stands for Mapping Network; h stands for hidden layer vector. Although the above method can achieve the restoration of occluded images, it may result in inconsistency between the restored image and the original image. Therefore, to solve the above problems, collaborative modulation methods are introduced in the research. The formula for generating style parameters for collaborative modulation is shown in equation (8).

$$s = A(E(x), M(h)) \tag{8}$$

In equation (8), E represents the image conditional encoder; x represents the input image. It is worth noting that the generator and discriminator of GAN need to be trained alternately. The goal of the generator is to generate restored images that are as close to the real image as possible, while the goal of the discriminator is to distinguish between the generated image and the real image. Therefore, in each iteration, the generator and discriminator update their parameters separately to minimize the adversarial LF. The Batch size of GAN is 24, with an initial learning rate of 0.002, and the learning rate is adjusted to 0.001 after 650000 iterations. The weight of reconstruction loss is 10, and the weight of identity preservation loss is 10. The above method can achieve the restoration of large-area occluded images, but due to the loss of texture details in the restored images, it seriously affects the success rate of FR. Therefore, to achieve the restoration of image texture details, the FFC module is introduced in the study. Although existing texture restoration methods, such as convolution-based restoration methods, can generate certain texture details, they have low efficiency in processing large-scale images and are difficult to effectively capture global features. FFC, through the fusion of global and local features, can generate higher quality texture details and improve the quality of restored images. The structure of FFC is shown in Figure 5.

As shown in Figure 5, FFC consists of global branches and local branches. FFC first splits the input features into global features and local features, where global features are processed through convolutional layers and Spectral Transformers, and local features are processed using two convolutional layers [20-21]. Next, the processed local features and global features are fused, and after batch normalization and ReLU processing, the output features can be obtained. The formula for calculating the local output features of FFC is denoted in equation (9).

$$Y^l = Y^{l \rightarrow l} + Y^{g \rightarrow l} = Fou_l(X^l) + Fou_{g \rightarrow l}(X^g) \tag{9}$$

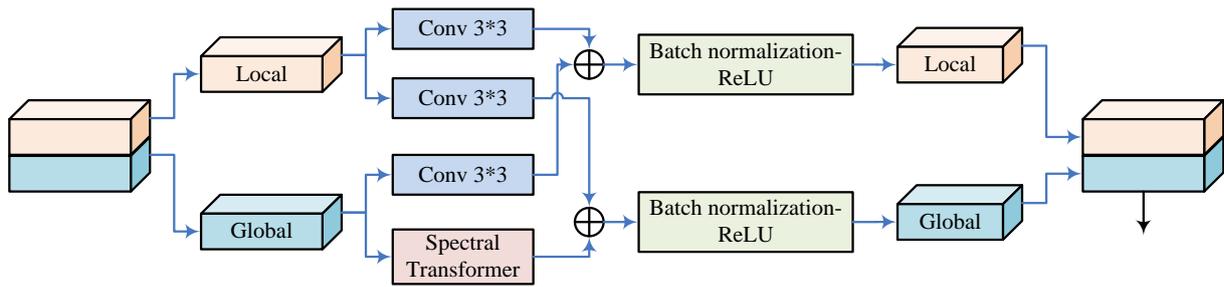


Figure 5: Structure of FCC. (Source from: Author's self drawn)

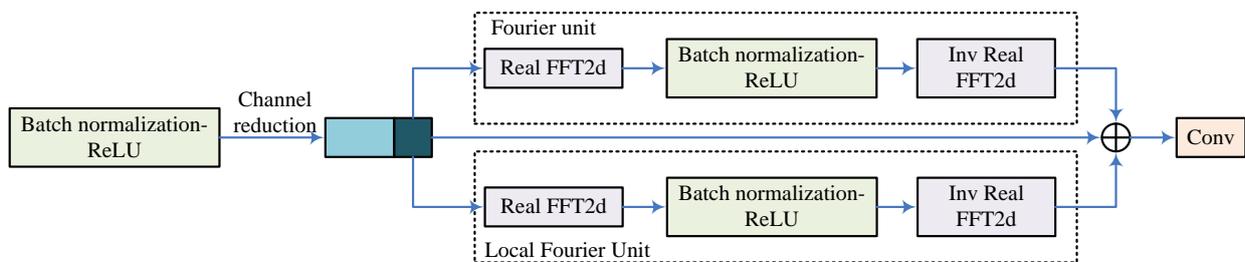
In equation (9), Y^l represents the output characteristics of local branches; $Y^{l \rightarrow l}$ represents small-scale feature components of local branches; $Y^{g \rightarrow l}$ represents the multi-scale receptive field components exchanged from global branches to local branches; Fou_l and $Fou_{g \leftarrow l}$ both represent fast Fourier transform (FFT); X^l and X^g represent the input features of local and global branches, respectively. The formula for calculating the global output characteristics of FCC is shown in equation (10).

$$Y^g = Y^{g \rightarrow g} + Y^{l \rightarrow g} = Fou_g(X^g) + Fou_{l \rightarrow g}(X^l) \quad (10)$$

In equation (10), Y^g represents the output feature of the global branch; $Y^{g \rightarrow g}$ represents the small-scale feature components of the global branch; $Y^{g \rightarrow g}$ represents the multi-scale receptive field components exchanged from local branches to global branches; Fou_g

and $Fou_{l \leftarrow g}$ both represent FFT. The structure of the Spectral Transformer in FCC is shown in Figure 6.

In Figure 6, the Spectral Transformer includes convolutional layers, Fourier units, and local Fourier units. Firstly, Spectral Transformer processes input information through convolutional and batch normalization layers, and then captures global and local features using Fourier units and local Fourier units, respectively, and fuses the features. Finally, the captured features can be output after being convolved again. The Fourier unit and local Fourier unit are both composed of real 2D FFT, convolutional layer, and inverse real two-dimensional FFT. The real 2D FFT is responsible for transforming spatial features into the spectral domain, the convolutional layer is responsible for updating spectral data, and the inverse real 2D FFT is responsible for restoring spatial features [22-23]. By using the above method, FCC is constructed, and after combining it with convolutional layers, a texture restoration module based on FCC can be constructed. The structure of the texture restoration module based on FCC is shown in Figure 7.



Note: Real FFT2d represents Real 2D Fast Fourier Transform, Inv Real FFT2d represents Inverse Real 2D Fast Fourier Transform

Figure 6: Structure of spectral transformer. (Source from: Author's self drawn)

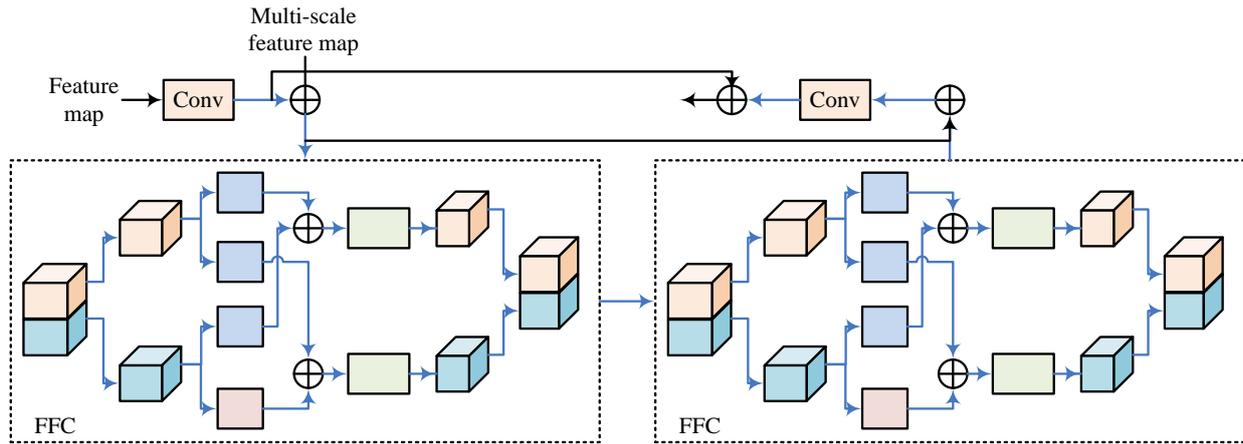


Figure 7: Texture repair module based on FFC. (Source from: Author's self drawn)

As shown in Figure 7, the FFC-based texture restoration module consists of convolutional layers and FFC residual structures. This module first processes the output features of the previous stage through convolutional layers and fuses them with the feature maps extracted by the image condition encoder. Next, the fused feature maps are subjected to contextual information extraction and fusion using the FFC residual structure, and the information is processed using convolutional layers. Finally, the processed information is fused with the features processed by the first convolutional layer to achieve texture restoration of the image. A monitoring FR model based on FFC-GAN is constructed using the above method, and the LF of the model is denoted in equation (11).

$$L_G = L_{gen} + L_{ref} + L_{id} \tag{11}$$

In equation (11), L_G denotes the overall LF of FFC-GAN; L_{gen} represents the adversarial LF of the generator; L_{gen} represents the reconstruction LF; L_{id} stands for identity preservation LF. The calculation formula for the adversarial LF is denoted in equation (12).

$$L_{gen} = -E_{I_{res}} [\log D(I_{res})] \tag{12}$$

In equation (12), $E_{I_{res}}$ represents the expected value of the restored image; D stands for discriminator; I_{res} represents the restored image. The calculation formula for the reconstruction LF is denoted in equation (13).

$$L_{ref} = \alpha \|I_{res} - I_{ori}\|_1 \tag{13}$$

In equation (13), α represents the weight of reconstruction loss; I_{ori} represents the original image. The identity preservation LF is shown in equation (14).

$$L_{id} = \beta \|F(I_{res}) - F(I_{ori})\|_1 \tag{14}$$

In equation (14), β represents the weight of identity preservation loss; $F(\cdot)$ represents the feature extraction process. The above method can achieve accurate recognition of faces with large areas of irregular occlusion.

3 Results

3.1 Small area occlusion face recognition test results

To test the recognition effect of the improved ResNet-U²Net proposed in the study for small area regular occlusion faces, it was tested and compared with the Fine-Grained Deep Feature Mask Estimation (FGDFME) occlusion FR algorithm and the Depth Image Priors and Robust Markov Random Fields (DIP-rMRF) occlusion FR algorithm based on depth image priors and robust Markov random fields. The datasets used in the experiment were the Labeled Faces in the Wild (LFW) dataset and the Masked Faces in Real World for Face Recognition (MFR2) dataset used for FR in the real world. The LFW dataset contains 13233 facial images, covering 5749 individuals of different identities. Each image is labeled with the name of the corresponding person, with 1680 individuals having two or more images. Meanwhile, each image has a size of 250 * 250 pixels, with the majority being color images, but there are also a few black and white facial images. The MFR2 dataset contains the identities of 53 celebrities and politicians, with a total of 269 images. The size of each image is 160 * 160 * 3. To ensure the reliability of the experimental results, a simulated occlusion dataset was constructed using the LFW dataset, which involves adding objects such as masks, sunglasses, and mobile phones to mask facial images. The CPU utilized in the experiment was Intel core i7 4720HQ, with 16GB of memory and GeForce RTX 4060Ti GPU. The Batchsize and initial learning rate of the model were 128 and 0.1, respectively, and the radius and spacing of the hypersphere were 64 and 0.48, respectively. For each evaluation metric, the mean and standard deviation of multiple experimental results was calculated to assess the stability and reliability of the model performance. The 95% confidence interval to evaluate the confidence level of the model performance. The recognition accuracy and F1-score of each model in the simulated occlusion dataset are shown in Figure 8.

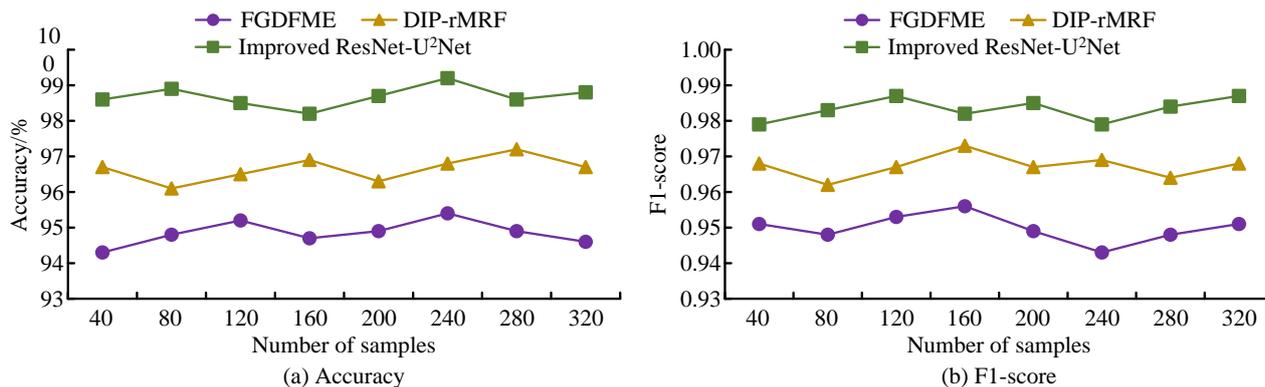


Figure 8: The recognition accuracy and F1-score of each model in the simulated occlusion dataset. (Source from: Author's self drawn)

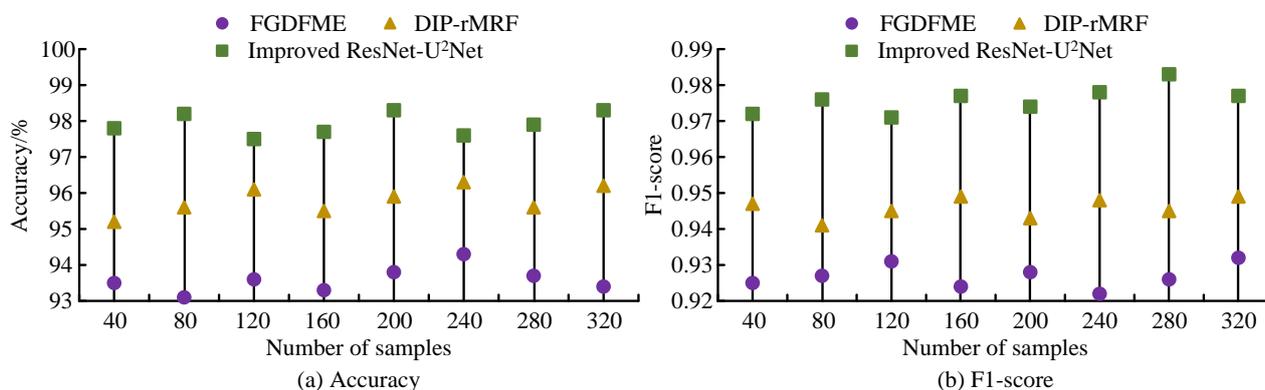


Figure 9: Recognition accuracy and F1-score of each model in MFR2 data set. (Source from: Author's self drawn)

In Figure 8 (a), in the simulated occlusion dataset, the facial recognition accuracy of FGDFME and DIP rMRF was the highest at 95.4% and 97.2%, the lowest at 94.3% and 96.1%, and the average accuracy was 94.9% and 96.7%, respectively. The improved ResNet-U²Net had a minimum FR rate of 98.2% and an average accuracy rate of up to 98.7%, which was higher than other algorithms. From Figure 8 (b), in the simulated occlusion dataset, the F1-score of FGDFME and DIP rMRF were the highest at 0.956 and 0.973, and the lowest at 0.943 and 0.962, respectively. The average F1-score was 0.950 and 0.967, respectively. The lowest F1-score of ResNet-U²Net improvement was 0.979, with an average F1-score of 0.983. The above outcomes denoted that the improved ResNet-U²Net had good performance in small area rule-based occlusion FR. The recognition accuracy and F1-score of each model in the MFR2 dataset are shown in Figure 9.

From Figure 9 (a), in the MFR2 dataset, the highest facial recognition accuracy of FGDFME and DIP rMRF was 94.3% and 96.3% respectively, the lowest was 93.1% and 95.2% respectively, and the average accuracy was 93.6% and 95.8% respectively. The improved ResNet-U²Net had a minimum FR rate of 97.6% and an average accuracy rate of 97.9%, which was higher than other algorithms. From Figure 9 (b), in the MFR2 dataset, the

highest F1-score for FGDFME and DIP rMRF were 0.956 and 0.973, and the lowest were 0.943 and 0.962, respectively. The average F1-score was 0.950 and 0.967, respectively. The lowest F1-score of ResNet-U²Net improvement was 0.979, with an average F1-score of 0.983. The True Acceptance Rate (TAR) of each model in different datasets is shown in Figure 10.

According to Figure 10 (a), in the simulated occlusion dataset, the highest TAR of FGDFME and DIP rMRF were 96.2% and 98.3%, respectively, and the lowest were 95.3% and 97.1%, respectively. The average TAR was 95.7% and 97.6%, respectively. The TAR of ResNet-U²Net was improved from a mini of 99.2% to a max of 99.9%, with an average TAR of 99.5%. According to Figure 10 (b), in the MFR2 dataset, the highest and lowest TARs for FGDFME and DIP rMRF were 95.3% and 96.8%, respectively, and 94.1% and 95.8%, respectively, with an average TAR of 94.6% and 96.3%. The TAR of the improved ResNet-U²Net ranged from 98.1% to 99.9%, with an average TAR of 98.5%. The above results indicated that the improved ResNet-U²Net had strong facial recognition capabilities and could effectively ensure the safety of scenic spots. To further analyze and improve the performance of ResNet-U²Net, ablation experiments were conducted on it. The outcomes of the ablation experiment are denoted in Table 1.

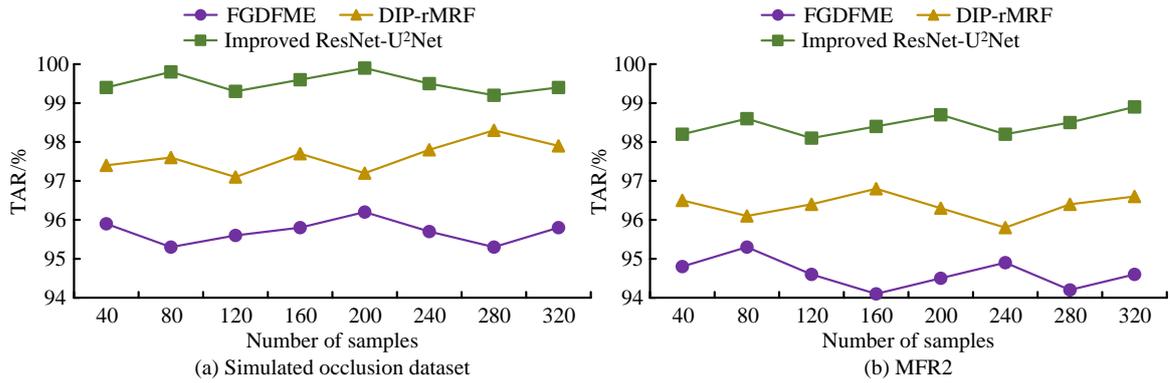


Figure 10: TAR of each model in different data sets. (Source from: Author's self drawn)

Table 1: Results of ablation experiments.

Model	ResNet	FPN	U²Net	Mask learning unit	Accuracy/%
1	√	×	×	×	92.2
2	√	√	×	×	93.1
3	√	×	√	×	94.2
4	√	×	×	√	94.5
5	√	√	√	×	95.6
6	√	√	×	√	96.7
7	√	×	√	√	97.4
8	√	√	√	√	98.7

According to Table 1, the facial recognition accuracy of the backbone network ResNet was only 92.2%. After introducing FPN, U²Net, and mask learning units, the facial recognition accuracy of the model significantly improved. Among them, U²Net and mask learning units had the most significant impact on model performance. After introducing the above two modules, the facial recognition accuracy of the model increased to 94.2% and 94.5%, respectively.

3.2 Large area occlusion face recognition test results

To test the effect of the proposed FFC-GAN in repairing and recognizing large-area irregularly occluded faces, it was tested and compared with the Partial Convolution and Multiscale Feature Fusion (PCMSF) facial image restoration model, Multiscale Feature Fusion U-Net (MSFFU-Net), Involution facial Feature Correction

Network (IFFR-Net), and Depth Separable Convolution and Hypersphere Loss (DSCHL) occlusion model. The software and hardware settings of the experiment are consistent with the above experiment and will not be repeated. The dataset utilized in the experiment was the CelebA HQ dataset, which contains 30000 facial images with a resolution of 1024 × 1024. To simulate irregular occlusion situations, various shapes were randomly used to occlude facial images, with an occlusion rate of over 50%. In the experiment, the reconstruction loss weight and identity preservation loss weight were both 10, and the initial learning rate and Batchsize were 0.002 and 24, respectively. Firstly, the facial image restoration performance of FFC-GAN was tested. The Structural Similarity Index Measure (SSIM) and Peak Signal-to-Noise Ratio (PSNR) of different models are shown in Figure 11.

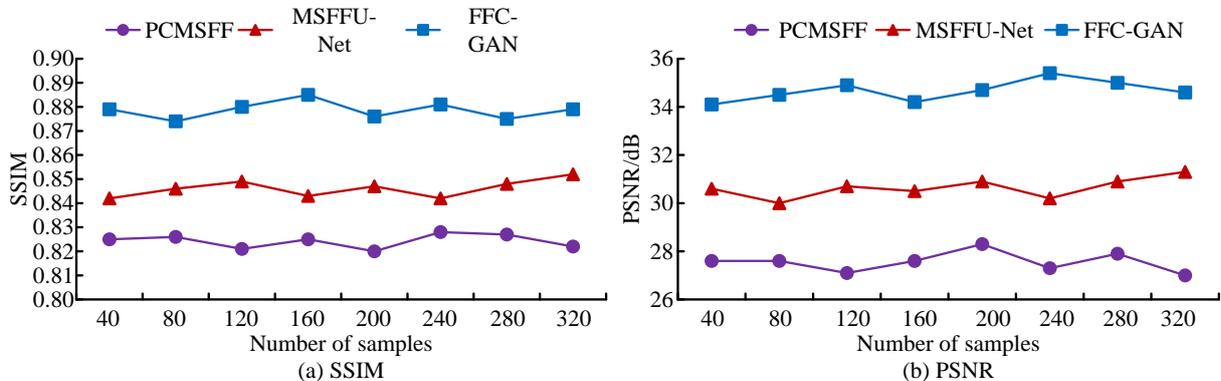


Figure 11: SSIM and PSNR of different models. (Source from: Author's self drawn)

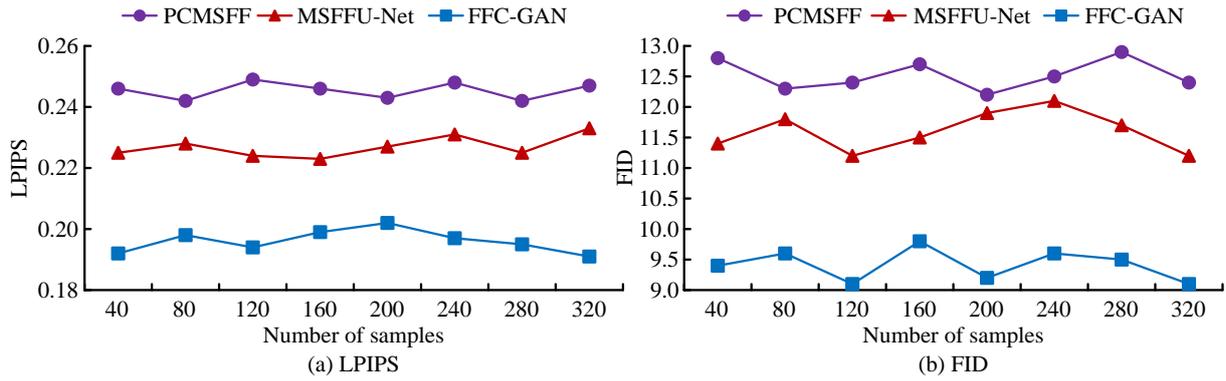


Figure 12: LPIPS and FID of different models. (Source from: Author's self drawn)

From Figure 11 (a), the SSIMs of PCMSFF and MSFFU-Net were the highest at 0.828 and 0.852, the lowest at 0.820 and 0.842, and the average SSIMs were 0.824 and 0.846, respectively. The SSIM of FFC-GAN was the lowest at 0.874, with an average SSIM of 0.878, which was higher than other methods. From Figure 11 (b), the PSNRs of PCMSFF and MSFFU-Net were the highest at 28.3dB and 31.3dB, and the lowest at 27.0dB and 30.0dB, respectively, with average PSNRs of 27.6dB and 30.6dB, respectively. The PSNR of FFC-GAN was the lowest at 34.1dB, with an average PSNR of 34.7dB, which was also higher than other algorithms. The above results indicated that the facial image restoration quality of FFC-GAN was superior to other algorithms. The Learned Perceptual Image Patch Similarity (LPIPS) and Frechet Inception Distance (FID) of different models are shown in Figure 12.

According to Figure 12 (a), the minimum and maximum LPIPS of PCMSFF and MSFFU-Net are 0.242 and 0.223, respectively, and 0.249 and 0.233, respectively. The average LPIPS was 0.245 and 0.227. The maximum LPIPS of FFC-GAN was 0.202, and the average LPIPS was 0.196, which was much lower than other methods. According to Figure 12 (b), the maximum FID of PCMSFF and MSFFU-Net were 12.9 and 12.1, and the minimum FID was 12.2 and 11.2. The average FID was

12.5 and 11.6, respectively. The maximum FID of FFC-GAN was 9.8, and the average FID was 9.4, which was also lower than other algorithms. The above results indicated that FFC-GAN could achieve high-quality restoration of large-area irregularly occluded face images. The facial recognition accuracy and TAR of different models are shown in Figure 13.

According to Figure 13 (a), the highest and lowest facial recognition accuracies of IFFR Net and DSCHL were 86.4% and 88.5%, respectively, and 84.7% and 87.3%, respectively. The average accuracies were 85.6% and 87.9%, respectively. The recognition accuracy of FFC-GAN was the lowest at 90.5%, with an average accuracy of 91.0%, which was higher than other algorithms. From Figure 13 (b), the TAR of IFF-Net and DSCHL were the highest at 88.3% and 90.3% respectively, the lowest at 87.2% and 89.1% respectively, and the average TAR was 87.7% and 89.6% respectively. The lowest TAR of FFC-GAN was 92.1, with an average TAR of 92.6%, which was also higher than other algorithms. The above results indicated that FFC-GAN could achieve accurate recognition of faces with large areas of irregular occlusion. To further analyze the performance of FFC-GAN, ablation experiments were conducted on it. The findings of the ablation experiment are denoted in Table 2.

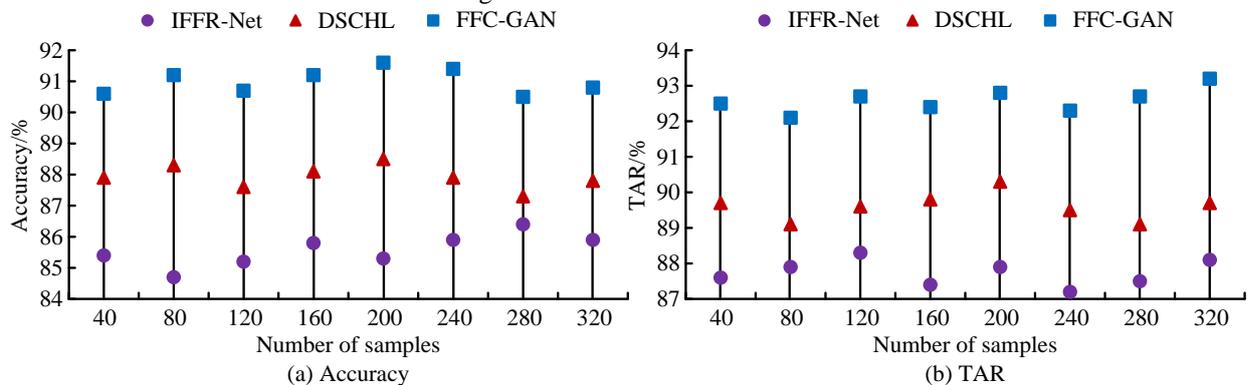


Figure 13: Face recognition accuracy and TAR of different models. (Source from: Author's self drawn)

Table 2: Ablation results.

Model	FFC	Spectral Transformer	Residual block	SSIM	Accuracy/%
1	×	×	×	0.725	84.4
2	√	×	×	0.796	88.2
3	×	√	×	0.771	87.5
4	×	×	√	0.795	87.9
5	√	√	×	0.827	88.9
6	√	×	√	0.846	89.2
7	×	√	√	0.859	89.8
8	√	√	√	0.878	91.0

According to Table 2, after introducing FCC, Spectral Transformer, and residual blocks, the SSIM and accuracy of the model significantly increased, reaching 0.878 and 91.0%, respectively. Among them, FFC and residual blocks had the most significant impact on model performance. After introducing FFC and residual blocks, the SSIM of the model increased to 0.796 and 0.795, respectively, and the accuracy increased to 88.2% and 87.9%, respectively.

4 Discussion

In recent years, with the booming development of the tourism industry, the number of tourists in scenic spots has been continuously increasing, which has brought many challenges to scenic spot management. The traditional management method of scenic spots has problems such as low efficiency, easy errors, and inability to monitor in real time, which not only affects the tourist experience but may also lead to safety hazards. The advent of artificial intelligence, computer vision, and deep learning technologies has precipitated a substantial enhancement in the security, convenience, and accuracy of facial recognition technology [24-25]. Real-time monitoring of personnel within the scenic area can be achieved through facial recognition technology, detecting abnormal behavior in a timely manner and issuing alerts. In addition, facial recognition systems can quickly locate missing persons or lost items, enhancing the emergency response capabilities of scenic spots. However, due to the complex environment and huge pedestrian flow in scenic spots, facial recognition is difficult [26]. Therefore, to achieve accurate recognition of faces in scenic area monitoring, an FR method based on improved ResNet-U²Net was proposed to address the problem of FR under small area rule occlusion such as sunglasses and masks. A recognition method based on FFC-GAN was proposed for the FR problem of large irregular occlusion.

For the improved ResNet-U²Net, experimental results showed that its average recognition accuracy and F1-score in simulated occlusion datasets were 98.7% and 0.983, respectively, with an average TAR of 99.5%, both higher than FGDFME and DIP rMRF. In the MFR2 dataset, the average recognition accuracy and F1-score of the improved ResNet-U²Net were 97.9% and 0.967, respectively, with an average TAR of 98.5%, which was also higher than other algorithms. Haider et al. designed a variational invariant FR method based on multi-task learning, which redefines FR by combining temporal dependence and temporal independence to decompose the

face into age and residual features. The experimental results showed that this method could achieve accurate recognition of faces of different races [27]. However, the above methods had low accuracy in recognizing occluded faces, while the proposed method could achieve accurate recognition of faces under objects such as masks and sunglasses. Akheel T S et al. proposed using optimized projection matrices in linear collaborative regression classification to improve recognition accuracy, and introduced a whale lion combination model to optimize the projection matrix. The findings denoted that the facial recognition accuracy of the model could reach 91.2% [28]. Compared to the above algorithms, the improved ResNet-U²Net proposed in the study had higher facial recognition accuracy. This is because the improved ResNet-U²Net introduces a global convolution module, allowing the model to capture a larger range of global information. Meanwhile, the model also introduced FPN, effectively enhancing its multi-scale feature extraction capability. In addition, the study also introduced a mask learning unit, which removes the features of occluded areas by generating multi-level masks to enhance feature representation.

For FFC-GAN, its average SSIM and average PSNR were 0.878 and 34.7 dB, respectively, which were higher than PCMSF and MSFFU-Net. The average LPIPS and FID were 0.196 and 9.4, respectively, which were lower than other algorithms. FFC-GAN could achieve accurate restoration of large-area irregularly occluded facial images. In terms of facial recognition performance, the average accuracy and TAR of FFC-GAN were 91.0% and 92.6%, respectively, both higher than existing advanced algorithms. Yan L. et al. proposed a methodology for optimizing image feature compensation coefficients. This methodology is based on an enhanced simulated annealing algorithm, the purpose of which is to enhance the recognition rate of facial recognition systems. The findings indicated that when the training image was designated as 6, the recognition rate attained a maximum of 100% [29]. Compared to the above methods, although the proposed method had lower recognition accuracy, it could effectively address the problem of large-scale irregular facial occlusion. Zaaaroui et al. put forward an FR method based on the mini value string, utilizing the mini value string as the face feature extractor for face representation. The findings demonstrated that the method exhibited high recognition accuracy and efficiency [30]. However, compared to the methods proposed in the research, the above methods significantly reduced the

accuracy of FR under large-scale irregular occlusion conditions. The reason why the proposed FFC-GAN can achieve accurate recognition of large-area irregularly occluded faces is that this method can accurately repair occluded images through GAN and accurately restore image texture details through FFC.

In summary, the improved ResNet-U²Net and FFC-GAN can achieve accurate recognition of occluded faces, among which the improved ResNet-U²Net has high recognition accuracy for small area regularly occluded faces. FFC-GAN can effectively repair large areas of irregularly occluded facial images, thereby achieving accurate facial recognition. The above two methods provide strong support for the development of facial recognition technology for scenic spot monitoring, which helps to achieve intelligent management of scenic spots. However, due to the high number of parameters and computational complexity of the proposed model, it requires high computing power from the server, making the deployment of the model difficult. Therefore, in the future, the model structure will be optimized to minimize the number of parameters and computational complexity of the model, so that it can be deployed on platforms with limited processing capabilities such as mobile devices and embedded devices.

5 Conclusion

A small area regular occlusion FR model based on improved ResNet-U²Net and a large area irregular occlusion FR model based on FFC-GAN were proposed to address the issue of FR in scenic spot monitoring. The improved ResNet-U²Net achieved accurate recognition of small area regularly occluded faces by introducing global convolution, FPN, and mask learning units. The findings denoted that the average recognition accuracy and F1-score of the improved ResNet-U²Net reached 98.7% and 0.983, respectively, with an average TAR of 99.5%. The FFC-GAN model utilized GAN and FFC modules to repair and recognize large-area irregularly occluded facial images. The findings denoted that the average SSIM and PSNR of the model were 0.878 and 34.7dB, respectively, and the average accuracy and TAR were 91.0% and 92.6%, respectively, which were better than existing advanced algorithms. The above results indicated that improved ResNet-U²Net and FFC-GAN could achieve accurate recognition of facial images under different occlusion conditions, providing strong support for the development of facial recognition technology for scenic spot monitoring. However, the model has high parameter count and computational complexity, which makes it impossible to deploy on mobile devices, greatly limiting its application scope. Therefore, in the future, the model will be lightweighted to reduce its complexity.

6 Funding

The research is supported by Research and Practice Project on Education and Teaching Reform of Henan Provincial Department of Education in 2024 (Project No. 2024SJGLX0837); The series of development

achievements of the 2024 Henan Province Higher Education Teaching Achievement "Practical Research on the Transformation and Application Mode of Tour Guide Service Skills Competition" Competition Teaching Post "Promotion of General Education" (Achievement Number: Yujiao [2024] 49961).

References

- [1] Hatef Otroshi Shahreza, and Sébastien Marcel. Template inversion attack using synthetic face images against real face recognition systems. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 6(3):374-384, 2024. <https://doi.org/10.1109/TBIOM.2024.3391759>
- [2] Gautam Srivastava, and Surajit Bag. Modern-day marketing concepts based on face recognition and neuro-marketing: A review and future research directions. *Benchmarking: An International Journal*, 31(2):410-438, 2024. <https://doi.org/10.1108/bij-09-2022-0588>
- [3] Volodymyr Mykolaevich Opanasenko, Shavkat Khayrullaevich Fazilov, Olimjon Nomazovich Mirzaev, and Shukrullo Sa'dullo ugli Kakharov. An ensemble approach to face recognition in access control systems. *Journal of Mobile Multimedia*, 20(3):749-768, 2024. <https://doi.org/10.13052/jmm1550-4646.20310>
- [4] Thai-Viet Dang. Smart attendance system based on improved facial recognition. *Journal of Robotics and Control*, 4(1):46-53, 2023. <https://doi.org/10.18196/jrc.v4i1.16808>
- [5] Lixiong Qin, Mei Wang, Chao Deng, Ke Wang, Xi Chen, Jiani Hu, and Weihong Deng. SwinFace: A multi-task transformer for face recognition, expression recognition, age estimation and attribute estimation. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(4):2223-2234, 2023. <https://doi.org/10.1109/TCSVT.2023.3304724>
- [6] Hind Moutaz Al-Dabbas, R. A. Azeez, and Akbas Eزالdeen Ali. Two proposed models for face recognition: Achieving high accuracy and speed with artificial intelligence. *Engineering, Technology & Applied Science Research*, 14(2):13706-13713, 2024. <https://doi.org/10.48084/etasr.7002>
- [7] Wenjing Gao, Jia Yu, Rong Hao, Fanyu Kong, and Xiaodong Liu. Privacy-preserving face recognition with multi-edge assistance for intelligent security systems. *IEEE Internet of Things Journal*, 10(12):10948-10958, 2023. <https://doi.org/10.1109/JIOT.2023.3240166>
- [8] Yun Xie, Peng Li, Nadia Nedjah, Brij B. Gupta, David Taniar, and Jindan Zhang. Privacy protection framework for face recognition in edge-based internet of things. *Cluster Computing*, 26(5):3017-3035, 2023. <https://doi.org/10.1007/s10586-022-03808-8>
- [9] Huifang Feng, Wen Li, Lingfei Ma, Yiping Chen, Haiyan Guan, and Yongtao Yu. Crack-U²Net:

- Multiscale feature learning network for pavement crack detection from large-scale MLS point clouds. *IEEE Transactions on Intelligent Transportation Systems*, 25(11):17952-17964, 2024. <https://doi.org/10.1109/TITS.2024.3436015>
- [10] Pengfei Shi, Fengting Zhu, Yuanxue Xin, and Shen Shao. U²CrackNet: A deeper architecture with two-level nested U-structure for pavement crack detection. *Structural Health Monitoring*, 22(4):2910-2921, 2023. <https://doi.org/10.1177/14759217221140976>
- [11] Yunchai Li, Run Fang, Nangang Zhang, Chengsheng Liao, Xiaochang Chen, Xiaoyu Wang, Yunfei Luo, Leheng Li, Min Mao, and Yunlong Zhang. An improved algorithm for salient object detection of microscope based on U²-Net. *Medical & Biological Engineering & Computing*, 63(2):383-397, 2025. <https://doi.org/10.1007/s11517-024-03205-w>
- [12] Zujia Zheng, and Kui Yang. Wall crack detection method based on improved YOLOv5 and U²-Net. *International Journal of Wireless and Mobile Computing*, 25(4):362-367, 2023. <https://doi.org/10.1504/ijwmc.2023.135405>
- [13] Jie Chen, Yong Kong, Dawei Zhang, Yinghua Fu, and Songlin Zhuang. Two-dimensional phase unwrapping based on U²-Net in complex noise environment. *Optics Express*, 31(18):29792-29812, 2023. <https://doi.org/10.1364/OE.500139>
- [14] Huahao Fan, and Yuan Li. Image recognition and reading of single pointer meter based on deep learning. *IEEE Sensors Journal*, 24(15):25163-25174, 2024. <https://doi.org/10.1109/JSEN.2024.3416436>
- [15] Liangzhe Liao, Zhenkun Lei, Chen Tang, Ruixiang Bai, and Xiaohong Wang. Performance of a U²-Net model for phase unwrapping. *Applied Optics*, 62(34):9108-9118, 2023. <https://doi.org/10.1364/AO.504482>
- [16] Zunmei Hu, Yuwen Huang, and Yuzhen Yang. Dual-feature and multi-scale fusion using U²-Net deep learning model for ECG biometric recognition. *Journal of Intelligent & Fuzzy Systems*, 45(5):7445-7454, 2023. <https://doi.org/10.3233/JIFS-230721>
- [17] Hebba Chandravva, and Mamatha hr. Comprehensive dataset building and recognition of isolated handwritten kannada characters using machine learning models. *Artificial Intelligence and Applications*, 1(3):179-190, 2023. <https://doi.org/10.47852/bonviewAIA3202624>
- [18] Kunhua Liu, Yunqing Zhang, Yuting Xie, Leixin Li, Yutong Wang, and Long Chen. SynerFill: A synergistic RGB-D image inpainting network via fast Fourier convolutions. *IEEE Transactions on Intelligent Vehicles*, 9(1):69-78, 2023. <https://doi.org/10.1109/TIV.2023.3326236>
- [19] Siavash Jafarzadeh, Farzaneh Mousavi, Longzhen Wang, and Florin Bobaru. PeriFast/Dynamics: A MATLAB code for explicit fast convolution-based peridynamic analysis of deformation and fracture. *Journal of Peridynamics and Nonlocal Modeling*, 6(1):33-61, 2024. <https://doi.org/10.1007/s42102-023-00097-6>
- [20] Xi Jia, Joseph Bartlett, Wei Chen, Siyang Song, Tianyang Zhang, Xinxing Cheng, Wenqi Lu, Zhaowen Qiu, and Jinming Duan. Fourier-net: Fast image registration with band-limited deformation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(1):1015-1023, 2023. <https://doi.org/10.1609/aaai.v37i1.25182>
- [21] Valeriy A. Buryachenko. Fast Fourier transform method in peridynamic micromechanics of composites. *ASME 2023 International Mechanical Engineering Congress and Exposition*, 29(9):1844-1878, 2024. <https://doi.org/10.1177/10812865241236878>
- [22] Nicholas F. Marshall, Oscar Mickelin, and Amit Singer. Fast expansion into harmonics on the disk: A steerable basis with fast radial convolutions. *Methods and Algorithms for Scientific Computing*, 45(5):2431-2457, 2023. <https://doi.org/10.1137/22M1542775>
- [23] Qiaosi Yi, Faming Fang, Guixu Zhang, and Tiejong Zeng. Frequency learning via multi-scale Fourier transformer for MRI reconstruction. *IEEE Journal of Biomedical and Health Informatics*, 27(11):5506-5517, 2023. <https://doi.org/10.1109/JBHI.2023.3311189>
- [24] Sefik Ilkin Serengil, and Alper Ozpinar. "A benchmark of facial recognition pipelines and co-usability performances of modules. *Bilişim Teknolojileri Dergisi*, 17(2):95-107, 2024. <https://doi.org/10.17671/gazibtd.1399077>
- [25] Shivam Gupta, Sachin Modgil, Choong-Ki Lee, and Uthayasankar Sivarajah. The future is yesterday: Use of AI-driven facial recognition to enhance value in the travel and tourism industry. *Information Systems Frontiers*, 25(3):1179-1195, 2023. <https://doi.org/10.1007/s10796-022-10271-8>
- [26] Rosalie A. Waelen. The struggle for recognition in the age of facial recognition technology. *AI and Ethics*, 3(1):215-222, 2023. <https://doi.org/10.1007/s43681-022-00146-8>
- [27] Abbas Haider, Guanfeng Wu, Ivor Spence, and Hui Wang. Residual feature decomposition and multi-task learning-based variation-invariant face recognition. *Neural Computing and Applications*, 36(32):20147-20166, 2024. <https://doi.org/10.1007/s00521-024-10234-x>
- [28] T. Syed Akheel, V. Usha Shree, and S. Aruna Mastani. Hybrid model for face recognition using optimized linear collaborative discriminant regression classification. *Mathematical Statistician and Engineering Applications*, 71(4):10916-10924, 2022. <https://doi.org/10.1007/s00521-018-3475-4>
- [29] Lijuan Yan, Yanhu Zhang, and Yanjun Zhang. A fast face recognition system based on annealing algorithm to optimize operator parameters. *The Imaging Science Journal*, 71(3):323-330, 2023. <https://doi.org/10.1080/13682199.2023.2182261>
- [30] Hicham Zaaoui, Samir El Kaddouhi, and Mustapha Abarkan. A novel face recognition approach based

on strings of minimum values and several distance metrics. *International Journal of Computer Aided Engineering and Technology*, 18(1):60-76, 2023. <https://doi.org/10.1504/ijcaet.2023.127787>

Enhanced IoT Intrusion Detection Using an Improved Autoencoder and Adversarial Convolutional Encoders

Yukun Peng, Yu Chen*

Zhangjiakou Open University, Zhangjiakou 075000, China

E-mail: Pengyukun150@163.com; windy518@163.com

*Corresponding author

Keywords: intrusion detection, IoT, network security, attention mechanism, encoder

Received: February 14, 2025

To develop an efficient and intelligent automated intrusion detection system for IoT, this study proposes a malicious network traffic recognition model based on an improved autoencoder and adversarial convolutional encoder (AECE). The model first uses mixed sampling and improved autoencoder for data augmentation. Then, convolutional neural networks and gated recurrent units are used to extract spatial and temporal features. AECE combines the idea of generative adversarial networks to enhance the model's adaptability to complex attack patterns. Finally, experimental validation was conducted on the NSL-KDD, UNSW-NB15, IoT-23, and CSE-CIC-IDS2018 datasets. The results showed that the designed data augmentation algorithm could effectively improve the clustering and classification performance of the dataset, with a minimum Xie Beni value of 0.259, a maximum decrease of 15.88% in Davidson Boudin index, and a maximum improvement of 0.214 in classification accuracy. In the IoT-23 dataset, the highest detection rate of the baseline model was 0.882, while the detection rate of the proposed intrusion detection model was 0.949, with an increase of about 7.6%. At the same time, the model had a minimum loss convergence value of 0.08, a response time of 368.16 ms, and the values of false alarm rate fluctuated between 0.10 and 0.20. The comprehensive values of data traffic per second and packet capture per second confirmed that the model had strong detection ability and efficiency for attack behavior. This study expands the application scope of deep learning in anomaly detection, providing new ideas and methods for improving the security and stability of Internet of Things systems.

Povzetek: Predlagan je model AECE za inteligentno zaznavanje vdorov v IoT omrežjih. Uporablja izboljšani avtoenkoder za povečanje podatkov (rešuje neuravnoteženost) ter konvolucijske in ponavljajoče se enote (GRU) za ekstrakcijo prostorskih in časovnih značilnosti. Na naboru podatkov IoT-23 je AECE dosegel odlične rezultate.

1 Introduction

The Internet of Things (IoT) can realize real-time collection, analysis and interaction of various data by connecting various devices, sensors, systems, etc. to the Internet. At present, IoT has been applied in smart homes, healthcare, transportation, logistics, etc. [1]. IoT contains numerous heterogeneous devices, protocols, and platforms, with complex and diverse interactions between components. IoT devices are typically distributed across a wide geographic area, and their highly interconnected and decentralized nature makes them a hotspot for network attacks, threatening the confidentiality and security privacy of IoT data [2-3]. Therefore, establishing effective IoT security monitoring and response mechanisms to promptly detect and respond to potential security threats is crucial. Traditional security defense techniques include deploying complex security mechanisms directly on devices, conducting regular security updates, and patch management. However, IoT devices are limited in computing power, storage space, and other aspects, and their diversity and dispersion make it difficult to identify and defend against potential threats from malicious attacks [4]. Intrusion Detection Systems (IDS) can detect and

report potential security threats by monitoring and analyzing data sources like network traffic and system logs. IDS has the advantages of real-time and proactive defense, and can be used to achieve security defense for IoT devices. However, malicious cyber attacks continue to emerge and develop, with increasingly diverse attack methods and strong concealment and destructive capabilities. This makes the current IDS relatively fragile and unable to effectively respond to new security threats. Ensuring IoT security requires more advanced and efficient IDS solutions [5]. In this context, how to build an efficient and intelligent automatic detection scheme for malicious network traffic intrusion in the IoT and improve the accuracy and efficiency of malicious network traffic identification, has become a key issue that urgently needs to be addressed. Therefore, this study focuses on the Malicious Network Traffic Identification (MNTI) algorithm in IDS. It introduces feature fusion, Attention Mechanism (AM), and improved Generative Adversarial Network (GAN) to construct the MNTI model, which can fully explore and utilize the spatiotemporal correlation of network traffic data, thereby more accurately detecting known and unknown network attacks. Firstly, a Data

Augmentation Algorithm (DAA) based on Mixed Sampling (MS) and Improved Autoencoder (IA) is designed to provide a higher quality data foundation for subsequent MNTI model training. Then, a Convolutional Neural Network (CNN), Gated Recurrent Unit (GRU), and AM are combined to build an MNTI model. A GAN-based Adversarial Convolutional Encoder model (AECE) is introduced to further enhance the MNTI's adaptability to complex attack patterns. This study innovatively combines oversampling and undersampling techniques for MS, and introduces Variational Autoencoders (VAEs) for dimensionality reduction of discrete data. This can effectively solve the problem of imbalanced number of normal and abnormal samples in IoT datasets and enhance the authenticity and richness of samples.

The study is structured into four main sections. Section 1 is a review of the current research status of network IDS-related technologies in the industry. Section 2 elaborates on the construction process of DAA and MNTI models. Section 3 involves performance testing and application analysis of the designed MNTI model. Section 4 summarizes the experimental results.

2 Related works

Intrusion Detection (ItruD) technology is an important guarantee in network security, playing a crucial role in responding to network attacks and protecting systems from malicious activities. Numerous scholars have conducted research on it. Machine Learning (ML) and Deep Learning (DL) have been widely applied in information security. Qazi et al. constructed a hybrid network IDS based on DL technology. The system used CNN to collect local features and utilized deep Recurrent Neural Networks (RNN) to extract features. The public dataset has confirmed the effectiveness of this method, with a mean accuracy of 98.90% when detect malicious attacks [6]. Improving network security for cloud computing and IoT was crucial. Kasongo first utilized the eXtreme Gradient Boosting (XGBoost) Feature Selection (FS) algorithm to lower down the feature space of the data, and then built an IDS framework based on ML. The experiment confirmed the performance of the research results [7]. IDS could effectively protect the security of IoT. Hazman et al. designed an integrated learning IDS framework based on IoT intelligent environment. This framework integrated Adaptive Boosting (AdaBoost), FS technique Boruta, mutual information, and correlation. In dataset validation, this method performed well in accuracy, recall, and precision, with a Detection Rate (DR) of approximately 99.9%, a learning computation time of 33.68 seconds, and a detection time of 0.02156 seconds [8]. Ghanbarzadeh et al. designed an IDS method based on the Horse Swarm Optimization Algorithm (HSOA) and K-nearest Neighbors (KNN), which mimics the behavior of horses and selects effective features for ItruD. This method used a base function to update HSOA into a

discrete algorithm and combined it with quantum computing to implement the transformation of a quantum inspired optimizer for improving population social behavior. This method has improved the average size and classification accuracy of FS by 6%, and the accuracy of ItruD has reached 99.8% [9].

Elnakib et al. designed an enhanced anomaly-based ItruD DL multi-class classification model based on ML. This method outperformed other DL models in accuracy in classifying network traffic behavior [10]. To enhance the security of IoT, Mohy Edine M et al. constructed an FS model using principal component analysis, univariate statistical testing, and genetic algorithm, and integrated KNN to build an IoT network ItruD model. This method had high accuracy and detection time of less than one minute [11]. In response to the increased security risks of data transmission caused by interconnected nodes in IoT, Alotaibi et al. constructed a binary classification model for IoT traffic using various supervised ML models and ensemble classifiers. The classifier's accuracy surpassed that of a single model, and the predictive classification was significantly reduced [12]. The current IDS still had a high level of false positives, so Al Ghuwairi et al. developed a method for early detection of cloud computing intrusions using time series data. This method included FS and FS-based prediction models, which could effectively solve the problem of misleading connections between time series anomalies and attacks. This method significantly reduced the use of predictive factors and improved the prediction error index, reducing training time, prediction time, and cross-validation time by about 85%, 15%, and 97% [13]. The security and privacy vulnerabilities of the Internet were very urgent. Ntizikira et al. used Federated Learning (FL), Differential Privacy (DP), and secure multi-party computation to enhance data confidentiality, and integrated Deep Neural Networks (DNN) to achieve real-time anomaly detection. This method had excellent accuracy, precision, and recall [14]. Omer N et al. used Firefly Algorithm (FA) to detect intrusions before evaluating IDS, and then used Probabilistic Neural Networks (PNN) for classification. This method performed well with an accuracy rate of up to 98.99% [15]. The summary table of the above related work is shown in Table 1.

In summary, although network IDS has received a lot of research, existing IDS models generally face problems such as imbalanced data samples, fragmented spatiotemporal features, and adaptability to unknown attack patterns. In response to the above issues, this study reconstructs the data distribution, uses CNN and GRU to jointly mine spatiotemporal features, and enhances the model's ability to recognize unknown attacks using AECE. It enhances the comprehensive defense effectiveness of the model in complex attack scenarios from three dimensions: data layer, feature expression, and detection mechanism.

Table 1: Summary table of related work.

Literature	Model	Data set	Result	Limitation
[6]	CNN + RNN	CSE-CIC-IDS2018	The average accuracy rate is 98.90%	High consumption of computing resources
[7]	RNN + XGBoost	NSL-KDD	The accuracy rate is 97.8%	Insufficient generalization ability for zero day attacks
[8]	AdaBoost + Boruta	UNSW-NB15	The accuracy rate is 99.9%	Weak robustness of adversarial samples
[9]	HSOA + KNN	CIC-IDS2017	The accuracy rate is 99.8%	Parameter tuning is complex
[10]	ML	IoT-23	The accuracy of multi class classification is 98.7%	Poor model interpretability
[11]	KNN + genetic algorithm	TON_IoT	The accuracy rate is 98.3%	Significant information loss
[12]	ML + ensemble classifier	BoT-IoT	The binary classification accuracy is 99.2%	Poor scalability in multiple attack scenarios
[13]	FS	AWS CloudTrail logs	85% reduction in training time	Restricted transferability
[14]	FL + DP + DNN	CIC-IDS2019	The accuracy rate is 96.5%	Slow convergence speed
[15]	FA + PNN	KDD Cup 99	The accuracy rate is 98.99%	Insufficient coverage of modern attack modes

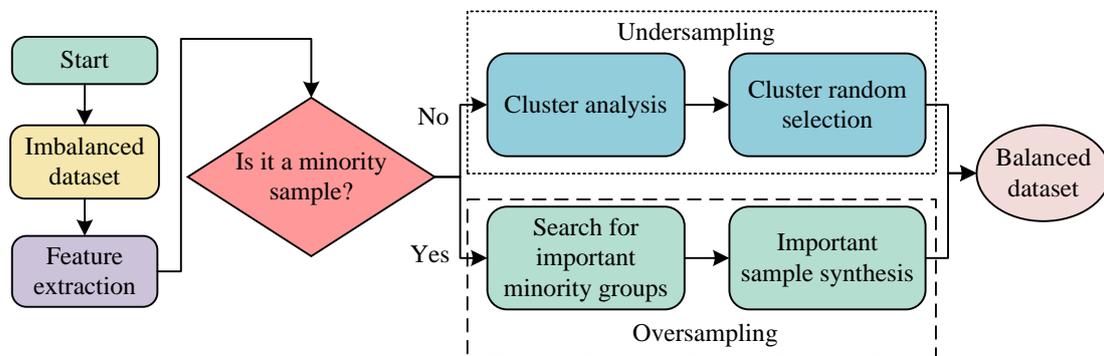


Figure 1: Schematic diagram of the workflow of DAA.

3 IoT security intrusion detection based on IA and AECE

ItruD technology can provide timely security alerts and response basis for network administrators. This study first designs DAA to improve the accuracy of traffic detection, and then integrates multiple DL technologies to construct the MNTI model.

3.1 Design of DAA based on MS and IA

With the popularity of IoT devices, the number of nodes has exploded, but normal network behavior accounts for the vast majority of traffic data, and there is a serious imbalance between the amount of abnormal samples and the normal samples. Unbalanced data samples can easily lead to a decrease in sample recognition accuracy [16]. Therefore, this study first designs DAA to address the sample imbalance, and Figure 1 shows the algorithm's framework.

In Figure 1, expansion operation is required for a few samples, while screening operation is required for most samples. Therefore, this study combines oversampling and undersampling techniques to construct a hybrid sampling system. Firstly, the category judgment threshold is determined, and the Synthetic Minority Oversampling Technique (SMOTE) is utilized to expand the imbalanced dataset and construct the training set for the classifier.

SMOTE changes the distribution of minority classes by searching for their neighbors in the feature space and generating new synthetic samples between these samples [17]. In this study, SMOTE is used to expand minority class samples, balance the class distribution in the dataset, and ensure that the model can fully learn the features of various attack categories during training, thereby improving the model's generalization ability. Then, based on ensemble thinking, multiple classifiers are used to complete ensemble training. Finally, an ensemble classifier is used to search for important minority class samples and divide them into oversampling objects. XGBoost belongs to the category of ensemble learning algorithm Boosting. This algorithm improves prediction accuracy by integrating multiple weak learners into one strong learner. In IoT datasets, normal network behavior samples often outnumber abnormal samples. XGBoost can assign higher weights to minority class samples during the training process, thereby improving the recognition ability of minority classes and effectively solving the problem of data imbalance. Therefore, this study adopts XGBoost as the basic classifier. The basic learner of XGBoost is decision tree $h(x; \theta_m)$. x is the input data. θ_m is the parameter. The weighting of all decision trees is the final prediction result. The calculation process of θ_m is equation (1).

$$\theta_m = \left\{ (R_j, c_j) \right\}_{j=1}^J \tag{1}$$

In equation (1), R_j is the leaf node region, and $J \in R$. c_j is a constant. XGBoost generates decision trees in the direction of reducing residual g_t . The calculation process of g_t is equation (2).

$$g_t = \frac{\partial L(y_i, \hat{y}_i^{t-1})}{\partial \hat{y}_i^{t-1}}, t = \{1, 2, \dots, N\} \tag{2}$$

In equation (2), y_i represents the true value of the i -th sample, \hat{y}_i^{t-1} represents the observed value of the sample at the $t-1$ -th iteration, t represents the number of iterations, and N represents the maximum number of iterations. The update process of the estimation function $F(x)$ is equation (3).

$$F_t(x) = F_{t-1}(x) + kh(x; \theta_t) \tag{3}$$

In equation (3), k is a constant. The objective function of XGBoost is the superposition of the loss function and the penalty function, as calculated in equation (4).

$$L(\phi) = \sum_i^n l(y_i - y_i) + \sum_p^P \Omega(f_p) \tag{4}$$

In equation (4), y_i and y_i are predicted values and true values, and $\sum_i^n l(y_i - y_i)$ is the loss function. $\Omega(f_p)$ is the regularization term, and the calculation process is shown in equation (5).

$$\Omega(f) = \gamma J + \frac{1}{2} \lambda' \|w\|^J = \frac{1}{2} \lambda' \sum_{j=1}^J w_j^2 \tag{5}$$

In equation (5), w is the leaf weight. γ and λ' are regular penalty terms for leaves and their weights. In ML, models may overfit training data, leading to a decrease in predictive ability on new data. Regularization reduces the risk of overfitting by adding additional penalty terms to the loss function to limit the complexity of the model. Equation (5) limits the model's complexity by comprehensively considering the number of leaf nodes and leaf weight sizes in the tree, which helps to improve the predictive performance of the model on new data. To improve the prediction accuracy of XGBoost, the training set is introduced as a new function f for greedy optimization of the objective function, as expressed in equation (6).

$$L^{(t)} = \sum_i^n \left(l\left(y_i^{(t-1)} - y_i\right) + f(x_i) \right) + \Omega(f_p) \tag{6}$$

After expanding equation (6) according to the second-order Taylor formula, the final objective function $L^{(t)}$ is obtained through training simplification, as shown in equation (7).

$$L^{(t)} = \sum_{j=1}^J \left[\left(\sum_{i \in I} G_i \right) w_j + \frac{1}{2} \left(\sum_{i \in I} H_i + \lambda' \right) w_j^2 \right] + \gamma J \tag{7}$$

In equation (7), G and H are the first and second derivatives of the loss function. IoT datasets typically contain a large number of numerical features, and some density or hierarchical clustering algorithms may have issues with not being intuitive or having high computational costs when processing numerical data. The K-means algorithm has a simple principle and good clustering effect on numerical data. It can quickly divide the data into different clusters and select representative samples, thereby improving training efficiency. Therefore, this study adopts the K-means clustering algorithm for undersampling operation, and the objective function is shown in equation (8).

$$J(X, \pi) = \sum_{j=1}^k \sum_{i \in \pi_j} \|x_i - m_j\|^2 \tag{8}$$

In equation (8), π_j is class j . m_j is the center of a certain category. x_i is a data point. The MS method compensates for the shortcomings of traditional sampling techniques, but IoT datasets typically involve data with discrete characteristics. The SMOTE algorithm has low applicability to discrete data. VAE can map input data to latent space through an encoder, obtain representation vectors, output parameters of the representation vectors, and generate diverse new samples. This will increase the richness of the dataset and help improve the model's generalization ability. VAE has good processing ability for discrete data. Therefore, the study introduces VAE for dimensionality reduction of discrete data. VAE contains an encoder and a decoder. The encoder maps the input data x to the latent space to gain the representation vector z , and outputs the parameters of the representation vector. The decoder maps the representation vector back to the data space to generate new samples and ensures that the new samples are as similar as possible to the original input data [18-19]. The training objective of VAE is to optimize the variational lower bound *ELBO*, as shown in equation (9).

$$ELBO(q) = Eq(z|x) [\log p(x|z)] - DKL(q(z|x) \square p(z)) \tag{9}$$

In equation (9), $p(x|z)$ is the generative model defined by the decoder. $p(z)$ is a standard Gaussian distribution. DKL is the Kullback Leibler divergence. $q(z|x)$ is the posterior distribution. The working principle of DAA based on MS/IA is shown in Figure 2.

In Figure 2, data augmentation is divided into two stages: model training and data synthesis. Firstly, VAE is used to learn data features during the training phase and convert them into representation vectors with rich information. Then, the representation vector and data labels are input into the MS module to achieve balanced processing of the data. Finally, the decoder completes the conversion of the data format. In summary, the proposed

DAA based on MS and IA mainly consists of four steps. Step 1 inputs network traffic data and preprocesses the raw data. Step 2 determines the majority class and minority class samples, applies SMOTE to generate new composite samples for the minority class samples, and uses K-means clustering algorithm to undersample the majority class samples. Step 3 trains the VAE using the training set data and uses the trained VAE to perform dimensionality reduction and feature extraction on minority class samples. Step 4 fuses the synthesized samples generated by SMOTE and VAE to obtain an enhanced dataset, and performs weighted fusion with the majority class samples to obtain a balanced dataset. This study uses Xie Beni Index (XBI) and Davies Bouldin Index (DBI) as indicators to evaluate the clustering quality of DAAs. XBI evaluates clustering performance by measuring the distance between cluster centers and the closeness of data points within clusters. The smaller the value of XBI, the more tightly clustered the sample points within the cluster are, and the better the separation between different clusters, resulting in better clustering performance. DBI takes into account both intra-class sample similarity and inter-class sample difference, with smaller values indicating better clustering performance.

3.2 Design of MNTI model based on feature fusion and AECE

IDS is usually segmented into two types of signature detection and two main technologies. Anomaly detection is a detection technique that identifies abnormal activity by analyzing the normal behavior patterns of network traffic. When network traffic deviates from normal behavior patterns, the system considers it a potential malicious activity and triggers an alert [20]. Network traffic data typically contain a mixture of multiple types of information, which are correlated in both temporal and spatial dimensions. Network traffic data have obvious temporal characteristics. For example, network attack behavior usually shows sudden growth of traffic in a short period. At the same time, network traffic data also have spatial correlations. In the same IoT network, data transmission between servers and multiple clients may exhibit synchronous or correlated trends, and devices in the same network often share certain common network configurations and security policies. Therefore, to capture information at different levels, this study constructs the basic framework of the MNTI model based on the concept of feature fusion, as shown in Figure 3.

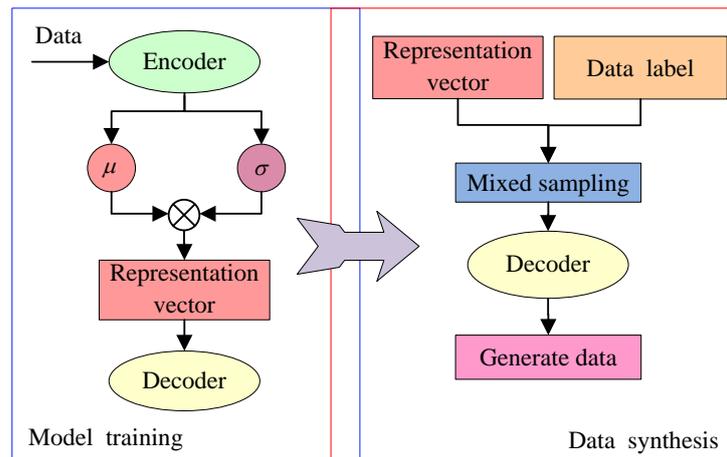


Figure 2: Schematic diagram of DAA based on MS/IA.

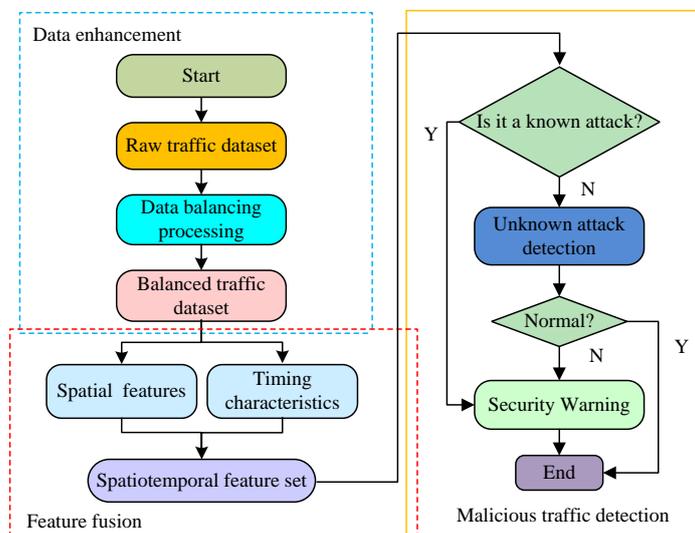


Figure 3: Basic framework structure of MNTI model based on feature fusion.

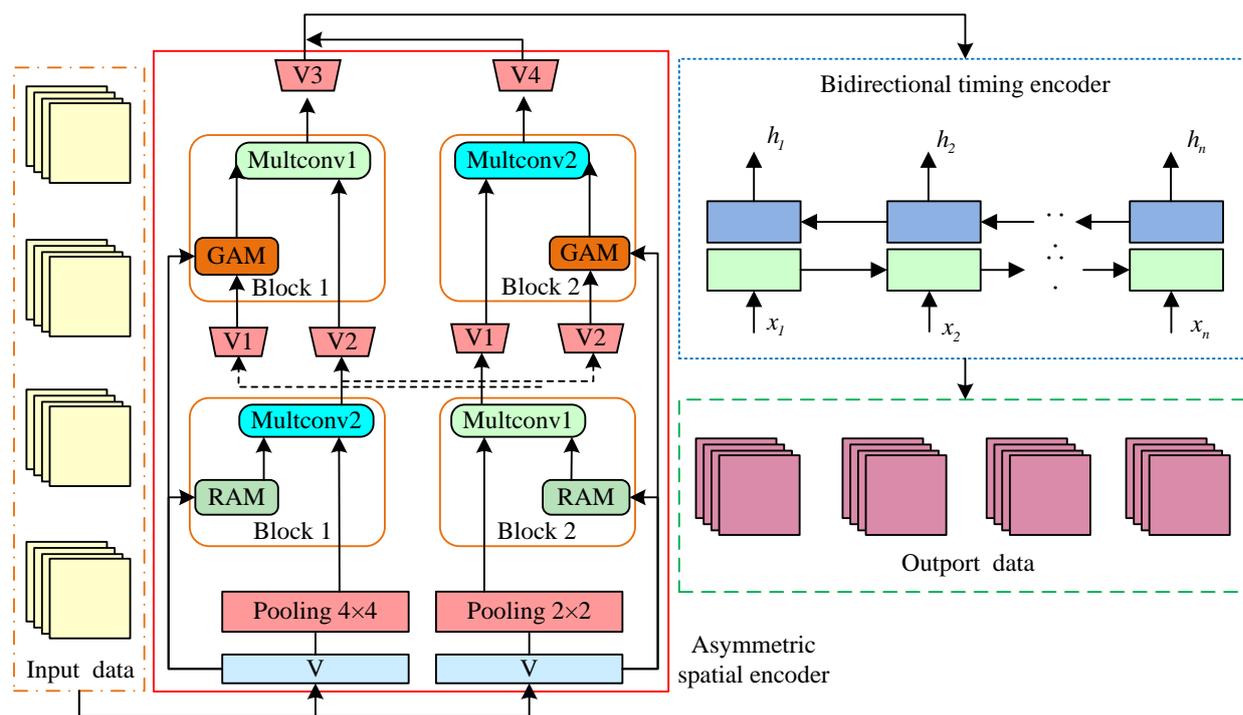


Figure 4: Schematic diagram of feature extraction encoder structure.

As shown in Figure 3, after data augmentation, the MNTI model mainly includes two modules: feature fusion and malicious traffic detection. Among them, the feature fusion module extracts spatial and temporal features of the balanced dataset, and disease fusion presents a spatiotemporal feature set. The malicious traffic detection module first determines whether it is a known attack. If it is a known attack, it directly initiates security warning measures. If not, further unknown attack detection will be conducted. If the traffic detection is abnormal, timely security warning measures should be taken and the network attack database should be updated. In the context of network traffic, spatial features reflect the combination relationship between different features in network traffic data, such as a combination pattern of features such as different IP addresses and different ports. Spatial features can also represent communication relationships between different devices or nodes. The temporal characteristics mainly describe the dynamic changes and patterns of network traffic data in the time dimension, reflecting the changing trends and periodic patterns of network traffic data in time, such as the peak and off peak periods of network traffic, periodic fluctuations in traffic, etc. This study uses CNN structure for spatial feature extraction and RNN suitable for sequence data processing for temporal feature extraction. The structure of the feature extraction encoder designed for the feature fusion module is shown in Figure 4.

In Figure 4, the feature extraction spatial encoder structure consists of an Asymmetric Spatial Encoder (ASE) and a Bidirectional Temporal Encoder (BTE). ASE is used to extract spatial features from raw data. BTE is used to extract temporal features from the extracted spatial features, achieving the effect of fusing features from

different dimensions. Finally, spatial and temporal features are fused to form a comprehensive feature representation. ASE is based on traditional CNN architecture, consisting of four blocks that integrate two different types of AMs and convolutional kernels of different scales. Four blocks use two types of multi-scale convolutional layers. Both Multichannel 1 and Multichannel 2 contain three convolutional path calculations and use three various sizes of convolution kernels, namely 3×3 , 5×5 , and 7×7 , to enhance the receptive field of the network. In addition, Block also introduces Global Attention Mechanism (GAM) and Residual Attention Mechanism (RAM). Firstly, RAM is used to fuse multi-scale inputs with the original image, and residual connections can be introduced to enhance the model's generalization ability. Then GAM is used to fuse the output of RAM with the original image. GAM can correlate and weight all positions in the input sequence, enhancing the model's overall understanding and processing ability of the input sequence. The selected basic RNN unit is GRU. GRU refers to a variant structure of RNN that can reduce gradient vanishing while preserving long-term sequence information. The BTE structure is shown in Figure 5.

In Figure 5, the BTE structure has undergone bidirectional improvement on the basis of traditional RNN and introduced multi head self AM. Bidirectional GRU (BiGRU) can extract forward and backward data and determine whether there is abnormal information in the current traffic data [21]. The merging strategy is used to fuse the forward and backward hidden states of BiGRU to generate the final sequence representation. To provide more comprehensive sequence feature information and improve the detection performance of the model, this

study adopts a concatenation strategy, directly concatenating the forward and backward hidden states into a vector. The update process of the forward update gate z_t and reset gate r_t in BiGRU is equation (10).

$$\begin{cases} r_t = \sigma(w_n x_t + u_n h_{t-1}) \\ z_t = \sigma(w_m x_t + u_m h_{t-1}) \end{cases} \quad (10)$$

In equation (10), σ represents the Sigmoid activation function, with an output value between 0 and 1. The larger the value, the more information from the previous time step is retained. w_n and w_m represent the weight parameters of the update gate and reset gate, respectively. u_n and u_m represent weight matrices. h_{t-1} represents the previous state. h is the hidden layer state. x_t is the input information at the current time. The calculation of output layer h_t is equation (11).

$$h_t = (1 - z_t)h_{t-1} + z_t h_t \quad (11)$$

In equation (11), h_t represents the updated value of the reset gate. The reverse calculation formula for BiGRU is equation (12).

$$\begin{cases} z_t^a = \sigma(w_m^a x_t + u_m^a h_{t+1}) \\ r_t^a = \sigma(w_n^a x_t + u_n^a h_{t+1}) \end{cases} \quad (12)$$

In equation (12), a is the reverse GRU. w_m^a and w_n^a represent weight parameters. u_m^a and u_n^a represent weight matrices. h_{t+1} represents the state at the next moment. Finally, the hidden layer states of the forward

and reverse GRUs are weighted and summed to obtain the final prediction result, as shown in equation (13).

$$y_t = \sigma(h_t \times w_y) \quad (13)$$

In equation (13), w_y is the weight between the hidden and output layers. Finally, the predicted temporal results are input into the multi head self-AM to achieve weighted summation of encoding. The detection objects of the traffic detection module include known and unknown network attacks. The known detector for network attacks is SoftMax. The SoftMax expression is equation (14).

$$SoftMax(x) = \exp(x) / \text{sum}(\exp(x)) \quad (14)$$

In equation (14), \exp is an exponential function. In the feature fusion module, this study achieves the extraction of spatial and temporal features through CNN and GRU, while reducing the impact of redundant features. The convolutional layer automatically filters local features through convolutional kernels of different scales, while the gating mechanism of GRU filters out irrelevant temporal information. In addition, to further improve the model's performance, this study also ranks the importance of features. In the data augmentation stage, XGBoost is used to rank features and select the top-ranked features for subsequent model training. Meanwhile, in the feature extraction encoder, Genetic Algorithm (GA) and RA are introduced to automatically focus on the more important features for ItruD by learning the weights of features, thus achieving feature importance ranking. The detection model for unknown network attacks is shown in Figure 6.

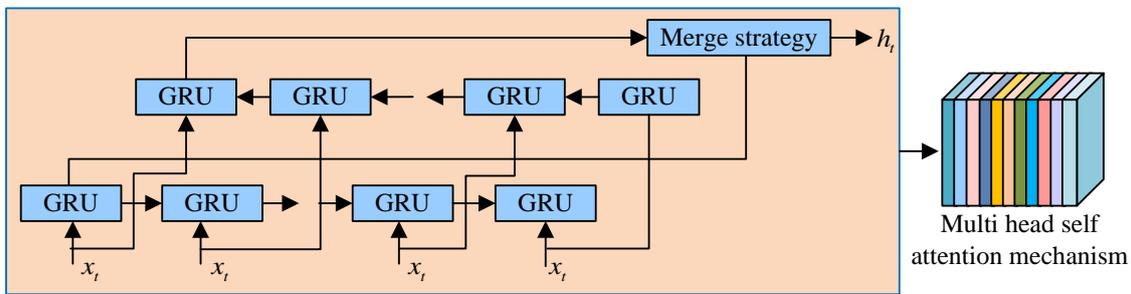


Figure 5: Schematic diagram of BBTE structure.

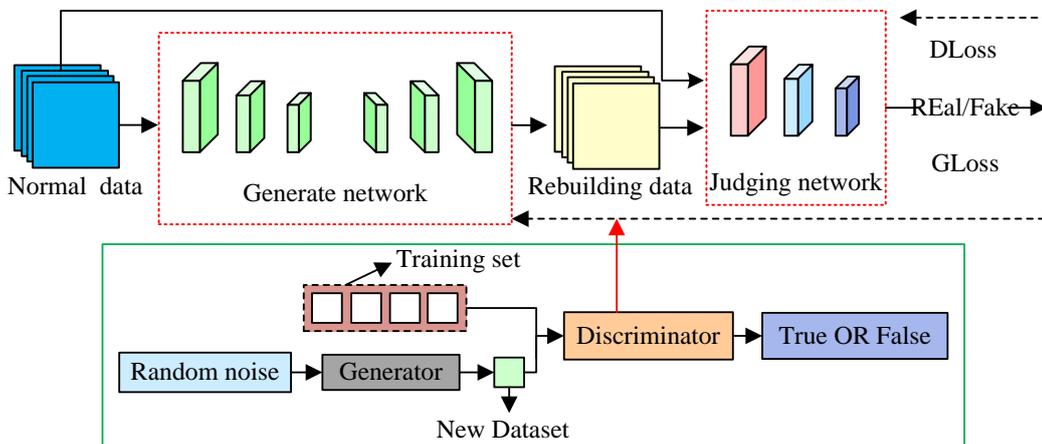


Figure 6: Unknown network attack detection model structure.

In Figure 6, the model is designed based on the concept of GAN and consists of two parts: the generative network and the judgment network. Generate models that produce fake data similar to real samples. The discriminative model is responsible for distinguishing and judging between real data and generated data [22]. In the early stages of training, the weights of the discriminator are randomly initialized. As training progresses, the GAN continuously learns how to generate more realistic data, while the discriminator also updates its parameters based on the feature differences between real and fake data. At the end of the training phase, the generator and discriminator reach Nash equilibrium, and the discriminator's loss tends to stabilize. The anomaly detection threshold is based on a dynamic adjustment strategy. In practical applications, if the False Alarm Rate (FAR) is too high, the threshold should be appropriately increased to reduce misjudgments of normal behavior. If the false alarm rate is too high, the threshold can be appropriately lowered to improve the detection ability of attack behavior. The training game process of GAN is equation (15).

$$\begin{aligned} \min_G \max_D V(D, G) = & \\ & E_{x \sim P_{data}(x)} [\log(D(x))] \\ & + E_{z \sim P_{model}(z)} [\log(1 - D(G(z)))] \end{aligned} \quad (15)$$

In equation (15), z represents noise. x is the real sample data. $P_{data}(x)$ is the probability distribution function of x . $P_{model}(x^{(i)}; \theta)$ is the probability distribution function for judging the network, and θ is the parameter. G and D are generative networks and discriminative networks. The detection model is defined as AECE. The generative network part includes encoders and decoders. The encoder and decoder both contain 3 convolutional layers and two pooling layers. The training process needs to make the reconstructed data of the decoder closest to the original data, and use the maximum reconstruction loss of normal traffic behavior as the threshold for detecting unknown attacks. The training process needs to maximize the probability of generating samples as real samples, consisting of convolutional, pooling, and fully connected layers. This study uses Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Mean Absolute Percentage Error (MAPE) as evaluation metrics. MAE, RMSE, and MAPE are commonly used indicators to evaluate the difference between predicted and true values in regression models. In ItruD, they can be used to measure the accuracy of predicting network traffic characteristics, indirectly reflecting the model's ability to distinguish between normal and abnormal traffic. MAE represents the average absolute error between predicted values and true values. RMSE emphasizes the impact of larger errors. MAPE

displays model accuracy in the form of relative errors. The calculation of MAE, RMSE, and MAPE is shown in equation (16).

$$\begin{cases} MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \\ RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \\ MAPE = MAE = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100\% \end{cases} \quad (16)$$

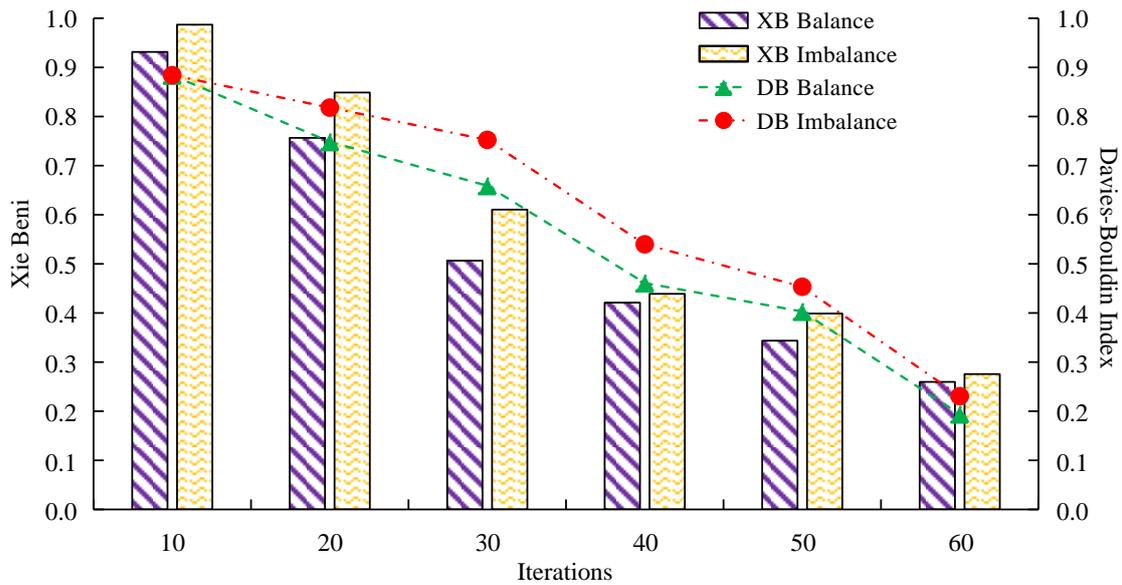
In equation (16), n represents the number of samples, y_i represents the true value, and \hat{y}_i represents the predicted value.

4 Performance testing and application effect analysis of IoT security ItruD model

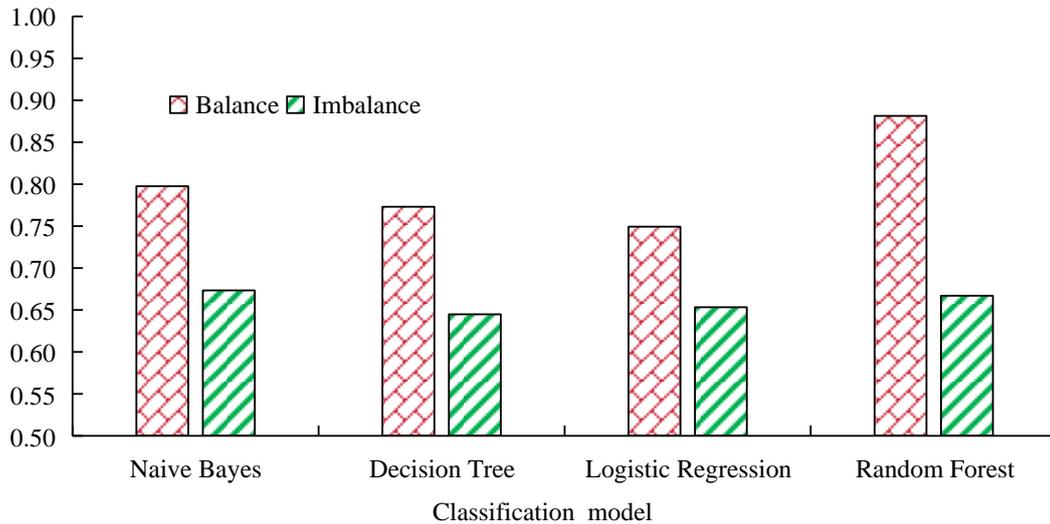
To verify the effectiveness of the designed DAA and MNTI models, this study conducts performance testing and application effect analysis, and discusses the results.

4.1 Performance testing of IoT intrusion detection model

The experiment is conducted using the CentOS 7 operating system and the DL framework is Pytorch 1.7. The central processing unit is Intel (R) Xeon (R) Silver 4214 2.20 GHz, with 128 GB of memory. The image processor is Ge Force RTX 2080Ti. The programming language is Python 3.8. The experiment selects Non-Intrinsic FS for KDD (NSL-KDD), UNSW-NB15, IoT-23, and CSE-CIC-IDS2018 datasets for performance testing. NSL-KDD includes normal connections and various types of attacks, covering multiple characteristics such as connection duration, source/destination ports, service type, protocol type, etc. UNSW-NB15 simulates network traffic in a real network environment, containing 175,341 network connection records, covering common network attacks and normal traffic. IoT-23 contains a large amount of device interaction data, sensor readings, and network communication records. CSE-CIC-IDS2018 contains network traffic data captured from multiple real network environments, covering various types of attacks and normal traffic. The eigenvalues are scaled to the range of 0-1 and divided into training, testing, and validation sets in an 8:1:1 ratio to standardize the data. The learning rate is set to 0.001, Epoch is 60, Batchsize is 32, hidden layer is 2, and Adam optimizer is used. Firstly, the performance of DAA is analyzed, and the clustering and comparison effects before and after data balancing are compared, as shown in Figure 7.



(a) Comparison of clustering effects



(b) Comparison of classification effects

Figure 7: Analysis of the effect of data enhancement algorithm.

Table 2: Results of ablation experiment.

Models	Detection rate	Precision	Recall	MAE
Without feature fusion module	0.856	0.865	0.846	0.214
Without DAA	0.824	0.834	0.813	0.245
Without AECE	0.879	0.887	0.871	0.198
Complete model	0.919	0.925	0.912	0.179

In Figure 7 (a), there is a significant difference in the clustering performance evaluation indicators of the dataset before and after data balancing. The XBI and the DBI both achieve better results on the balanced dataset, with a minimum XBI of 0.259 and a minimum DBI of 0.194, with a decrease of 5.78% and 15.88%, respectively. After DAA processing, the intra cluster compactness and inter class separation of the dataset are improved, and the clustering effect is improved. In Figure 7 (b), four different baseline classification models achieve better classification accuracy on the balanced dataset after data

augmentation, with a maximum accuracy improvement of 0.214. To demonstrate the contribution of each component of the model to overall performance, ablation experiments are designed and studied. The ablation experiment uses the NSL-KDD dataset to compare the DR and error metrics of the complete model with models without feature fusion modules, DAA, and AECE. The results of the ablation experiment are shown in Table 2.

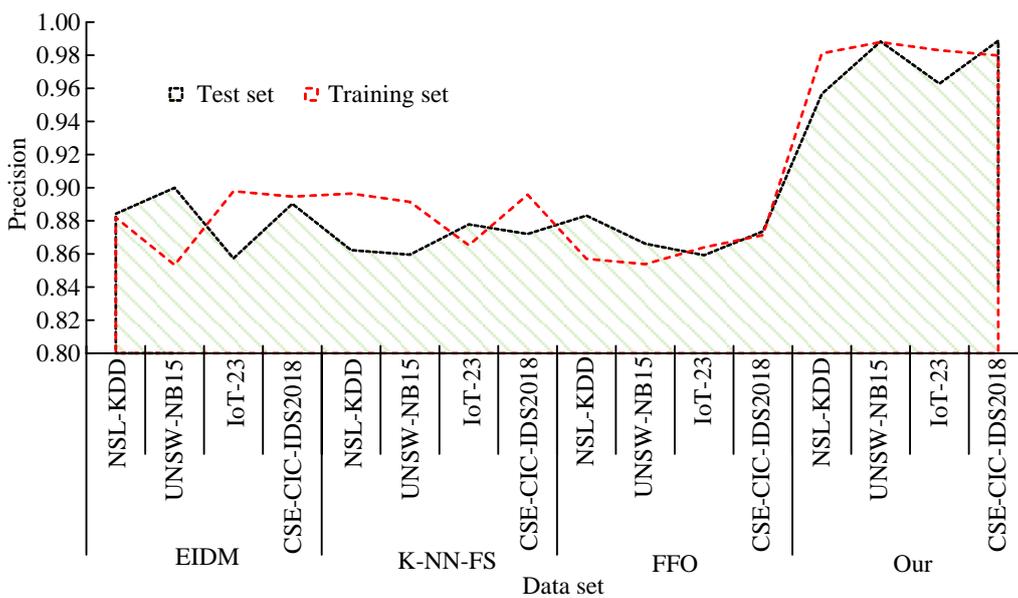
From Table 2, the DR, precision, and recall rate of the complete model are the highest, while the MAE is the lowest, indicating that the proposed improvement

strategies can effectively improve the ItruD performance. Among them, the model without DAA performs the worst in terms of metrics, indicating that DAA contributes the most to model performance and can significantly improve the model's ability to identify attack samples by solving the problem of data imbalance. The MNTI model is compared with the Enhanced Anomaly-based ItruD DL Multi-class Classification (EIDM) proposed in reference [10], the KNN classifier and FS-based ItruD model (K-NN-FS) in reference [11], and the Firefly Optimization (FFO) detection model in reference [15]. Wilcoxon signed rank test is used to evaluate the performance difference between the proposed model and the baseline model, with a $p < 0.05$ indicating statistical significance of the difference. To ensure the reliability and stability of the results, each model is independently run 5 times on each dataset. The performance indicators reported are the average of these 5 runs, presented in the form of mean \pm standard deviation. The classification performance of different ItruD models is shown in Table 3.

In Table 3, the performance of the proposed model on all four datasets is significantly better than the other three baseline models ($p < 0.001$). The research model has the smallest value in the ItruD classification error index, with the minimum values of MAE, RMSE, and MAPE being 0.179, 0.236, and 0.197. The model detection errors of the other three literature are all greater than 0.3. This means that the designed model has the smallest classification error in ItruD and accurately distinguishes traffic between attack behavior and normal behavior. In addition, the DR of the proposed model is the highest, reaching 0.949. The maximum DR values for EIDM, K-NN-FS, and FFO models are 0.885, 0.882, and 0.853. High detection precision means that the model can effectively identify malicious traffic from a large amount of network traffic data, which is crucial for timely detection and response to network attacks. The F1 index is the harmonic mean of precision and recall, used to comprehensively evaluate the performance of a model. Figure 8 compares the scalability of different models.

Table 3: Classification performance of diverse ItruD models.

Model	Index	NSL-KDD	UNSW-NB15	IoT-23	CSE-CIC-IDS2018	p -value (vs research model)
Research model	MAE	0.179±0.015	0.198±0.018	0.269±0.022	0.199±0.017	-
	RMSE	0.236±0.020	0.273±0.024	0.286±0.023	0.284±0.026	-
	MAPE	0.199±0.019	0.197±0.019	0.200±0.021	0.261±0.023	-
	DR	0.919±0.013	0.921±0.011	0.949±0.008	0.900±0.015	-
Reference [10]	MAE	0.337±0.033	0.424±0.038	0.325±0.032	0.304±0.029	<0.001
	RMSE	0.430±0.039	0.394±0.037	0.400±0.037	0.392±0.036	<0.001
	MAPE	0.338±0.035	0.442±0.042	0.349±0.033	0.404±0.033	<0.001
	DR	0.729±0.026	0.885±0.018	0.796±0.023	0.823±0.017	<0.001
Reference [11]	MAE	0.419±0.042	0.419±0.040	0.426±0.042	0.325±0.035	<0.001
	RMSE	0.338±0.033	0.437±0.041	0.362±0.035	0.319±0.031	<0.001
	MAPE	0.465±0.045	0.497±0.048	0.359±0.034	0.450±0.043	<0.001
	DR	0.827±0.020	0.842±0.021	0.882±0.018	0.820±0.023	<0.001
Reference [15]	MAE	0.345±0.034	0.490±0.045	0.339±0.036	0.461±0.042	<0.001
	RMSE	0.478±0.047	0.411±0.038	0.416±0.041	0.451±0.043	<0.001
	MAPE	0.351±0.036	0.400±0.039	0.447±0.042	0.379±0.036	<0.001
	DR	0.841±0.023	0.853±0.022	0.828±0.020	0.833±0.022	<0.001



(a) Comparison of precision

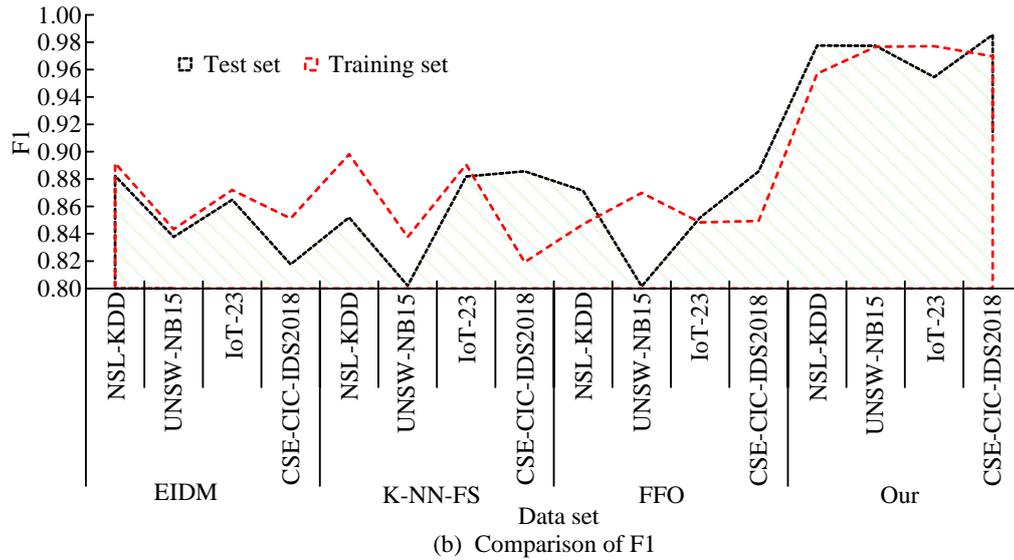


Figure 8: Scalability comparison of different ItruD models.

Table 4: The training time of the model on different datasets (s).

Model	NSL-KDD	UNSW-NB15	IoT-23	CSE-CIC-IDS2018
EIDM	158.25	183.49	204.96	198.42
K-NN-FS	92.33	105.56	120.71	112.98
FFO	143.75	165.42	192.04	178.64
Research model	182.43	210.76	244.67	226.28

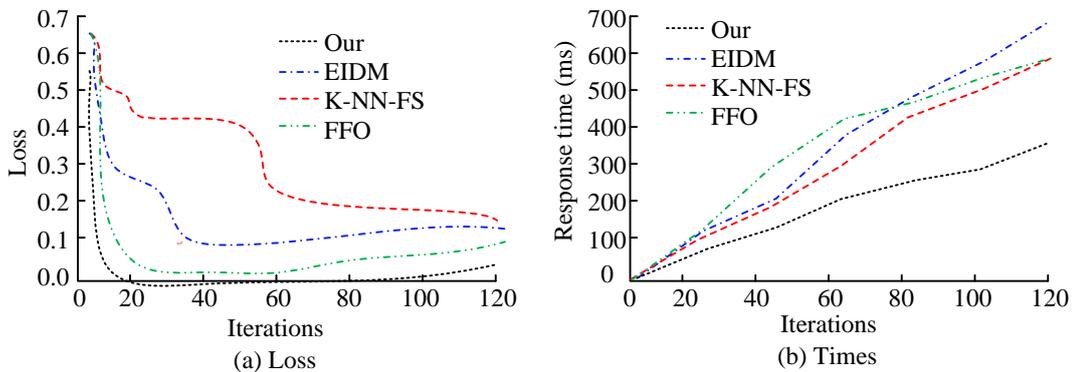


Figure 9: Comparison of loss function curves and response time for different ItruD models.

In Figure 8, the research model has significant advantages in detection accuracy and F1 index values, and performs well on four different datasets. The maximum precision values on each dataset are 0.981, 0.988, 0.983, and 0.989. The maximum values of F1 index are 0.978, 0.977, 0.977, and 0.985. The difference in values between the test and training sets of the research model is small, and the data fluctuation is not significant. The results indicate that the proposed model can maintain high detection precision and F1 index on different datasets, and has good generalization ability, balance, and stability. This is mainly due to the introduction of data augmentation, feature fusion and extraction, and adversarial training techniques in the model, which significantly improve the performance and scalability of IDS. The training time of the above model on different datasets is shown in Table 4.

From Table 4, compared to the comparison model, the proposed model has a longer training time on all four datasets, with the longest being 244.67 seconds. This is because the architecture of the proposed model is more complex, including feature fusion, IA, AECE, and other components, which increases the complexity of the model and leads to an increase in training time. The baseline model architecture is relatively simple, so the training time is relatively short.

4.2 Performance testing and application effect analysis of IoT ItruD model

It continues to compare the application effects of different ItruD models in practice. The loss function curve and response time of the model are shown in Figure 9.

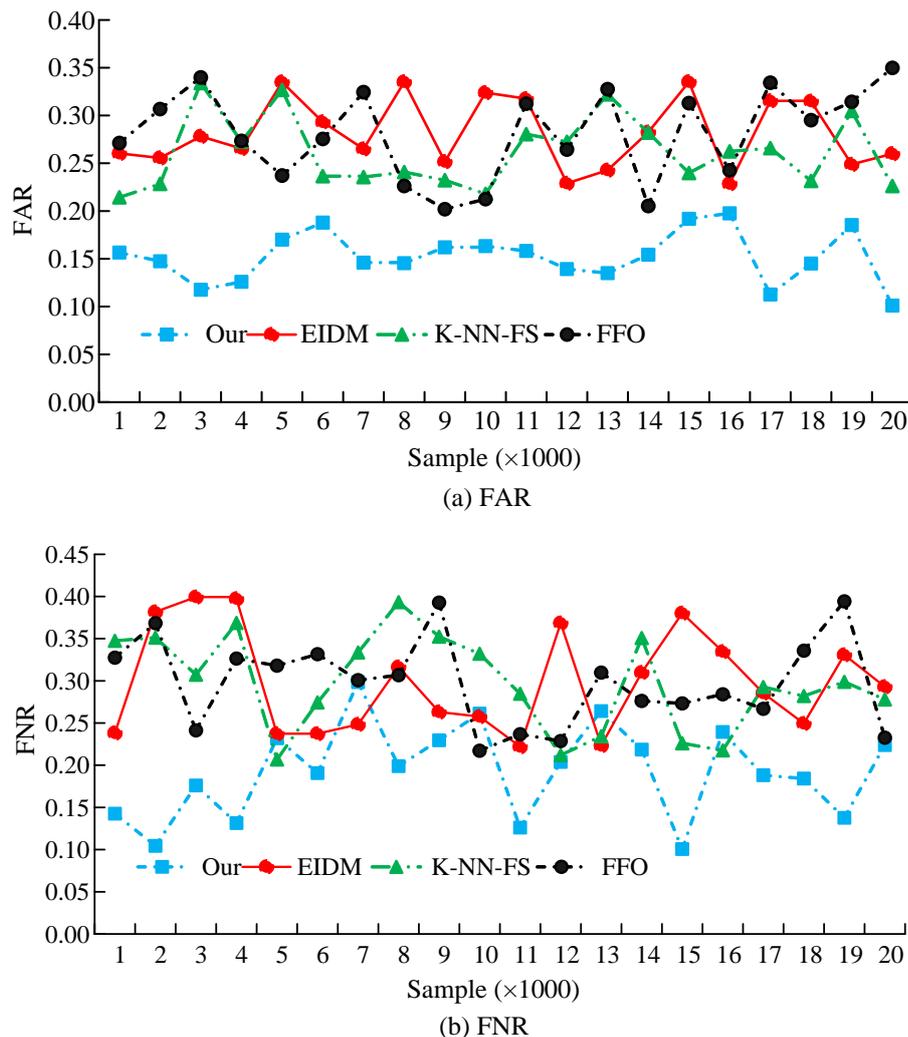


Figure 10: FNR and FAR of various models.

In Figure 9 (a), the research model has the fastest convergence speed on the loss function curve, can converge early in the iteration, and has a minimum convergence value of 0.08, which has a significant convergence advantage over other models. The fast convergence loss function curve indicates that the model can learn features and patterns in the data faster, and does not overfit the training data during the training process, but learns the general features of the data well. In addition, rapid convergence also indicates that the optimization process of the model is more efficient and can achieve the expected performance level in fewer iterations. In Figure 9 (b), the response time of the research model is 368.16ms. The response times of EIDM, K-NN-FS, and FFO models are 684.1 ms, 589.3 ms, and 598.4 ms. A shorter response time means that the model can detect and respond to network traffic faster in practical applications, which is crucial for IoT security IDSs with high real-time requirements. The False Negative Rate (FNR) and FAR of different models in application are displayed in Figure 10.

In Figure 10 (a), the FAR values of the proposed model fluctuate in the range of 0.10-0.20 under different

sample sizes. The FAR of other models fluctuates within the range of 0.20-0.35. FAR reflects the tendency of the model to misjudge normal traffic as attack traffic. The proposed model has good recognition performance for normal behavior, with fewer false alarms. In Figure 10 (b), the research model achieves excellent FNR performance, with values fluctuating between 0.10-0.20. The proportion of actual attack samples that can be detected is relatively high compared to all actual attack samples. The results of data traffic per second and packet capture per second for different models are shown in Figure 11.

Figures 11 (a) and (b) show that the research model has the highest values in both data traffic per second and packet capture per second. Overall, the model is capable of processing a large number of data packets per second and has strong packet processing capabilities, reflecting the strong detection ability and efficiency of the research design for attack behavior. Based on Figure 10, this method has a low rate of missed attacks. In Figure 11, there is no packet loss phenomenon for all methods.

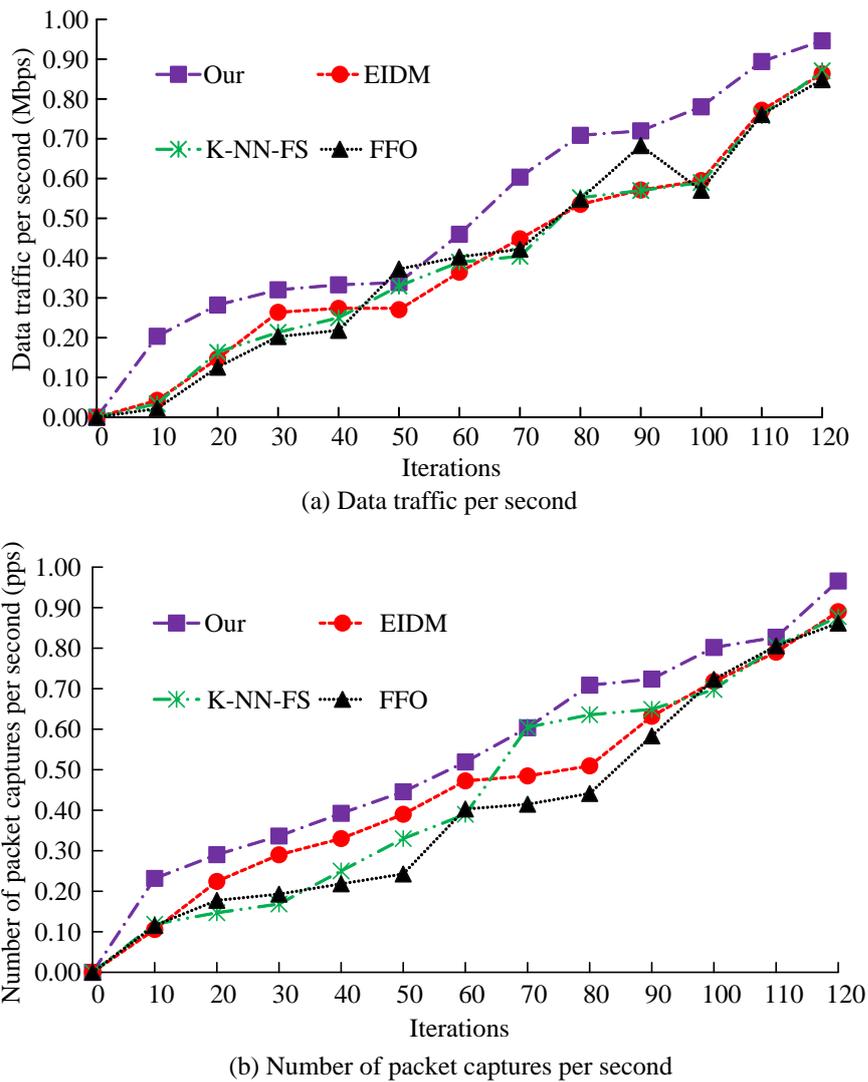


Figure 11: Comparison of data traffic per second and packet capture per second for different models.

5 Discussion

To cope with malicious attacks on IoT devices, this study conducted data augmentation based on hybrid sampling and auto-encoder, and constructed an MNTI model using feature fusion on this basis. The experiment showed that after balancing the DAA dataset, the minimum XBI value was 0.259, the minimum DBI value was 0.194, and the decrease was 5.78% and 15.88%, respectively. The classification accuracy of different classification base models has been improved. The minimum values of MAE, RMSE, and MAPE for the research model were 0.179, 0.236, and 0.197, and the maximum DR value was 0.949. The maximum accuracy of this model on four datasets was 0.981, 0.988, 0.983, 0.989, and the maximum F1 index was 0.978, 0.977, 0.977, 0.985. In the application process, the convergence speed on the loss function curve was the fastest, the convergence value was the smallest, and the response time was 368.16 ms. In addition, compared with the baseline models EIDM, K-NN-FS, and FFO, the proposed MNTI model also showed significant advantages in false positives and false negatives. The FAR

and FNR values of the proposed MNTI model fluctuated within the range of 0.10-0.20, which could more accurately distinguish between normal and abnormal traffic, thereby reducing the false positive rate. In contrast, the FAR values of the baseline model fluctuated within the range of 0.20-0.35, indicating a relatively high false positive rate.

The DR of EIDM proposed in reference [10] on the NSL-KDD dataset was 0.729, while the MNTI model proposed in the study reached 0.919. On the UNSW-NB15 dataset, the DR of EIDM was 0.885, while the proposed MNTI model was 0.921. The DR of the proposed MNTI model was superior to that of the EIDM model on various datasets. Similarly, the DRs of the ItruD models proposed in references [11] and [15] were also lower than those of the proposed MNTI model. This was mainly attributed to the integration of various advanced DL techniques and ideas in this study, including feature fusion, AMs, and improved GANs. Combining CNN and GRU to extract spatiotemporal features and introducing VAE for dimensionality reduction and feature extraction of data can effectively capture spatial correlations and local

features, and better process sequence data. It can also enrich feature information, enabling the model to more accurately capture key features in network traffic data. The AM can automatically learn the importance of different features, making the model more focused on key features related to ItruD, thereby improving the model's discriminative ability. In addition, the AECE introduces the idea of GAN and utilizes adversarial training between the generative network and the judgment network to further enhance the model's ability to detect unknown attacks.

In practical applications, the proposed MNTI model demonstrates good scalability through its flexible design and modular structure. The feature extraction module can be adjusted according to the type of input data, such as replacing CNN with a network structure suitable for processing specific data types, or adding new feature extraction components to adapt to new data sources. The feature fusion mechanism can effectively integrate feature information from different modules, thereby enhancing the model's ability to process multi-source data. In addition, the depth and breadth of the model can be expanded according to actual needs to further improve its performance and application scope. For example, increasing the number of network layers to capture more complex feature patterns, or adopting multi task learning strategies to simultaneously process multiple related tasks.

6 Conclusion

This study aims to improve the accuracy of identifying malicious network traffic in the IoT environment to cope with malicious attacks on IoT devices. By using MS and VAE for data augmentation, the problem of data imbalance has been effectively solved, providing a high-quality data foundation for model training. On this basis, multiple technologies such as CNN, RNN, AM, and GAN are integrated to construct the MNTI model, which can comprehensively capture the characteristics of network traffic data. Experimental studies have shown that the proposed model has good detection performance and stability, can accurately distinguish between attack behavior and normal behavior of traffic, and has high security protection efficiency and real-time performance. However, the computational complexity of the proposed model is relatively high, and deployment on resource constrained IoT devices may pose certain difficulties. Therefore, in future research, the model structure should be further optimized by using techniques such as model compression and quantization to reduce the computational complexity of the model, making it more suitable for resource constrained IoT environments.

References

- [1] Arash Heidari, and Mohammad Ali Jabraeil Jamali. Internet of Things intrusion detection systems: A comprehensive review and future directions. *Cluster Computing*, 26(6):3753-3780, 2023. <https://doi.org/10.1007/s10586-022-03776-z>
- [2] Oluwadamilare Harazeem Abdulganiyu, Taha Ait Tchakoucht, and Yakub Kayode Saheed. A systematic literature review for network intrusion detection system (IDS). *International Journal of Information Security*, 22(5):1125-1162, 2023. <https://doi.org/10.1007/s10207-023-00682-2>
- [3] Sampath Rajapaksha, Harsha Kalutarage, M. Omar Al-Kadri, Andrei Petrovski, Garikayi Madzudzo, and Madeline Cheah. Ai-based intrusion detection systems for in-vehicle networks: A survey. *ACM Computing Surveys*, 55(11):1-40, 2023. <https://doi.org/10.1145/3570954>
- [4] Ankit Thakkar, and Ritika Lohiya. A review on challenges and future research directions for machine learning-based intrusion detection system. *Archives of Computational Methods in Engineering*, 30(7):4245-4269, 2023. <https://doi.org/10.1007/s11831-023-09943-8>
- [5] Noor Aldeen Alawad, Bilal H. Abed-alguni, Mohammed Azmi Al-Betar, and Ameera Jaradat. Binary improved white shark algorithm for intrusion detection systems. *Neural Computing and Applications*, 35(26):19427-19451, 2023. <https://doi.org/10.1007/s00521-023-08772-x>
- [6] Emad Ul Haq Qazi, Muhammad Hamza Faheem, and Tanveer Zia. HDLNIDS: Hybrid deep-learning-based network intrusion detection system. *Applied Sciences*, 13(8):4921-4936, 2023. <https://doi.org/10.3390/app13084921>
- [7] Sydney Mambwe Kasongo. A deep learning technique for intrusion detection system using a Recurrent Neural Networks based framework. *Computer Communications*, 199(2):113-125, 2023. <https://doi.org/10.1016/j.comcom.2022.12.010>
- [8] Chaimae Hazman, Azidine Guezzaz, Said Benkirane, and Mourade Azrour. IIDS-SIoEL: Intrusion detection framework for IoT-based smart environments security using ensemble learning. *Cluster Computing*, 26(6):4069-4083, 2023. <https://doi.org/10.1007/s10586-022-03810-0>
- [9] Reza Ghanbarzadeh, Ali Hosseinalipour, and Ali Ghaffari. A novel network intrusion detection method based on metaheuristic optimisation algorithms. *Journal of Ambient Intelligence and Humanized Computing*, 14(6):7575-7592, 2023. <https://doi.org/10.1007/s12652-023-04571-3>
- [10] Omar Elnakib, Eman Shaaban, Mohamed Mahmoud, and Karim Emar. EIDM: Deep learning model for IoT intrusion detection systems. *The Journal of Supercomputing*, 79(12):13241-13261, 2023. <https://doi.org/10.1007/s11227-023-05197-0>
- [11] Mouaad Mohy-eddine, Azidine Guezzaz, Said Benkirane, and Mourade Azrour. An efficient network intrusion detection model for IoT security using K-NN classifier and feature selection. *Multimedia Tools and Applications*, 82(15):23615-23633, 2023. <https://doi.org/10.1007/s11042-023-14795-2>
- [12] Yazeed Alotaibi, and Mohammad Ilyas. Ensemble-learning framework for intrusion detection to enhance internet of things' devices security. *Sensors*,

- 23(12):5568-5587, 2023.
<https://doi.org/10.3390/s23125568>
- [13] Abdel-Rahman Al-Ghuwairi, Yousef Sharrab, Dimah Al-Fraihat, Majed AlElaimat, Ayoub Alsarhan, and Abdulmohsen Algarni. Intrusion detection in cloud computing based on time series anomalies utilizing machine learning. *Journal of Cloud Computing*, 12(1):127-143, 2023. <https://doi.org/10.1186/s13677-023-00491-x>
- [14] Ernest Ntuzikira, Wang Lei, Fahad Alblehai, Kiran Saleem, and Muhammad Ali Lodhi. Secure and privacy-preserving intrusion detection and prevention in the internet of unmanned aerial vehicles. *Sensors*, 23(19):8077-8104, 2023. <https://doi.org/10.3390/s23198077>
- [15] Nadir Omer, Ahmed H. Samak, Ahmed I. Taloba, and Rasha M. Abd El-Aziz. A novel optimized probabilistic neural network approach for intrusion detection and categorization. *Alexandria Engineering Journal*, 72(6):351-361, 2023. <https://doi.org/10.1016/j.aej.2023.03.093>
- [16] Yujun Wang. Deep learning models in computer data mining for intrusion detection. *Informatica*, 47(4):555-568, 2023. <https://doi.org/10.31449/inf.v47i4.4942>
- [17] Zhenpeng Zhang. SD-WSN network security detection methods for online network education. *Informatica*, 48(21):51-66, 2024. <https://doi.org/10.31449/inf.v48i21.6257>
- [18] Nour Moustafa, Nickolaos Koroniotis, Marwa Keshk, Albert Y. Zomaya, and Zahir Tari. Explainable intrusion detection for cyber defences in the internet of things: Opportunities and solutions. *IEEE Communications Surveys & Tutorials*, 25(3):1775-1807, 2023. <https://doi.org/10.1109/COMST.2023.3280465>
- [19] James Halvorsen, Clemente Izurieta, Haipeng Cai, and Assefaw Gebremedhin. Applying generative machine learning to intrusion detection: A systematic mapping study and review. *ACM Computing Surveys*, 56(10):1-33, 2024. <https://doi.org/10.1145/3659575>
- [20] S. Sivamohan and S. S. Sridhar. An optimized model for network intrusion detection systems in industry 4.0 using XAI based Bi-LSTM framework. *Neural Computing and Applications*, 35(15):11459-11475, 2023. <https://doi.org/10.1007/s00521-023-08319-0>
- [21] Ngamba Thockchom, Moirangthem Marjit Singh, and Utpal Nandi. A novel ensemble learning-based model for network intrusion detection. *Complex & Intelligent Systems*, 9(5):5693-5714, 2023. <https://doi.org/10.1007/s40747-023-01013-7>
- [22] Md. Alamin Talukder, Selina Sharmin, Md Ashraf Uddin, Md Manowarul Islam, and Sunil Aryal. MLSTL-WSN: machine learning-based intrusion detection using SMOTETomek in WSNs. *International Journal of Information Security*, 23(3):2139-2158, 2024. <https://doi.org/10.1007/s10207-024-00833-z>

A Hybrid Mamdani Fuzzy Inference System and Generalized Regression Neural Network for Cost and Time Overrun Prediction in Expressway Construction Projects

Tong Yao¹, Xiao Luo^{2*}

¹School of economics and management, Sichuan Tourism University, Sichuan, 610100, China

²Overseas Education College, Chengdu University, Chengdu, Sichuan, 610106, China

E-mail: love_luoxiao1989@126.com

*Corresponding author

Keywords: cost and time overrun, mamdani fuzzy inference system, generalized regression neural network model, expressway construction and risk factors

Received: June 18, 2025

In construction projects like expressways, risks and uncertainties are unavoidable and have the potential to significantly alter the expected outcome, which would be detrimental to the project's success. Risk is one of the main causes of productivity and efficiency losses in the construction sector that results in project demise, disputes, costs, and time overruns. The failure to complete the construction project within the stipulated time and estimated cost due to various risk factors is a major problem nowadays. This work analyzes the risk factors resulting in cost and time overrun in expressways construction projects using the Mamdani Fuzzy Inference System and Generalized Regression Neural Network (H-MFIS-GRN2) hybridization. Initially, the MFIS is used to find and measure potential dangers in the building process by analysing expert-made fuzzy rules. Defuzzification of the outputs yields clear risk severity values, which are further weighted with the use of mean scores and the Relative Importance Index (RII). Learning to anticipate budget and schedule overruns, MFIS doesn't provide the final product risk variables as normalized and prioritized inputs. The GRN2 model trains the observed risk factors with a pattern using the Gaussian activation function in a single pass. The model was validated using data from 27 areas from completed highway building projects. The trials show that the H-MFIS-GRN2 model outperforms baseline models. These baseline models are HRF-GA, H-AHP-ANN, and F-MRA. The H-MFIS-GRN2 model has 92.5% accuracy and 5.3% MAPE. Comparison analysis has increased forecast accuracy and interpretability, helping prioritize and minimize key risk variables. Fuzzy logic and neural networks can be used together to detect early risk in major infrastructure projects due to their strengths in uncertainty and learning.

Povzetek: V članku je opisan sistem za napovedovanje prekoračitev stroškov in časa pri gradnji avtocest. Hibridni algoritem H-MFIS-GRN2 združuje Mamdanijev sistem zamegljenega sklepanja (MFIS) in regresijsko nevronske mrežo (GRN2) za obravnavo tveganj in negotovosti.

1 Introduction

The planning, engineering, and construction of high-speed roadways specifically intended for quick and effective transit is called expressway construction. Expressways, sometimes called motorways or highways, are built to manage heavy traffic, enabling efficient vehicle movement and uninterrupted travel. The construction of expressways ranks with the highest frequency of mishaps, fatalities, delays, and cost overruns, largely because of unmanaged risks. With grade partitions at key intersections, an expressway is a segment of an urban highway with variable restrictions on entry. The term "overrun risk factors" refers to the variables that may cause schedule and expense overruns, resulting in delays and greater expenditures. These variables may change according to the project, the region, and other particulars. Risk is characterized as an unpredictable, unforeseen event or

circumstance that could have an unanticipated impact on a minimum of one project objective like schedule, expense, quality performance, etc. The risk might come from the bidding process, weather changes, job site effectiveness, political circumstances, market rivalry, etc., in constructing expressway fields. An efficient risk management system can reduce its negative effects on project objectives. In expressway construction projects, overrunning the risk factors and uncertainties are unavoidable and has the potential to significantly alter the planned result, which would be detrimental to the project's success.

Understanding the root of any uncertainty in managing expressway construction projects demands going above a simple inquiry into the expenses and schedule. Key aspects like "the imprecision, ambiguity, and uncertainty of the risk variables are essential for an expressway side construction to properly cope with a contractor's risks to

the project using Fuzzy Set Theory (FST). Several risk factors impacting construction project's time and cost overrun been found utilizing various statistical and computational techniques. Fuzzy logic is one such concept that can be applied. There are many unknown risk factors involved. In the case of the construction sector, it has the potential to produce reliable outcomes. Fuzzy techniques have been widely embraced as hybridized tactics for construction project risk analysis because they are incredibly effective at quantifying the risks experienced in complicated construction endeavors. The usage of neural networks is justified by their adaptability and propensity to anticipate and categorize all types of data more accurately than any other type of classifier. Hence, various neural network approaches are analyzed for predicting overrunning risk factors in expressway construction projects.

The main contributions of the article include

1. A new hybrid model (H-MFIS-GRN2) combines Mamdani-type Fuzzy Inference Systems (MFIS) with Generalized Regression Neural Networks to anticipate expressway building project cost and time overruns.
2. The model overcomes fuzzy-only models' static constraint and neural networks' black-box character by combining data-driven learning with expert-driven fuzzy rule logic, improving interpretability and prediction accuracy.
3. The authors propose a triangular fuzzy membership function, language-variable, and expert-review-derived Relative Importance Index (RII) risk prioritization method.
4. Single-pass training and Gaussian kernel activation help the GRN2 model beat baseline models like F-MRA, H-AHP-ANN, and HRF-GA with 92.5% prediction accuracy and 5.3% MAPE.
5. Validated on 27 Indian expressway projects, the model provides operational, financial, regulatory, investment, and climatic risk information. These parameters enable practical risk minimization.

The automation, productivity, and dependability of the construction sector are significantly increasing, and the sector is changing throughout the whole project life cycle, including the planning of projects, development, operation, and servicing [1]. The construction business projects, which greatly boost a nation's gross domestic product, include the transportation sector as a crucial component. Highways, primary intercity roads, toll roads, and other significant roadways with bridges, expressways, ducts, and tunnels are all included in road networks [2]. A very challenging and complex process is to analyze the time, costs, and risks that accompany a construction project while accounting for the specifics of all investments and the differing nature of its completion situations; hence fuzzy logic is applied to solve the above issues and imprecise information in [3]. Tiruneh et al. [4] studied

system training, assessment, and decision-making with neuro-fuzzy hybrid systems utilized in the predictive modeling of construction project challenges and concentrated information about the model's precision and accessibility performance. The applied neural network model was used to predict cost coated at highway construction projects' bidding stages [5]. Correlation analysis is used to identify the project attributes associated with a cost for lowering a bid, a technical, budget, etc. This method affects the dependent variables due to its limited number of factors (four) for analysis. Gondia et al. [6] analyzed and predicted the project's postponement risks in construction using Decision Tree and Naive Bayes models utilizing objective data from prior projects' time and delay influencing factors. Alawad et al. [7] proposed a computer vision with a Convolutional neural network model for safety risk assessment in railways to avoid high-risk possibilities. However, it does not focus on factors of risk overrun. To concentrate on risk factors like cost, quality, delay, and high economic risk, Andrić et al. [8] applied a combined approach of probabilistic model, fuzzy logic, and matrices with sensitivity analysis for belt and road initiative projects. The above-said method failed to concentrate on the dynamic nature of risk throughout the project. For implementing the dynamic nature of risk categories related to public and private partnership projects [9], an applied fuzzy interpretative structural model was used to create a link between the risks. Although fuzzy theory lessens the subjective aspect of expert judgment, it also makes processing information more time-consuming and expensive, necessitating more funding for modeling and analysis that affects the project's success. Hence, Isah and Kim [10] proposed a study for risk evaluation in expressway projects and incorporated an expert choice into the deep neural network to forecast the optimal cost and project allocation performance for successful outcomes. Risk assessment now includes a new source of uncertainty and unpredictability due to the expert judgment's subjectivity. Afzal [11] indicated that to more accurately represent complexity-based risk relationships under uncertainty, a mixed approach combining fuzzy logic and an enhanced form of the Bayesian belief network (BBN) may be used in the cost-risk analysis. This study exclusively addresses the individual risk evaluation methodologies used in expressway construction management to address the issue of cost overruns.

To avoid single cost factor estimation, Petroutsatou et al. [12] introduced a probabilistic approach for calculating the life-cycle expenses of road tunnels. Subsequently, tries to comprehensively capture their innate ambiguity, facilitating more trustworthy making choices at the start of a project. Once the input data are stored, a variant of the ANN model may generalize from them and learn the risk overruns in a single data pass [13]. Due to the local minimums of the variance requirement, the neural network method can converge to good solutions and is relatively easy to simulate the relationship between the risk overrun

factors in construction projects. Hence, Elbashbishy et al. [14] created a set of genetic algorithms (GA2) and ANN models to evaluate the effect of construction project risks on cost overruns rather than a single-point calculation of risk factors. Even though the existing works perform well related to the analysis of overrun risk factors, there is room for advancement in each work mentioned with its research gaps that need to be solved in the proposed scheme. Inspired by the existing works, to solve the limitations of single risk factor analysis of cost, static nature of risk analysis, solely based on expert choice, and region-based analysis, the proposed model handles the dynamic nature of various risk factors based on different dimensions in all areas starting from the initial stage of the project.

The study's main objective is i) to create a model that can accurately assess the intricate relationships between numerous project risk components and their impact on cost and time overruns using the relative importance index and mean score method in the MFIS model. ii) to develop a predictive model based on project risk parameters and historical data that can precisely predict future cost and time variations using the Gaussian kernel activation function from the project plan using GRN2. iii) The results are obtained by calculating accuracy, MAPE, and influencing factors-based cost and time overrun factors.

Although the H-MFIS-GRN2 framework suggests a new way to combine neural network and fuzzy logic techniques, the study is still led by well-defined goals and testable hypotheses. The main goal of this research is to create and assess a hybrid risk prediction model for expressway building projects that uses data-driven learning (GRN2) and expert-driven fuzzy reasoning (MFIS) to anticipate when and how much money will go over budget. As a result, we arrive at the following important research topics:

(RQ1:) Is it possible for the hybrid MFIS-GRN2 model to forecast time and cost overruns based on risk variables better than current models like F-MRA, H-AHP-ANN, and HRF-GA?

(RQ2): When working with subjective and uncertain construction data, can integrating fuzzy rule-based reasoning with GRN2 increase interpretability and generalization? (RQ3): What is the efficacy of MFIS in dealing with complex and unclear risk indicators and converting them into weighted inputs for predictive modeling?

In line with this, the study postulates that

(H1): When tested on actual project data in benchmarking trials, the H-MFIS-GRN2 model will attain a prediction accuracy more than 90% and a MAPE less than 6%.

(H2) By maintaining rule-based decision traceability via fuzzy inference, the model will offer increased interpretability compared to black-box machine learning methods. At the outset, the MFIS part of the goal is to

capture and quantify expert knowledge using fuzzy rules and trapezoidal membership functions so that uncertain risk dimensions (such as operational, financial, and climatic risks) may be effectively categorized. The second stage is completed by GRN2, which uses Gaussian kernel-based regression to generalize the learned patterns from previous project data. It does this by passing the defuzzified outputs, which have been weighted using the Relative Importance Index (RII). All of these parts work together to help project managers make better decisions by balancing the need for explanations with the need for accurate predictions. To make sure the methodology is in line with the study's goals, these hypotheses and objectives also serve as a foundation for comparing the hybrid model's performance to previous efforts.

The discussion of the remaining work is structured in the following way: Sect II discusses the existing research articles relevant to construction project overrun risk factors. Sect. III implements the detailed procedure of a hybrid combination of FMIS-GRN2. Section IV discusses the experimental analysis, including a comparative study analysis of various metrics. Sect. V summarizes the overall work procedure and contributions related to the research idea with the addition of future enhancement.

2 Literature background

Different fuzzy combinations have tackled individual uncertainty and unpredictability within construction project management situations. Risk factors were identified and classified through this literature review.

Chattapadhyay et al. [15] applied Genetic Algorithm with K-means clustering (GA+Kmeans), with the Euclidean distance method and silhouette coefficient for centroid optimization, to identify high-risk factors correlated with sub-risk categories. Applied statistics method for analyzing the most important features and achieved desired performance through cost-effectiveness, prompt project completion, high quality, and improved project scope. However, this research had the restriction of effectively acquiring inputs from experts since it was frequently difficult to persuade the participants regarding the relevance of their inputs, involving input from humans as expert opinion and discussions.

Yaseen et al. [16] created the Hybrid Random Forest classifier with Genetic Algorithm optimization (HRF-GA) to predict construction project challenges especially delay problems. The algorithm supports uncertain, complex, and dynamic environmental tasks and predicts the project performance based on these risk factors collected from Iraq. HRF-GA achieved 92%, 87.2%, and 0.833% of the measured accuracy, kappa, and classification error while identifying the delay of the construction project. It also identified precision and recall. Considering how these outside influences, like environmental changes, may affect the delay prediction model's accuracy is crucial.

El-Kholy [17] investigated the most effective models for forecasting delays and overruns in cost percentages for

expressway projects using different ANNs-based algorithms: Principal Component Analysis (PCA), modular N2, Radial Basis Function (RBF) generalized regression model/ probabilistic N2, and time-lag recurrent network. The research results produce a MAPE of 25% of cost overrun. The research limitation is that it is not applicable globally due to the construction's rapid growth environment and requires frequent updates regarding factors that impact cost and delay overrun.

Hung [18] applied the Artificial Neural Network (ANN) method for risk evaluation in construction projects to help contractors at the initial stage. The failure mode and effect analysis approach is used for assessing risk factors, while the ANN technique measures the effect of risk on contractor profit. The result shows that the minimum error rate during the training period is 0.004, and for testing is 0.035. Although ANNs are good at capturing complicated correlations and patterns in the data, it can be difficult to comprehend the reasoning underlying the predictions made by the model and pinpoint the precise elements that go into the risk assessment.

Lin et al. [19] employed a hybridized Analytic Hierarchy Process (AHP), and ANN(H-AHP-ANN) was employed to develop a framework for predicting the 19 important risk factors of the construction quality of Taiwan projects. The weightage of significant risky elements to confirm their impact on construction was determined using AHP. The ANN was used to identify the characteristics of significant risk indicators and estimate the caliber of a construction project. The suggested approach produced an accuracy of 0.852% with ANN and construction audit data to forecast outcomes related to project quality; nevertheless, a weakness of the study is that the choice of risk indicators relied on experts or authorities.

Sharma et al. [20] analyzed the level of risk that overrunning of expense elements has been estimated using a new fuzzy-based approach, which considers the problems' ambiguity, uncertainty, and subjective nature. The degree of severity and likelihood ratio determines the risky elements contributing to expense overrun. An innovative measure for overruns in cost elements indicates a particular factor's risk magnitude, known as the fuzzy Index for cost overrun. The research drawbacks are there may have been more respondents picked for the study. There were also very few professionals to provide their opinions based on the sector's size.

Ashtari et al. [21] employed a Bayesian Network (BN) classifier model to forecast cost overrun and evaluate risks by considering potential interactions between predictors. The current analysis determined that inflation, a rise in the cost of resources, and a lack of experience and knowledge among the workforce were the three most important threats with risk as an input factor. According to the findings, the 18 BN models had an average accuracy of predicting 79.1% with 10fold cross-validation of training and evaluation methods. Before using the model in other contexts, one

should carefully assess its relevance to various construction projects, locales, or periods.

Yun et al. [22] gathered socio-geospatial features from various data sets and built a Random Forest (RF) model to find their relationships with cost overrun. The developed models identify extremely important characteristics that affect overrun of cost, such as the original sums, unique time frame, management regions, total number of lane sets, commuting sequences, technological terrain, and temperatures on average, demonstrating that social and economic circumstances have a significant impact on real project expenditures. Lack of focus on other elements, including project management techniques, contractor efficiency, and decision-making processes.

Zafar et al. [23] demonstrated the Fuzzy Synthetic Evaluation (FSE) of the construction site's time overrun risk factor. The most important factor categorizing is stakeholder influence and threats to security, in which unproductive workers and workers, inefficient scheduling and agreements, and shortages of construction materials follow. Due to the safety circumstances, the study sought the perspectives of fewer professionals and business specialists, and time overrun risk factors are applied only to a specific region; hence it is not a generalized approach.

Sharma et al. [24] employed the average medium Fuzzy influencing score and the relevant frequency index result by Multiple Regression Analysis(F-MRA) models. The time overrun-inducing aspects were categorized into risk groups, notably red, yellow, and green zone. The created model can predict the anticipated time overrun of any forthcoming construction project with a confidence level of 79.6%. The model identified the top 5 risk influencing factors for cost overrun through the survey: land usage, material shift, under traffic, lack of a proper plan, and blueprint change. Due to the changing nature of construction projects, the main causes of time delays need to be updated regularly.

Informatica has been employing fuzzy systems to facilitate infrastructure decisions. The MFIS module's fuzzy inference is congruent with Chatterjee Performance results are and Banerjee's (2023) expert system for dynamic risk assessment in infrastructure building, highlighting the requirement for adaptive rule logic in uncertain situations. Kaur & Kaur (2023) use fuzzy-AHP and TOPSIS to make multi-criteria assessments, like our MFIS-GRN2 hybrid model. This supports merging expert-driven and data-driven methodologies.

Due to its complementarity, the hybrid MFIS-GRN2 model is employed in limited data, expert knowledge, and uncertainty settings. Mamdani-type Fuzzy Inference System (MFIS) lets domain specialists embed qualitative rules using language variables, making risk assessment more apparent and explicable. This improves interpretation. This is especially beneficial for infrastructure projects since subjective judgment and risk

choices must be traceable. However, data-driven Adaptive Neuro-Fuzzy Inference Systems (ANFIS) hide the taught rule structure, making them difficult to grasp. In low-data areas, ANFIS needs larger training datasets to converge. Due to its non-parametric nature, Generalized Regression Neural Networks (GRN2) with strong generalization are chosen over Genetic Algorithm-based models like HRF-GA for small to medium datasets. Unlike backpropagation

networks, GRN2 learns quickly and approximates functions smoothly without weight adjustment. For sparse data cost and time overrun prediction, this is essential. By combining MFIS for fuzzy rule-based reasoning with GRN2 for data-driven regression, the hybrid framework addresses the drawbacks of opaque (GA-based) or dense data (ANFIS) models. The consequence is interpretability and forecast accuracy.

Table 1: Literature survey

Method	Model Type	Dataset / Domain	Performance (Accuracy / MAPE)	Key Features	Limitations
HRF-GA (Yaseen et al., 2020)	Hybrid Random Forest + Genetic Algorithm	Risk delay data from Iraqi construction projects	Accuracy: 92% Kappa: 87.2%	Handles complex and uncertain delay factors	Region-specific; lacks interpretability of rule-based systems
H-AHP-ANN (Lin et al., 2022)	Analytic Hierarchy Process + Artificial Neural Network	19 quality risk factors in Taiwanese projects	Accuracy: 85.2%	Combines expert ranking with data-driven learning	Risk indicators rely solely on expert judgment
F-MRA (Sharma et al., 2021)	Fuzzy Multiple Regression Analysis	Survey-based time overrun data	Accuracy: Not reported Confidence: 79.6%	Zones (red/yellow/green) to prioritize time risks	Static rules; lacks adaptive learning from data
BN Classifier (Ashtari et al., 2022)	Bayesian Network	Cost-risk interdependencies in various construction sites	Accuracy: 79.1% (10-fold CV)	Captures probabilistic relationships between risks	Model calibration is complex and dataset-specific
ANN Models (El-Kholy, 2021)	PCA, RBF, Time-lagged ANN	Expressway delay and cost %	MAPE: 25%	Evaluates time-lag impact in project delays	Frequent updates needed; not globally generalizable
FSE (Zafar et al., 2022)	Fuzzy Synthetic Evaluation	Highway projects in conflict zones	Accuracy: Not quantified	Emphasizes stakeholder influence and security risks	Limited to regional, socio-political factors
H-MFIS-GRN2 (proposed model)	Hybrid MFIS + GRN2	27 Indian expressway projects	Accuracy: 92.5%, MAPE: 5.3%	Combines expert rules and learning; handles dynamic risk profiles	Limited sample size; more projects needed for generalization

3 Proposed scheme

Risks associated with construction projects include environmental effects, financial overruns, delays, and safety risks. Due to their distinctive features, such as independence, diversity, and uncertainty, construction projects are complicated and fraught with risk. Construction projects interact with the surrounding environment rather than running independently. If all of the data is unknown, the issue is uncertain. Because most of the risks associated with building projects are stochastic and unclear, developing an effective algorithm and decision-making process is impossible. The success of a project depends on early cost estimation and cost management during the construction phases. Unfortunately, cost overruns in expressway construction projects are widespread and occur under various economic and regulatory circumstances, which frequently harm achieving project objectives.

GRN2 avoids local minima traps and converges faster than iterative backpropagation-based ANN models. It excels in

infrastructure risk assessments and other sectors with minimal training data and noisy inputs. GRN2 is computationally cheap, generalizable, and accurate, unlike GA-based models like HRF-GA, which need extensive tuning.

Figure 1 analyzes the possible risk factors from various expressway construction projects. The collected input variables are passed to the MFIS model to identify risks and uncertainty in the construction investment process and those related to construction initiatives timing and cost problems. An accurate estimate of potential risks guarantees the accomplishment of the project's objectives. The analysis can include linguistic factors and expert knowledge to produce a rule-based system using an MFIS. Those experts were chosen based on their professional backgrounds, and administration positions in the construction industry may involve assistant engineers and deputy engineers. This system can include a variety of input factors. Both fuzzy sets and functions for

membership can be used to define these variables. The fuzzy inference system then processes the input variables. It establishes the membership level for every result factor, specifically the cost and time overrun, using fuzzy logic operations like fuzzy rules and implications. The output variables have distinct fuzzy sets and functions of membership defined for them. A knowledge base is a body of information compiled through human experience and used to simulate how the inputs and results of a process related to one another. Rules are used to communicate knowledge, and using linguistic variables is the most prevalent rule framework in MFIS.

MFIS converts expert language evaluations into numerical risk weights that the GRN2 model may use as structured input features to anticipate numerical cost and time overruns.

A particular kind of neural network that works well for tasks involving regression analysis is the GRN2. It makes use of a network architecture known as radial basis function (RBF), which is able to generalize well by learning from the correlations between input and output data. The non-linear correlations between the input factors, which include project scope, operation detail, technical knowledge, weather conditions, etc., and the related cost and time overruns can be captured efficiently by the GRN2.

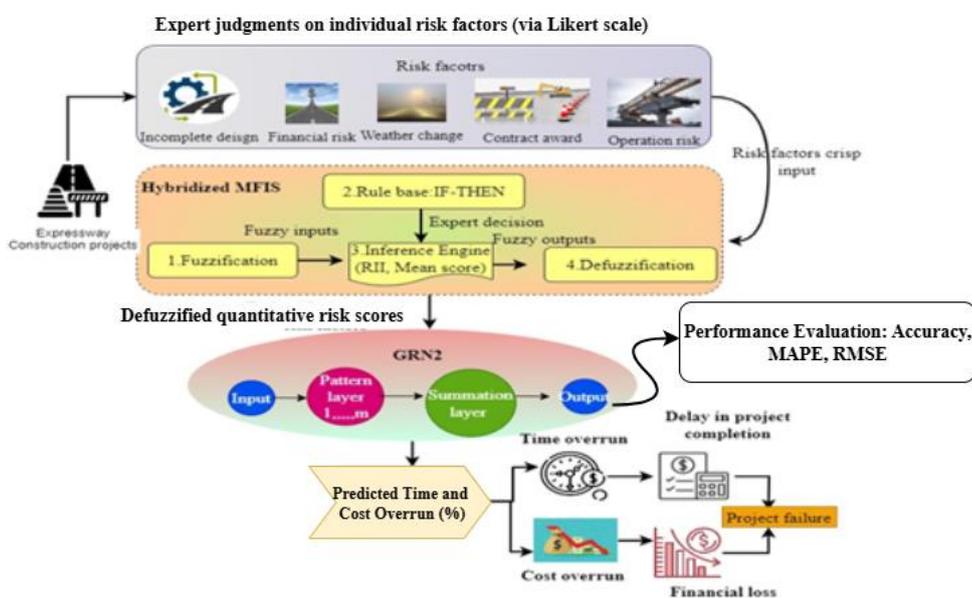


Figure 1: Systematic approach of H-MFIS-GRN2 Algorithm

With this GRN2 model, the four layers are employed to predict the cost and time overrun due to the influencing of risk factors applied in the training of GRN2. If a cost overrun occurs, it leads to financial loss that affects the economic condition, and similarly, the time overrun delays the completion of the project schedule. If both these conditions occur in construction projects, it leads to project failure. This emphasizes the importance of finding suitable risk reduction techniques in the future. The MFIS module utilizes relative significance indexing, mean score aggregation, and fuzzy inference based on expert-defined rules to provide a collection of defuzzified, weighted risk factors, as shown in Figure 1. These outputs are structured input characteristics used to train the GRN2 prediction model; they are prioritized representations of construction hazards.

Figure 1 clearly labels the GRN2 output node as "Prediction output (accuracy, MAPE)" for simplicity. The GRN2 model does, however, directly forecast percentages of both time and expense overruns. The image displays the post-prediction assessment metrics that contrast the

model's predictions with the dataset's actual overruns. In this context, the terms "accuracy" and "MAPE" are utilized.

Despite its broad definitions of "risk overrun" and "project success," the MFIS is not utilized alone to forecast project results. Instead, a semantic pre-processing layer rates risk factor severity using expert-based fuzzy inference. The Mean Score and Relative Importance Index refine these numerical risk severity values from these language assessments after defuzzification. Final time and cost overrun predictions are made by the GRN2 model using a feature vector constructed with defuzzified weighted risk variables. MFIS doesn't provide the final product, but it gives GRN2 standardized and interpretable risk representations.

3.1 Mamdani type of fuzzy inference system (MFIS)

The MFIS model's functionality is based on three consecutively running blocks: fuzzification, inference engine, and defuzzification.

Step 1 shows that the MFIS can handle discrete risk factor inputs such as scope quality, climate effect, technical competency, etc. Inserting input variables and linguistic terms into a normalized [0,1] domain.

Step 2 a triangle membership functions fuzzify risk factors. It appears that the MFIS turns human risk severity evaluations into fuzzy integers, which are then defuzzified. Many fuzzy rules utilize expressions like "time overrun" or "project success," they better represent language judgments of risk factor severity than project results. To prioritize inputs for the GRN2 model, the MFIS solely

$$\mu_{T(a)} = \begin{cases} 0, & a \leq x \\ \frac{a-x}{y-x}, & x \leq a \leq y \\ \frac{z-a}{z-y}, & y \leq a \leq z \\ 0, & z \leq a \end{cases} \quad (1)$$

Where a represents the input variable, the remaining variables x , y , and z denotes the influencing conditions to which the fuzzy rules have been derived for evaluating overrun factors of construction projects.

Figure 2 shows that the membership function plotted for risk factors as input variables with a range of [0:1]. The linguistic terms of the input variables are categorized into five levels, "very low," "low," "Moderate," "high," and "very high," to reflect the influencing level. Every point in

fuzzifies, processes, and defuzzes risk factor variables, resulting in normalized severity ratings. The hybrid design isolates MFIS (rule-based semantic reasoning).

Although "time overrun" and "project success" are employed to express fuzzy rules, they truly refer to language judgments of risk factor severity rather than project outcomes. All input risk factor variables are fuzzified, analyzed, and defuzzified by MFIS. Normalized severity ratings are priority GRN2 model inputs. GRN2 and rule-based semantic reasoning (MFIS) are separated in the hybrid design.

the given input space is mapped by a fuzzy membership function curve to the extent of membership that ranges from 0 to 1. The description of the fuzzy variable with the linguistic term is explained as [Very Low: above/below 0.1; Low: above 0.3; Moderate: above 0.5; High: above 0.7; Very High: above 0.9] for the risk factors. These triangular membership functions specify each input variable's membership level to various fuzzy sets, representing meaningful risk levels.

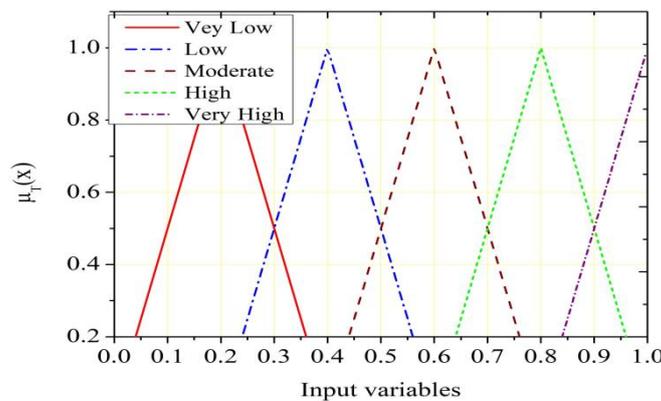


Figure 2: Membership function plots for risk factors

These $\mu_{T(a)}$ can be represented for various risk factors dimensions like operational, financial, regulation, investment, and other factors categorized for analyzing the cost and time overrun risks in the construction projects.

Step 3: Fuzzy rule-based systems use fuzzy sets and fuzzy logic to describe the humanistic knowledge about an issue and the interactions between its variables. The problem-related information that is currently known is added to the knowledge base as fuzzy IF-THEN rules. The

processing component performs the inference process on the risk inputs based on the knowledge base. Managers of construction projects can receive a quantitative evaluation or prediction of the possibility of overruns in time and

costs based on various influencing variables by using a fuzzy rule basis. To lessen the risk and effects of overruns, they can use this information to make wise decisions, prioritize their resources, and put the necessary mitigation measures in place.

IF x_1 is L_{t1} AND x_2 is L_{t2} OR x_3 is L_{t3} THEN Y is L_{t5}

Here the term x_1, x_2, x_3 are represented as input variables affecting risk factors, then Y is given as the output variable, then L_{t1}, L_{t2}, L_{t3} and L_{t5} are linguistic terms represented based on these variables.

Rule 1: IF the project scope is poor, THEN the risk overrun is high.

Rule 2: IF the climate is poor, THEN time overrun is high.

Rule 3: IF technical knowledge is low, AND contract design approval is poor, THEN the possibility of project success is very low.

Rule 4: IF delay in bank guarantees OR liquidity damage is high, THEN the project outcome is risk overrun is very high.

Rule 5: IF contract award is late, AND excavation works delay, THEN risk overrun is High.

Fuzzy rules that link the input parameters to output values were established in the third stage to carry out the inference. The Relative Importance Index RII_{Ec} cost and time overrun were discovered, and the relative proportion was assigned as weightage to the fuzzy rules to construct the assessment model. Give each detected risk indicator a score based on its estimated importance or possible influence on cost and time overrun. In a fuzzy membership function, Equation 2 can be used to calculate the weights from the average score of variables.

$$W_i = \frac{M_s}{\sum M_{s(\mu)}} \tag{2}$$

Where M_s represents the mean score of risk overrun factors and the $M_{s(\mu)}$ denotes the sum of all mean scores of identified risk factors in a membership function. W_i represents the weighting score for risk factor i. The percentage scores of the replies for each risk factor over the 5-point Likert scale provided by the experts are used to calculate the membership function of every risk factor.

The relative importance measure Index is a statistical method for identifying risk rank factors predicted in the expressway construction projects that impacts the project's overall success using the weighted average method. It usually ranges from 0 to 1. Project administrators can identify and concentrate on the most important risks by determining the RII_{Ec} . The priorities and issues of various stakeholders can be evaluated using the RII_{Ec} . The RII_{Ec} can direct project managers in efficiently meeting stakeholders' requirements, strengthening interaction, and enhancing final project outcomes by allocating weights to their priorities or those risk parameters, which enables efficient resource allocation and measures to mitigate risks. MFIS' fuzzy rule basis and risk factor weight allocations were based on structured feedback from six subject-matter experts. The 10-year-experienced specialists build highways and expressways. All were commercial or government infrastructure assistant engineers, project managers, or deputy engineers. Candidates excelled in building project planning, construction, and quality. The researchers used a 5-point Likert scale to assess risk factor severity. Qualitative evaluations yielded final outcomes, relative significance indices (RIIs), and fuzzy language rules. We maintained consistency using expert feedback loops since there weren't enough experts to eliminate statistical outliers. This verifies and tracks MFIS rule base data.

3.1.1 Relative importance measure index

Table 2: Likert scale type expert analysis of risk factors

No.	Causes of risk occurrences	1	2	3	4	5
1	Design changes	x	✓	x	x	x
2	Insufficient bidding method	x	x	x	✓	x
3	Failures in the design model	x	x	✓	x	x
4	Financial constraints and inadequate fund allocation from the client	x	x	x	x	✓
5	Delay in getting approvals	✓	x	x	x	x
6	Incomplete drawings	x	✓	x	x	x
7	Lack of skilled technical staff	x	x	x	✓	x
8	Deficiency of materials, equipment, and tools	x	✓	x	x	x
9	Price variation of materials	x	x	✓	x	x
10	Unfavorable climate changes	x	✓	x	x	x
11	Delay in obtaining government permits and approvals	✓	x	x	x	x

Table 1 uses symbolic indicators (✓ or X) to visualize expert selections on a 5-point Likert scale from 1 (low relevance) to 5 (high significance). Equation (2) computes the mean score (Ms) and weighted score (W) by numerically encoding replies, with each ✓ mark representing its column value (e.g., ✓ under column '4' = 4 rating). The table structure is reduced for readability because all survey responses were numerical Likert values. These numeric ratings were used to calculate the relative relevance of risk factors in MFIS fuzzification and rule prioritizing across experts.

From Table 2, construction administration experts analyze the risk factors that lead to time and cost overrun of expressways, which can be easily evaluated with the help of the mean score method. The attributes named causes of risk occurrences are evaluated on the Likert scale range from [1 to 5] where the representation of scales given by experts represented that 1 as "very low" significance for risk occurrence, 2 represents the "low" importance of risk occurrence, 3 represents the "moderate" possibilities of risk, 4 represents "high" chances of risk, and 5 represents "very high" chances of risk overrun they choose the

response option that most accurately reflects their opinion or evaluation of the specific risk factor. The pertinent risk factors that affect a particular risk assessment are identified and defined by experts. Assigning exact numerical figures to these risk factors might be difficult because they can be evaluative in form. Based on their skill, experience, and understanding, each expert assigns a rating or score to each

$$RII_{E_c} = \frac{\sum_{j=1}^A R_{ij}}{A \times H} \tag{3a}$$

In Equation 3a, R_{ij} denotes Likert score (1-5) assigned by expert j to factor i , A total number of experts and H is the highest possible score on the scale. The Relative Importance Index (RII) assigns a number to each risk factor. The sum of all experts' assessments for each element may be normalized by multiplying A by the maximum Likert score ($H = 5$). Risk factors can be prioritized when RII values are guaranteed to be $[0, 1]$.

risk factor. Each risk factor's weight is considered after analysis when estimating the scale scores. The mean of the scores supplied by each participating expert is then used to get the mean score of risk variables. Understanding the total risk connected to each element is made easier with the help of the mean score, which offers an aggregated measurement of the expert's opinions.

To get the weighted scores, increase the scores for all risk factors by the associated weights. A higher relative relevance index (RII) suggests that the associated risk factor is more important. As illustrated in Figure 3, the number of expressway projects taken for analysis is limited to 5, with risk dimensions prioritizing the scores based on the mean score method.

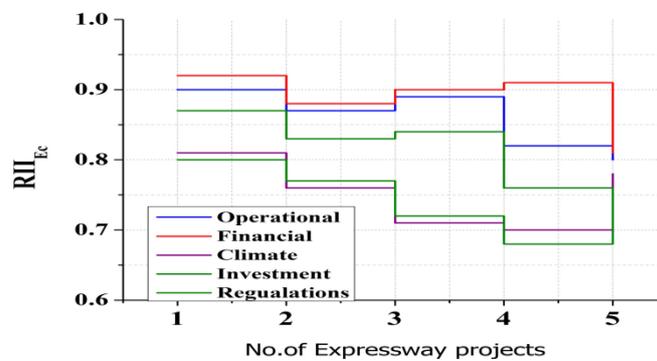


Figure 3: Analysis of projects using RII_{E_c}

Table 2 explains the ranking value given for each risk factor with its weightage assigned by experts with a clear analysis of fuzzy rules based on the assigned linguistic variables. Systems with several rules to express human reasoning and knowledge in the produced system are known as rule-based systems. Domain specialists manually built Table 2's larger risk dimensions (Operational, Financial, Climate, etc.) based on functional similarity and effect area. They then grouped Table 1's 27 granular risk indicators into five themes. The RII scores of each dimension's risk factors are summed.

Domain specialists manually built Table 2's larger risk dimensions (Operational, Financial, Climate, etc.) based on functional similarity and effect area. They then grouped Table 1's 27 granular risk indicators into five themes. The RII scores of each dimension's risk factors are summed. Specifically, for each dimension D_j , the aggregated RII is computed as in Equation 3b

$$RII_{D_j} = \frac{1}{n_j} \sum_{i=1}^{n_j} RII_i \tag{3b}$$

Where n_j is the number of risk factors belonging to dimension D_j , RII_i is the relative importance index of individual risk factor i , Qualitative reasoning and rule design were aided by MFIS linguistic mapping of these aggregated scores (e.g., Low, Medium, High influence). For the interpretation given in Table 2, Figure 7, and Figure 8, this mapping is crucial.

Table 3: RII_{E_c} based ranking

No.	Dimensions of Risk Factors	Linguistic variable	Ranking
1	Operational	High	2

2	Financial	Very High	1
3	Climate	Low	4
4	Investment	Moderate	3
5	Regulations	Very low	5

The established rule base foundation allows for the classification of cost and time/delay overrun risk into the following levels: "very low," "low," "Moderate," "high," and "very high." Various fuzzy rules were framed in the inference block to analyze the risk overrun category levels (Table 3).

Step 4: The aggregated model collected from fuzzy rules is considered fuzzy output defuzzified by the model in the fourth stage to provide crisp output values. The cost overrun percentage was calculated in the range of [0:100].

The resultant variable represents the likelihood that the costs and time of a certain construction project unit will end up being exceeded. These defuzzified risk variables are given as input to the next GRN2 prediction model to measure construction projects' cost and risk overruns accurately. This study's fuzzy rule foundation relied on domain experts' manual input, not optimization or automated rule extraction. No post-generation rule pruning or weight learning was done. The rules were constructed using language characteristics and expert consensus to assure semantic interpretability and domain knowledge alignment. The method improves transparency but isn't suitable for larger or more diversified datasets. Evolutionary algorithms and subtractive clustering may be helpful data-driven fuzzy rule learning and pruning strategies for model refinement. These approaches may dynamically optimize the rule base, reduce redundancy, and preserve interpretability.

3.2 Generalized regression neural network training model

The obtained defuzzified risk values are used for the next purpose, mainly for training the neural network model to recognize the cost and time overrun of an expressway construction project.

A generalized regression neural network (GRN2) is frequently utilized for feature estimation. It is a type of one-pass feed-forward type ANN. GRN2 represents a better method for creating neural networks based on nonparametric analysis. The GRNN uses neural networks to approximate or estimate functions and predicts the output from input data. Every training sample is supposed to represent an average corresponding to a radial basis neuron. It has a specific linear layer and a radial basis layer, whereas the accuracy and speed of the GRN2 training process are advantages. GRN2 has four processing units: Input, hidden/pattern, summation /add, and output neurons. The reason for choosing this neural network model is it can handle noisy information in the input and requires only a small amount of data for training purposes. Regression analysis is used to predict time and cost overruns in this study. Comparisons to baseline estimates show ongoing percentage overruns. GRN2 converges rapidly, performs well with nonlinear regression, and approximates complex input-output mappings without repetitive backpropagation. The model converts the MFIS-derived risk factor vector to % overrun values for more accurate time and cost deviation estimates.

3.2.1 GRN2 training

Considering the fuzzy-encoded input information and the related risk overrun factors, train the GRNN portion of the Fuzzy GRNN. Analyzing the initial training data, the GRNN discovers the relationships and trends among the input factors and the cost and schedule overrun probability. While training, it memorizes every pattern of risk overrun factors in a single time pass; hence, backpropagation is unnecessary.

Table 4: Parameters of developed Grn2 model

Setting up of parameters in the GRNN model	
Input layers	i1, i2, i3, i4.
Activation function	Gaussian kernel
Type of analysis	Regression

Smoothing factor	Σ
No. of neurons in the model	21 in hidden(pattern) layer
Learning cycles	1000

From Table 4 above, the initial parameter settings of the GRN2 model are assigned, and the internal behavior of GRN2 with the obtained fuzzy rule-based output values is trained well to easily predict the various dimensional behavior of construction projects.

i) The input layer of the GRNN is the scaled representation of the fuzzy risk overrun factors. A defuzzified crisp set of a certain risk overrun factor corresponds to each input neuron in the source layer. Then it forwards the feature variables to the hidden layer.

$$G(i, i_g) = e^{-(i-i_g)^T(i-i_g)/2\sigma^2}, \sigma > 0$$

Where i represents the input and i_g denotes the training samples. The squared Euclidean distance between the i and i_g is given as $-(i - i_g)^T(i - i_g)$. The standard deviation value or Gaussian activation function spread with a particular neuron is noted as σ if it ranges from 0.01 to 1, then it has a good analysis result. A bigger distance of results causes the term $G(i, i_g)$ smaller, leading to other training samples of risk overrun factors also being small.

The Gaussian kernel is excellent for this application because it simulates input pattern localization. Smooth function approximation is needed to predict construction risk overruns with sparse and faulty data. Regression concerns benefit from Gaussian kernel activation functions, which decrease overfitting and increase generalization. Instead of figuring out and altering the weights and the bias in each pattern/hidden layer as input is loaded into the model, it keeps the training information as the parameter's value. The model will total the scores of the other variables weighed by the radial basis function when the request is made, then estimate the value. Fuzzy-based GRNN architecture is illustrated in Figure 4.

ii) Following the input layer, the pattern layer performs a weighted summation of the received input features, then employs the activation function as Gaussian function kernel (G) and generates the output passed to the third layer. The hidden layers allow the network to understand complicated associations and produce precise predictions or estimates by applying non-linear modifications to the information provided. The activation function value for a neuron in the pattern layer is calculated by using Equation 4 as follows:

$$(4)$$

The GRN2 architecture was calibrated on the training set using grid search and 5-fold cross-validation. Five to twenty neurons were counted in the primary hyperparameter, the radial basis (pattern) layer. Accuracy and MAPE were used to validate performance. After evaluating different settings, twelve neurons balanced model complexity and prediction inaccuracy. The smoothing parameter (σ) of the Gaussian kernel function was adjusted simultaneously, with values from 0.1 to 1.0 considered in increments of 0.1. We got optimal results at $\sigma = 0.6$ by reducing validation MSE without overfitting. The GRN2 model generalizes successfully even with small training data changes thanks to these considerations. Fast convergence without iterative weight updates is made possible by the GRN2's organized design and single-pass learning. Overrun forecasts in terms of both time and money are guaranteed by its capacity to generalize from sparse data using the Gaussian kernel. Hyperparameter tweaking and k-fold cross-validation are used to optimize the architecture.

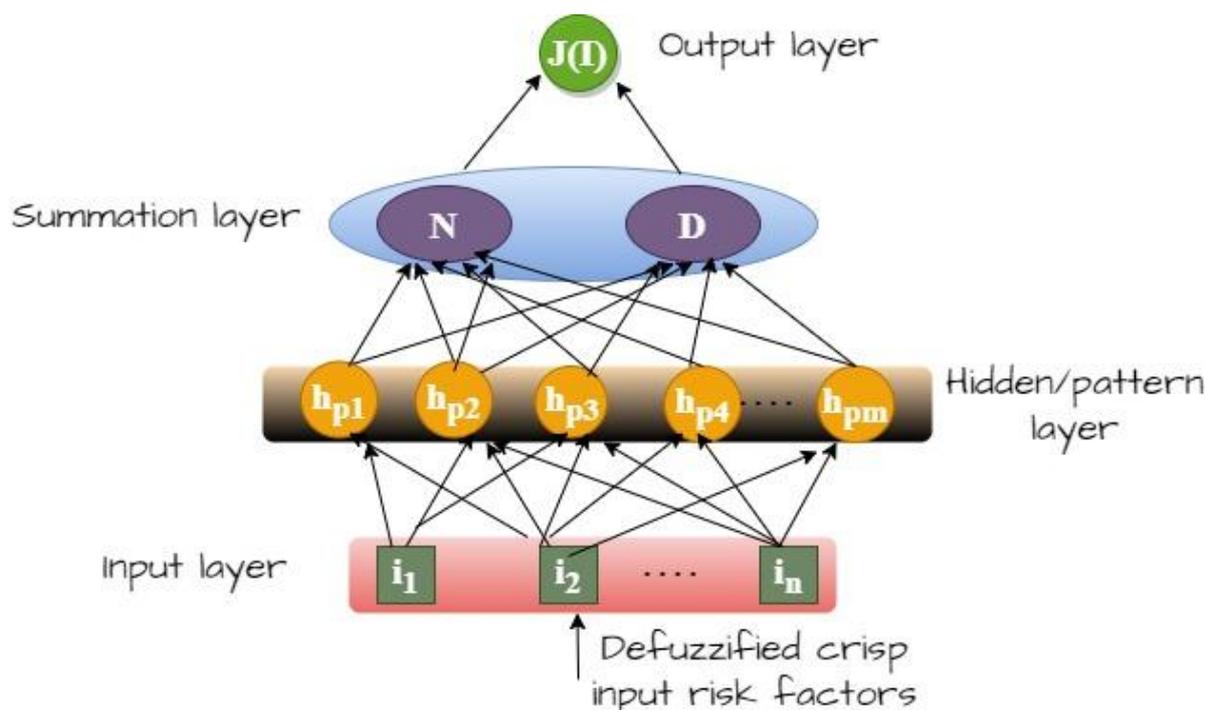


Figure 4: Fuzzy-based GRNN architecture

iii) Summation/Add layer: The summation layer combines the pattern layer's weighted outputs into one output value. Based on the input risk variables, this output value shows the anticipated risk overrun or the level of risk related to the expressway construction project in the form of 2 neuron nodes.

$$J(I) = \frac{\sum_{g=1}^N j_g G(i, i_g)}{\sum_{g=1}^N G(i, i_g)} \quad (5)$$

Where $J(I)$ denotes the analyzed prediction result of overrun risk variables from input i (predicted cost or timeoverun percentage). j_g denotes the Gaussian activation weight utilized in the pattern layer of neurons. $J(I)$, representing the target variable's anticipated value, is generated by the output layer. Here, $J(I)$ represents the predicted percentage of cost overrun or time overrun for the input vector of the project.

Once trained, the hybrid model can forecast risk overruns in new highway construction tasks. The validation and analysis have to be done. The model receives the fuzzified risk factors and predicts the risk overrun using the discovered associations. In expressway, construction projects, evaluate how the model collects and anticipates overrun risk aspects. To determine the project's risk level, analyze the expected risk overrun and compare it to established criteria. The decision-making procedures, risk-minimization tactics, and project execution can all benefit from this implementation. To interpret the connections between risk factors and overrunning parameters, analyze the learned model. The model's interpretability is made possible by fuzzy logic, which offers linguistic conventions that connect risk indicators to probable overruns from MFIS to GRN2. Finally, because

iv) Output layer: The layer receives information from numerator and denominator nodes and divides the two gathered data to analyze the predicted risk overrun factors for cost and time using Equation 5.

of its hybridized behavior, the suggested algorithm H-MFIS-GRN2 helps predict and analyze construction projects' cost and risk overrun.

Defuzzification in the MFIS module used the centroid (center of gravity) technique. The weighted average of all membership function results is used to create the crisp output to represent the inferred fuzzy region. It was chosen because of its sensitivity to the output fuzzy set's structure and distribution and its widespread application in engineering decision systems. Defuzzified outputs—defined risk severity levels—provide standardized, weighted inputs to the GRN2 prediction algorithm.

4 Experimental analysis

4.1 Data collection

The input data source taken from the expressway construction of 27 completed projects in the various regions is described with the year of start and follow-ups of the total project duration with cost during the contract award process from [27]. The project owner or customer assesses competing bids or proposals from several contractors during the contract award process. The choice of contractor is typically made based on some factors, such as projected costs, qualifications, experience, and the

intended project timetable. The collection comprises official sources ([27]- [28]) that cover 27 separate highway development projects in various areas of India. Metadata pertaining to the risk dimension, actual spending, scheduled vs. real length, and contract award cost are all included in each record. Before data was processed, it was

checked by comparing project reports and expert interviews, and numerical and categorical values pertaining to risk events and overrun metrics were extracted. Although there isn't a ton of data, it does cover a wide range of operational, regulatory, and financial risks.

Table 5: Risk dimensions with their types for analyzing risk importance

Risk dimensions	Influencing categories
operational impediments	Environment and forest clearance, loan restructures, rescheduling design criteria, lack of supervision.
financial risk impediments	Contract award problems, lack of equity, timely renewal of bank guarantees
regulation acts	Land acquisition, design approval, state support agreement, shifting utilities, and judicial interventions.
investment plans	Bidding, quality, contractual basis, excavation work, technical does
others	Toll, liquidated damages, an extension of time, technical follow-ups, weather conditions

From Table 5, the analysis of individual risk categories with their group dimensions related to construction projects is categorized into several types to identify the importance of risk. Among all attributes, the sample of a construction project in the region of Barwa Adda-Barakar is listed with a length of 42.69km with submission of bid on 5/1/1996 with actual data of construction awarded to initiate progress on 20/9/1996 with a delay in award of working is two months with the drop of 6months of government schedule exclusion. As per the data source, the scheduled month of completion is June 2000, but the actual month of completion is December 2001, with a delay of 18 months lacking from the scheduled month. The cost as per award is 155.00cr. The actual expense calculated up to August 2004 is 208.54cr. From this analysis, the cost overrun of the scheduled construction project is 53.54cr. The major challenges/risks dimension factors considered are operational impediments, financial risk impediments, regulation acts, investment plans, and others collected from [26].

4.2 Comparative analysis

The performance of the proposed algorithm H-MFIS-GRN2 is validated with various metrics, including accuracy, MAPE, Risk influencing factors, cost, and time overrun. For this implementation, the proposed concept is compared with existing algorithms like HRF-GA [15], H-AHP-ANN [19], and F-MRA [24]. The HRF-GA, H-AHP-ANN, and F-MRA baseline models were reimplemented and tested on the same dataset of 27 highway development projects using the same risk factor inputs, preprocessing methods, and evaluation metrics (accuracy and MAPE). Each model's hyperparameters were set to the original

publications' values for optimal performance. To eliminate biases from dataset size, data quality, or evaluation criteria, these results were reproduced under controlled experimental conditions. This controlled setting will show the model's true capabilities rather than experimental disparities.

Regression metrics like MAPE, MSE, and R³ excel at assessing continuous predictions like cost and time overruns. However, an accuracy-like metric was calculated by categorizing overrun values into predefined bins (e.g., Low: <10%, Moderate: 10-20%, High: >20%). A forecast was "right" if it matched the ground truth value. How effectively a regression model can be comprehended and compared to previous categorical models is more important than its classification accuracy. The H-MFIS-GRN2 model's robustness is assessed using confidence intervals, MAPE, and MSE.

4.2.1 Accuracy

An accuracy statistic from Equation 6 is used to assess a prediction model's total accuracy and efficacy. The study forecasts continuous factors, such as percentage cost and time overruns, therefore classification metrics are not appropriate. The predictions are typically evaluated using regression metrics like MSE and MAPE. To compare results, divided the projected overrun values into risk categories (e.g., low: <10%, moderate: 10-20%, high: >20%) and assessed classification accuracy accordingly. This thresholding was just used for interpretive visualization, not assessment. MAPE is used to validate predictions with continuous-value measurements. (Refer fig 5).

$$Accuracy_{Threshold} = \text{No. of correctly classified overrun levels} / \text{No. of total construction projects} \quad (6)$$

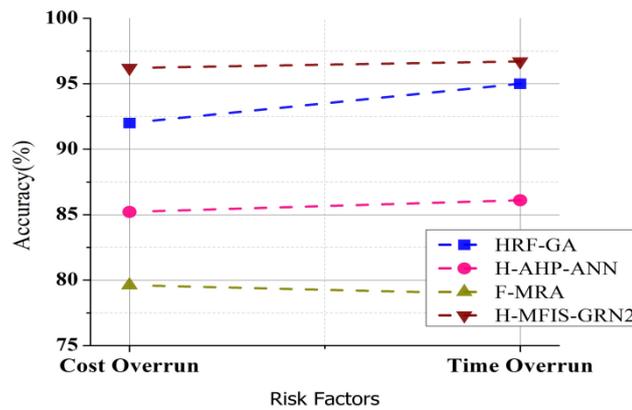


Figure 5: Comparison based on prediction accuracy

4.2.2 Mean absolute percentage error (MAPE):

The MAPE shows the average amount of prediction error as the average % variance between predicted and actual

$$M = \frac{1}{N} \sum \frac{|A_v - P_v|}{|A_v|} * 100 \tag{7}$$

From Equation 7, M represents the MAPE, A_v denotes the actual value of cost and time overrun and P_v

values. It is a widely used statistic, especially in predicting regression activities, to evaluate the precision and dependability of forecasts.

denotes the predicted value observed by the model for the same risk factor, and N denotes the total number of projects. In Figure 6, the smaller MAPE values suggest forecasts of risk overrun have become more accurate.

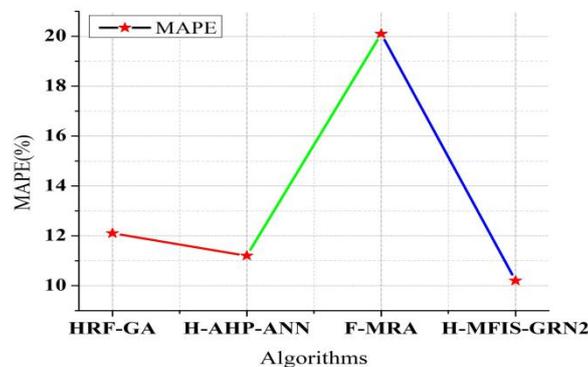


Figure 6: MAPE analysis

4.2.3 Time overrun calculation (T_o):

The project completion time/schedule deviation or delay about the initial predicted completion time act_T is calculated using the time overrun T_o calculation in response to identified risk factors.

$$T_o(\%) = ((act_T - adj_{eT}) / adj_{eT}) * 100 \tag{8}$$

$$adj_{eT} = eT + (eT * r_f) \tag{9}$$

Equations 8 and 9 indicate how much the project's timeline has fallen off plan.

r_f represents a aggregated risk factor weight allocated according to its possible influence on the project timeline, a risk score is given to each detected risk factor. A scalar number between 0 and 1 indicates how the identified risks affect the project timetable. Expert judgment and risk assessment ranking methods can be used to identify the risk factor. The adj_{eT} The expected duration/estimated schedule was modified to consider the risk variables.

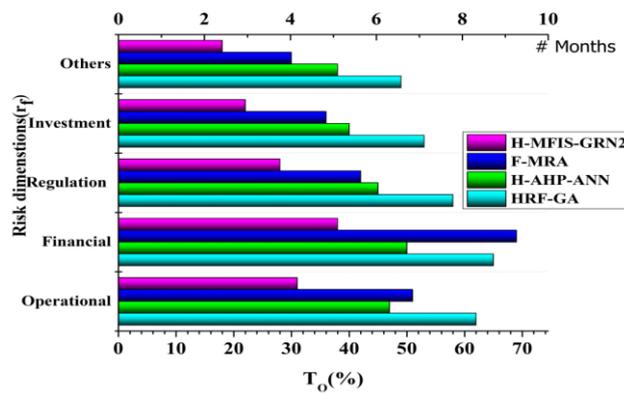


Figure 7: Time overrun $T_o(\%)$ analysis

As shown in Figure 7, the calculation is derived by adding the projected time compounded by the risk factor. r_f to account for any potential delays brought on by the risks identified. The symbol usage \exists represents there exists a presence of risk factors to calculate the T_o . The main advantage of this proposed H-MFIS-GRN2 algorithm is to accurately analyze the time overrun of project schedule by incorporating the possible risk factors obtained from expert decision with the rank obtained by RII_{Ec} using MFIS and GRN2 training schedules to consider risk impacts on project timelines.

4.2.4 Cost overrun calculation

The cost overrun (C_o) method from Equations 10 and 11 determines the proportion of the rise in a project's actual expenses above the original estimation. It indicates how much the project's expenses have increased from the projected budget. It's crucial to remember that the anticipated cost of construction projects comprises a variety of charges associated with the project, including resources, technical staff, machinery, permits, and other costs

$$C_o(\%) = ((act_c - B_c) / B_c) * 100 \tag{10}$$

Contract Score (Cts) is a derived measure that we establish to describe risk exposure unique to contracts, in addition to individual risk component scores. To account for cost overruns (Co) and the relative importance of important contract-related concerns, it uses:

$$Ct_s = (1 - C_o) \times \sum_{k=1}^n \omega_k \cdot c_k \tag{11}$$

Where c_k is the value obtained through expert likert scoring and normalized, ω_k is the weights are derived via RII analysis. The remaining contract risk after cost variations are adjusted is captured by Cts, which is an aggregated input to the GRN2 model. Here is the contract score Ct_s given to other risk evaluation factors, such as investment plans, the technical skill of experts, timeline, operational maintenance cost, etc., during the contract awarding process is represented by the term weighted score (W) of other criteria. Based on their respective importance, these criteria are given varying weights.

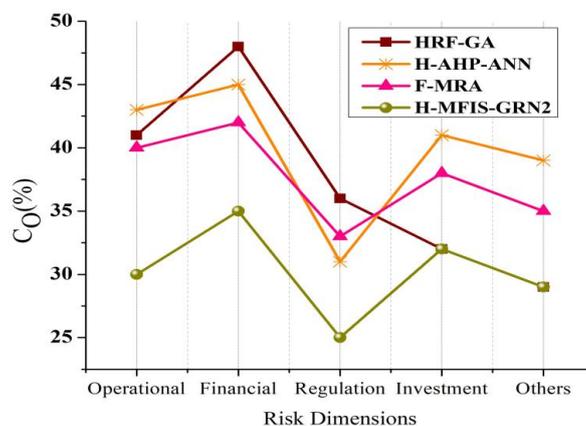


Figure 8: Cost overrun $C_o(\%)$ analysis.

Figure 8 shows the C_o analysis, others estimate the extent of cost overruns as the disparity between the cost at the award of the contract and the ultimate completion expenses. At the same time, some measure it as the cost at the moment of decision for construction vs. final

completing expenses. Ref Figure 9 for analyzing risk measurements of cost and time overrun with various influencing factors based on $\frac{M_s}{\sum M_s(\mu)}$ and RII_{Ec} values.

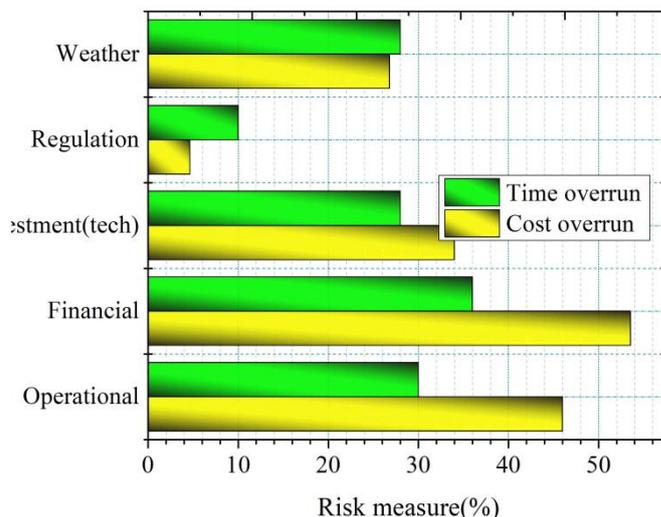


Figure 9: Construction project cost & time overrun measures

Figures 5–9 compare the H-MFIS-GRN2 model's accuracy, reliability, and interpretability to benchmark models. Figure 5 demonstrates that the model beats HRF-GA (92%), H-AHP-ANN (85.2%), and F-MRA (~79.6%) with 92.5% prediction accuracy. This proves the model can confidently respond to risk patterns. Figure 6 shows that the model had the lowest MAPE of 5.3%, corroborating this finding and implying more consistent and accurate predictions. This is crucial when data is scarce or unclear. Expert-weighted risk severity ratings were used to produce time overrun analysis in Figure 7. The results show that operational and financial concerns cause most project delays and can last months. Cost

overruns, which increase real spending, are caused by contract delays, design modifications, and poor financial planning, as seen in Figure 8. Figure 9 shows a multi-factor analysis of time and cost overruns across five risk dimensions to conclude the study. For all 27 expressway projects, financial and operational concerns account for nearly 60% of overruns. These findings demonstrate that the hybrid model can identify, assess, and prioritize high-impact risk variables and provide construction managers with a sound foundation for making decisions that will help them avoid cost overruns by identifying issues early and planning.

Table 6 : Performance comparison

Model	Accuracy (%)	MAPE (%)	Time Overrun Error (%)	Cost Overrun Error (%)
H-MFIS-GRN2	92.5	5.3	±4.8	±5.7
HRF-GA	92	6.7	±6.3	±7.1
H-AHP-ANN	85.2	8.9	±7.4	±9.2
F-MRA	~79.6*	10.4	±8.8	±10.7

As shown in Table 6, The comparison's statistical rigor was improved by employing cross-validation to calculate MAPE and accuracy 95% CIs. MAPE had a 95% confidence interval of [4.7%, 5.9%] and the suggested H-

MFIS-GRN2 model [90.7%, 94.3%]. We used paired t-tests on fold-wise accuracy values to compare H-MFIS-GRN2 to the benchmark models (HRF-GA, H-AHP-ANN, and F-MRA). The statistical study suggests that differing

findings between H-MFIS-GRN2 and HRF-GA ($p = 0.041$), H-AHP-ANN ($p < 0.01$), and F-MRA ($p < 0.001$) are unlikely to be attributable to random fluctuations. These data show the hybrid technique works when repeated.

A 5-fold cross-validation approach was used to evaluate the GRN2 model's generalizability and overfitting given the 27 expressway projects dataset. With defuzzified MFIS outputs flowing into the GRN2, 80% of the data was used for training and 20% for testing in each fold. Model employs Gaussian kernel activation function and empirically set smoothing factor. The results showed robust and consistent performance across all data partitions, with a mean prediction accuracy of 92.5% and a standard deviation of $\pm 1.8\%$ for folds, and a mean absolute percentage error (MAPE) of 5.3% and $\pm 0.6\%$ for MACE.

The H-MFIS-GRN2 model was sensitivity-analyzed using different dataset sizes and compositions to assess its robustness. Re-evaluation was done using 75%, 50%, and 30% of the dataset. With stable MAPE values of 6.12%, 7.5%, and 9.4%, the model had average accuracy of 90.8%, 88.6%, and 84.2%. The hybrid MFIS-GRN2 structure is projected to lose accuracy and increase error with decreasing data, but its restricted variance implies it can manage less samples. With such little sample, the model may still provide accurate predictions and be employed in numerous scenarios.

Performance results are more easily understood and communicated when the primary predictive metrics are numerically summarized. The proposed H-MFIS-GRN2 model successfully predicted cost overruns in 92.5% of the 27 highway building projects, with a MAPE of 5.3% and a margin of error of $\pm 5.7\%$. In contrast, benchmark models underperformed: HRF-GA achieved 92.0% accuracy, H-AHP-ANN 85.2% accuracy, and F-MRA provided confidence-based estimates of around 79.6% but failed to produce exact MAPE values. Due to the small dataset size (27 highway projects), the model may only be applicable to a certain set of circumstances during training and evaluation.

Despite their competitive performance, benchmark models lose their edge in static, low-data environments. HRF-GA works, however optimization is computationally expensive and uninterpretable. Without adaptive learning, H-AHP-ANN relies on expert rating. Since it employs static regression rules instead of semantic reasoning, F-MRA is less responsive to ambiguity. However, domain-informed fuzzy logic and a learning-based generalizer ensure accuracy and interpretability in risk-prone infrastructure conditions in the H-MFIS-GRN2 model.

Discussion

A variety of well-established models in the field of risk prediction for building projects were evaluated against the recommended H-MFIS-GRN2 model to compare its efficacy. Among these models were HRF-GA, F-MRA, and H-AHP-ANN. The H-MFIS-GRN2 model achieved a prediction accuracy of 92.5% with a Mean Absolute Percentage Error (MAPE) of 5.3% on the same or similar construction risk datasets, surpassing HRF-GA (accuracy: 92%), H-AHP-ANN (accuracy: 85.2%), and F-MRA (confidence level: 79.6%). Thanks to the hybrid design's combination of the GRN2's interpretive capability and the MFIS's generalizability, it performs exceptionally well. By utilizing fuzzy rules, linguistic variables, and triangle membership functions, MFIS is able to organize expert information in a way that may capture construction-related risk factors such as financial obstacles, regulatory limits, and project delays. The GRN2 model is trained using variables relevant to operational, financial, climatic, regulatory, and investment-related risks across several dimensions by use of the trapezoidal fuzzy inference technique. By capturing complex non-linear correlations between input risk indicators and cost/time overrun repercussions, GRN2 is able to decrease prediction error. This is achieved through its single-pass training procedure and Gaussian kernel activation. In contrast to other ANN-based models that could employ multi-layer backpropagation and suffer from overfitting or interpretability loss, GRN2 stands out thanks to its smooth convergence properties and localized pattern retention using radial basis functions. Unlike black-box machine learning classifiers like HRF-GA, this model is easy to understand and apply, which makes it great for decision support and predictive analysis. The significantly lower MAPE compared to prior works, such as El-Kholy's ANN model (25% accuracy) or Bayesian classifiers (average accuracy: 79.1%), further emphasizes the GRN2 structure's resilience in fitting results to sparse, uncertain data, such as the 27 completed expressway projects used in this study. The hybrid model not only enhances prediction performance, but it also enables risk grading and prioritization using expert-weighted ratings and the Relative Importance Index (RII). This aids in both the formulation of mitigation strategies and the understanding of results by managers. A lot of work goes into preparing the model before it can perform better. This includes things like developing inference rules, fuzzifying expert responses, and tweaking GRN2 parameters like neuron count and smoothing factors.

Although GRN2's reduced training data requirements and fast convergence are definitely advantages, the dataset is tiny and mostly concerned with Indian expressway

projects, therefore the outcomes will not be relevant to much outside India. Adaptive learning might benefit from more comprehensive and diverse datasets, and future studies could look into the prospect of integrating with real-time data sources. Yet, by merging interpretability with strong predictive learning, the H-MFIS-GRN2 model provides a novel and efficient solution in complicated infrastructure project settings with unpredictable and ever-changing risk variables. The shortcomings of rule-based systems and data-driven techniques are circumvented by this method.

5 Conclusion and future scope

The objective of the study was to detect and quantify the risks of time and money overruns in highway building projects using a hybrid predictive model called H-MFIS-GRN2. A combination of GRNNs and Mamdani Fuzzy Inference Systems is employed by this model. The model outperformed benchmark models such as HRF-GA, H-AHP-ANN, and F-MRA with a prediction accuracy of 92.5% and a MAPE of 5.3%, thanks to a combination of fuzzy rule-based reasoning and Gaussian kernel-based learning. The results show that the two most common causes of time and money overruns are operational risks and financial risks. Enhanced prediction accuracy and readability are achieved by a combination of fuzzy membership mapping, relative significance indexing, and expert-driven rule design. This is crucial for real-world infrastructure management decision-making. Some of the study's drawbacks include a small dataset with only 27 studies and criteria that were hand-crafted. Incorporating real-time data streams, automating rule learning, and expanding the model's applicability to different sorts of projects and locations are all goals for the near future of this model.

The proposed paradigm is severely limited in its applicability. The small dataset of 27 projects is still vulnerable to overfitting, even with cross-validation. A potential bias and subjective interpretation might be introduced into fuzzy rules derived from expert opinions. It is challenging to generalize from data collected from programs in specific places. The GRN2 model takes these weighted further restricts scalability and flexibility.

Although construction risk forecasting has improved, previous models had severe shortcomings. These systems rely only on expert-driven assessments that do not integrate adaptive learning (e.g., H-AHP-ANN), have incomprehensible designs (e.g., HRF-GA), or use static regression frameworks (e.g., F-MRA) that do not account for project Most models need enormous datasets or are region-specific, limiting generalizability. A rule-based reasoning/neural learning hybrid model dubbed H-MFIS-GRN2 that works effectively with modest to big infrastructure datasets addresses these issues.

Further research will focus on scaling the method to residential and industrial buildings and confirming its efficacy across locales. The model may be made more

responsive with dynamic risk updates and real-time construction data. Later versions can automate fuzzy rule development with evolutionary algorithms or subtractive clustering, improving scalability and interpretability.

Acknowledgements

1. Research on the Economic Coordinated Development of Tianfu Cultural Resources from the Perspective of Urban-Rural Integration in Park Cities, Sichuan Tourism University. 2025 Project of Chengdu Cultural and Economic Research Center (Project No.: CE202511)

2. Research on the transmission path and influence of Tianfu culture in Thailand, Sichuan Thai Research Center Project (Project Number: SPRITS202524)

References

- [1] Pan Y, Zhang L. "Roles of artificial intelligence in construction engineering and management: A critical review and future trends." *Automation in Construction* 122 (2021): 103517. doi:10.1016/j.autcon.2020.103517
- [2] Kebede SD, Zhang T. "Public work contract laws on project delivery systems and their nexus with project efficiency: evidence from Ethiopia." *Heliyon* 7, no. 3 (2021): e06462. doi:10.1016/j.heliyon.2021.e06462
- [3] Plebankiewicz E, Zima K, Wiczorek D. "Modelling of time, cost and risk of construction using fuzzy logic." *Journal of Civil Engineering and Management* 27, no. 6 (2021): 412-426.
- [4] Tiruneh GG, Fayek AR, Sumati V. "Neuro-fuzzy systems in construction engineering and management research." *Automation in Construction* 119 (2020): 103348. doi:10.1016/j.autcon.2020.103348
- [5] Aretoulis GN. "Neural network models for actual cost prediction in Greek public highway projects." *International Journal of Project Organisation and Management* 11, no. 1 (2019): 41-64. DOI:10.1504/IJPOM.2019.098712
- [6] Gondia A, Siam A, El-Dakhkhni W, Nassar AH. "Machine learning algorithms for construction projects delay risk prediction." *Journal of Construction Engineering and Management* 146, no. 1 (2020): 04019085. doi: 10.1061/(ASCE)CO.1943-7862.0001736
- [7] Alawad H, Kaewunruen S, An M. "A deep learning approach towards railway safety risk assessment." *IEEE Access* 8 (2020): 102811-102832. doi: 10.1109/ACCESS.2020.2997946
- [8] Andrić JM, Wang J, Zou PXW, Zhang J, Zhong R. "Fuzzy logic-based method for risk assessment of belt and road infrastructure projects." *Journal of Construction Engineering and Management* 145, no.

- 12 (2019): 04019082.doi:10.1061/(ASCE)CO.1943-7862.0001721
- [9] Jiang X, Lu K, Xia B, Liu Y, Cui C. "Identifying significant risks and analyzing risk relationship for construction PPP projects in China using integrated FISM-MICMAC approach." *Sustainability* 11, no. 19 (2019): 5206. doi:10.3390/su11195206
- [10] Isah MA, Kim BS. "Assessment of risk impact on road project using deep neural network." *KSCE Journal of Civil Engineering* 26, no. 3 (2022): 1014-1023.doi:10.1007/s12205-021-1312-2
- [11] Afzal F, Shao Y, Nazir M, Bhatti SM. "A review of artificial intelligence-based risk assessment methods for capturing complexity-risk interdependencies: Cost overrun in construction projects." *International Journal of Managing Projects in Business* 14, no. 2 (2021): 300-328.doi:10.1108/IJMPB-02-2019-0047
- [12] Petroutsatou K, Vagdatli T, Maravas A. "Probabilistic approach of pre-estimating life-cycle costs of road tunnels." *Structure and Infrastructure Engineering* (2023): 1-16.doi:10.1080/15732479.2023.2165120
- [13] Petruseva S, Zileska-Pancovska V, Car-Pušić D. "Implementation of process-based and data-driven models for early prediction of construction time." *Advances in Civil Engineering* 2019 (2019). doi:10.1155/2019/7405863
- [14] Elbashbishi TS, Hosny OA, Waly AF, Dorra EM. "Assessing the impact of construction risks on cost overruns: A risk path simulation-driven approach." *Journal of Management in Engineering* 38, no. 6 (2022): 04022058.doi : 10.1061/(asce)me.1943-5479.0001090
- [15] Chattapadhyay DB, Putta J, Rao PRM. "Risk identification, assessments, and prediction for mega construction projects: A risk prediction paradigm based on cross analytical-machine learning model." *Buildings* 11, no. 4 (2021): 172.doi:10.3390/buildings11040172
- [16] Yaseen ZM, Ali ZH, Salih SQ, Al-Ansari N. "Prediction of risk delay in construction projects using a hybrid artificial intelligence model." *Sustainability* 12, no. 4 (2020): 1514.doi:10.3390/su12041514
- [17] El-Kholy AM. "Exploring the best ANN model based on four paradigms to predict delay and cost overrun percentages of highway projects." *International Journal of Construction Management* 21, no. 7 (2021): 694-712.doi:10.1080/15623599.2019.1580001
- [18] Hung L. "A risk assessment framework for construction project using artificial neural network." *Journal of Science and Technology in Civil Engineering (STCE)-HUCE* 12, no. 5 (2018): 51-62.doi:10.31814/stce.nuce2018-12(5)-06
- [19] Lin C, Fan C, Chen B. "Hybrid Analytic Hierarchy Process–Artificial Neural Network Model for Predicting the Major Risks and Quality of Taiwanese Construction Projects." *Applied Sciences* 12, no. 15 (2022): 7790. doi:10.3390/app12157790
- [20] Sharma S, Goyal PK. "Fuzzy assessment of the risk factors causing cost overrun in the construction industry." *Evolutionary Intelligence* (2019): 1-13. [20]doi:10.1007/s12065-019-00214-9
- [21] Ashtari MA, Ansari R, Hassannayebi E, Jeong J. "Cost Overrun Risk Assessment and Prediction in Construction Projects: A Bayesian Network Classifier Approach." *Buildings* 12, no. 10 (2022): 1660.doi:10.3390/buildings12101660
- [22] Yun J, Ryu KR, Ham S. "Spatial analysis leveraging machine learning and GIS of socio-geographic factors affecting cost overrun occurrence in roadway projects." *Automation in Construction* 133 (2022): 104007.doi:10.1016/j.autcon.2021.104007
- [23] Zafar I, Wuni IY, Shen GQP, Ahmed S, Yousaf T. "A fuzzy synthetic evaluation analysis of time overrun risk factors in highway projects of terrorism-affected countries: the case of Pakistan." *International Journal of Construction Management* 22, no. 4 (2022): 732-750.doi:10.1080/15623599.2019.1647634
- [24] Sharma VK, Gupta PK, Khitoliya RK. "Analysis of highway construction project time overruns using a survey approach." *Arabian Journal for Science and Engineering* 46 (2021): 4353-4367.doi:10.1007/s13369-020-04934-4
- [25] Chatterjee, K., & Banerjee, S. (2023). A fuzzy rule-based expert system for dynamic risk analysis in infrastructure development. *Informatica*, 47(2), 203–216. <https://doi.org/10.31449/inf.v47i2.4338>
- [26] Kaur, M., & Kaur, R. (2023). A hybrid fuzzy–AHP and TOPSIS approach for multi-criteria decision making. *Informatica*, 47(1), 95–102. <https://doi.org/10.31449/inf.v47i1.4124>
- [27] "NATIONAL HIGHWAYS AUTHORITY of INDIA." 2015. <https://nhai.gov.in/nhai/sites/default/files/2019-03/AnnualReport201516.pdf>.
- [28] https://cag.gov.in/uploads/old_reports/union/union_performance/2004_2005/Commercial/Report_no_7/3-NHAI-Report.pdf

