

Software

CATD: a reproducible pipeline for selecting cell-type deconvolution methods across tissues

Anna Vathrakokoili Pournara ^{1,*}, Zhichao Miao ^{1,2,3,*}, Ozgur Yilmaz Beker ^{1,4},
Nadja Nolte ^{1,5}, Alvis Brazma¹, Irene Papatheodorou ^{1,2,6,*}

¹European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SD, United Kingdom

²Open Targets, Wellcome Genome Campus, Hinxton CB10 1SD, United Kingdom

³GMU-GIBH Joint School of Life Sciences, Guangzhou Laboratory, Guangzhou Medical University, Guangzhou, 511436, China

⁴Faculty of Engineering and Natural Sciences, Sabanci University, Tuzla 34956, Turkey

⁵Department of Biotechnology and Systems Biology, National Institute of Biology, Ljubljana, 121-1000, Slovenia

⁶Earlham Institute, Norwich Research Park, Norwich NR4 7UZ, United Kingdom

*Corresponding author. European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SD, United Kingdom. E-mail: annavp@ebi.ac.uk (A.V.P.); GMU-GIBH Joint School of Life Sciences, Guangzhou Laboratory, Guangzhou Medical University, Guangzhou, 511436, China. E-mail: miao_zhichao@gzlab.ac.cn (Z.M.); Earlham Institute, Norwich Research Park, Norwich NR47UZ, United Kingdom. E-mail: irene.papatheodorou@earlham.ac.uk (I.P.)

Associate Editor: Dr. Marieke Kuijjer

Abstract

Motivation: Cell-type deconvolution methods aim to infer cell composition from bulk transcriptomic data. The proliferation of developed methods coupled with inconsistent results obtained in many cases, highlights the pressing need for guidance in the selection of appropriate methods. Additionally, the growing accessibility of single-cell RNA sequencing datasets, often accompanied by bulk expression from related samples enable the benchmark of existing methods.

Results: In this study, we conduct a comprehensive assessment of 31 methods, utilizing single-cell RNA-sequencing data from diverse human and mouse tissues. Employing various simulation scenarios, we reveal the efficacy of regression-based deconvolution methods, highlighting their sensitivity to reference choices. We investigate the impact of bulk-reference differences, incorporating variables such as sample, study and technology. We provide validation using a gold standard dataset from mononuclear cells and suggest a consensus prediction of proportions when ground truth is not available. We validated the consensus method on data from the stomach and studied its spillover effect. Importantly, we propose the use of the critical assessment of transcriptomic deconvolution (CATD) pipeline which encompasses functionalities for generating references and pseudo-bulks and running implemented deconvolution methods. CATD streamlines simultaneous deconvolution of numerous bulk samples, providing a practical solution for speeding up the evaluation of newly developed methods.

Availability and implementation: https://github.com/Papatheodorou-Group/CATD_snakemake.

1 Introduction

There is a growing interest in understanding the level of heterogeneity and the importance of cell-type abundances in healthy and diseased tissues. Transcriptomic heterogeneity and cell-type composition reveal critical features of tissue functionality. For example, in cancer, immune cells can either be recruited in the vicinity of the tumor to contribute to cancer progression, or they can play a tumor-suppressive role by recognizing and killing abnormal cells (Hanahan and Coussens 2012, Dumont *et al.* 2013, Taube *et al.* 2018, Jorge *et al.* 2020). High infiltration from specific immune cells can be predictive of the disease outcome, stage and aid the treatment selection (Zhang *et al.* 2013). However, determining the composition of a tissue sample can present significant challenges. While single-cell RNA sequencing techniques offer valuable insights into tissue heterogeneity, challenges such as technical handling, sorting of

specific cell types, sub-optimal dissociation and tissue vulnerability may affect the accurate determination of cell type proportions in single-cell data (Denisenko *et al.* 2020). Furthermore, conventional methods for quantifying cell abundances, such as flow cytometry and immunohistochemistry, rely on prior knowledge of cell markers, depend on interpretation bias and are not easily scalable (Matos *et al.* 2010, Patrick *et al.* 2020). Additionally, there are only a handful of this type of datasets available and mostly deriving from non-solid tissues such as blood (Finotello *et al.* 2019, Monaco *et al.* 2019). Taking all the above into account, computational deconvolution methods show a great advantage, since cell proportions from hundreds of samples can be obtained computationally very fast—in a matter of seconds to a few hours depending on the selected method. As a result, cell-type deconvolution methods—also called decomposition methods—provide a cost-effective way of measuring cell proportions when compared to experimental

methods. Moreover, deconvolution methods can be very useful for inferring cell type compositions from historical bulk RNA-seq data from large databases such as GTEx (GTEx Consortium 2017), TCGA (The Genome Atlas Cancer Program) (The International Cancer Genome Consortium, 2010), TARGET (<https://ocg.cancer.gov/programs/target>) and GEO (Edgar *et al.* 2002) in which it would be impossible to repeat experiments on.

Cell-type deconvolution can contribute significantly in answering critical biological questions and understanding disease mechanisms to a greater extent. Specifically, it can unravel how cell proportions can affect certain phenotypes or specific clinical features available within genomic databases (Donovan *et al.* 2020). Additionally, it can serve as a powerful tool for mitigating the confounding effect of cell-type proportions in bulk differential expression analysis and as a result help in the identification of reliable biomarkers for diseases and perturbed processes and pathways with cell-type resolution (Inkeles *et al.* 2016). Moreover, in the context of transcriptome-wide association studies deconvolution can rectify cell abundance bias in the expression quantitative trait loci (eQTL) analysis and facilitate the identification of cell-type specific eQTLs (Lowe and Rakyen 2014). Finally, the abundance of bulk RNA-seq samples present a unique opportunity for training machine learning models tailored to phenotype prediction, a task that is often challenging with single-cell data due to limited sample sizes. Importantly, by integrating cell composition, we effectively enrich the information content that bulk samples are missing due to the experimental method principles.

Deconvolution methodologies can be formally categorized into three main groups based on their input requirements: supervised (reference-based), semi-supervised, and unsupervised (complete). Supervised deconvolution methods necessitate an expression matrix (typically from bulk RNA-sequencing or microarray data) to be deconvolved as well as cell-type-specific information. This specific information can be derived from sources like flow cytometry, single-cell RNA-seq data, or marker gene lists. The above supervised methods can further be divided into two subcategories: bulk when the reference is either a signature of sorted cell types or a marker gene list and single-cell when the reference is a single-cell dataset. Well-known ‘bulk’ deconvolution methods include CIBERSORT (absolute) (Newman *et al.* 2015), FARDEEP (Hao *et al.* 2019), NNLS (non-negative least squares) and EpiDISH (Teschendorff *et al.* 2017). Concurrently, there are contemporaneous deconvolution methods that harness available single-cell datasets to extract cell type-specific features and subsequently perform deconvolution. Examples of such methods include Bseq-sc (Baron *et al.* 2016), DWLS (Tsoucas *et al.* 2019) and MuSiC (Wang *et al.* 2019) and others (Frishberg *et al.* 2019, Jew *et al.* 2020, Dong *et al.* 2021). These approaches incorporate feature selection methods to choose informative genes for deconvolution, such as DWLS-MAST, MuSiC-informative genes, and AutogeneS. All the above methods output cell proportions, enabling both intra-sample and inter-sample comparisons. In the supervised category, gene set enrichment methods such as xCell (Aran *et al.* 2017), MCP-counter (Becht *et al.* 2016), and SaVant (Lopez *et al.* 2017) can also be included. However, these methods use an enrichment-based approach and output enrichment scores, not percentages. These scores can be compared across

samples, but do not allow intra-sample comparison. Unsupervised methods, on the other hand, such as deconf (Repsilber *et al.* 2010) and CDSeq (Kang *et al.* 2019) require only a bulk sample and the number of components (representing the cell types to be identified) as input. These methods estimate both the cell-type-specific expression and cell fractions simultaneously. Semi-supervised methods, such as ssKL and ssFrobenius, can also estimate both the cell-specific expression matrix and the matrix of cell proportions at the same time. They utilize a list of pre-determined markers, the expected number of components, and the bulk matrix as input. These semi-supervised methods employ a block-descent approach to estimate both matrices (Gaujoux and Seoighe 2012). In this study, our primary focus is on comparing supervised and semi-supervised methods, and we do not extensively explore unsupervised (complete) deconvolution methods. Unsupervised methods are particularly valuable when the cell composition is entirely unknown or cell signatures for the cell types of interest are not available. These methods have been thoroughly investigated in previous benchmarks to assess their performance and underlie their potential (Jaakkola and Elo 2021, Jin and Liu 2021, Sutton *et al.* 2022). Lastly, there is a separate class of methods known as ‘expression deconvolution methods’. These methods tackle the inverse problem; they aim to find the cell type-specific expression matrix given prior knowledge of cell proportions and bulk expression data. Methods falling under this category include Rodeo (Jaakkola and Elo 2021), csSAM (Shen-Orr *et al.* 2010), Deblender (Dimitrakopoulou *et al.* 2018), cs-lsfite (Gaujoux and Seoighe 2012), and cs-prog (Gong *et al.* 2011). Previous benchmark studies have evaluated the performance of these methods, but due to their distinct design and research focus, we have not included those in our study.

The field of deconvolution methods has witnessed significant growth in recent years, driven by the complex nature of the deconvolution problem and the ongoing development of new and more sophisticated techniques. Over the past decade, there have been notable efforts to independently evaluate deconvolution methods (Avila Cobos *et al.* 2020, Jin and Liu 2021). These efforts explore various aspects, such as the effectiveness of simulations, method performance, the importance of reference data, the impact of data normalization, the influence of different technologies, and the granularity of cell-type information in the data (Teschendorff *et al.* 2017, Jiménez-Sánchez *et al.* 2019, Sturm *et al.* 2019, Nadel *et al.* 2021, Sutton *et al.* 2022, White *et al.* 2022, Alonso-Moreda *et al.* 2023, Cobos *et al.* 2023). However, as our knowledge in this field continues to expand, and new methods are introduced each year, there is a growing need for comprehensive, reproducible, and versatile benchmarking pipelines that can accommodate both existing and newly developed methods (Garmire *et al.* 2024). This need is especially pressing because different deconvolution methods may have varying input formats, design methodologies, preprocessing steps, and output results. Therefore, a centralized and efficient platform is essential to streamline the evaluation process, facilitating a quicker and more unified assessment of deconvolution methods. Furthermore, while previous studies have primarily focused on deconvolution in peripheral blood mononuclear cell (PBMC) samples or in a limited number of tissues each time, there is a growing demand for research that covers more

tissues and datasets as more and more cell type-specific information becomes available (Maden *et al.* 2023).

To address these challenges, we have designed a robust and reproducible Snakemake pipeline for evaluating 31 available deconvolution methods and obtaining results on cell abundances from real data. Our evaluation involves nineteen single-cell datasets from diverse sources, covering five different tissues with multiple datasets per tissue, enabling us to generate a wide range of simulated and reference data for method evaluation. Our pipeline has also incorporated various simulation approaches, leveraging single-cell datasets and suggesting parameter settings to generate realistic simulated data across different tissue types. Additionally, we study the impact of preprocessing methods and the effect of distinct technologies in the reference data. Deconvolution performance is finally evaluated on real data from PBMC samples with ground truth. Moreover, we introduced a consensus deconvolution approach by combining well-performing methods and validated the results on stomach samples from GTEx using abundances from tissue slides as our ground truth. Finally, we assess the spillover effect of the consensus method on purified samples from cerebral cortex and lung tissue. One notable feature of our pipeline is its adaptability to newly developed deconvolution methods that require independent benchmarking. Furthermore, researchers can employ our pipeline to simultaneously deconvolve multiple biological samples, obtaining results from implemented methods and benefiting from the pipeline's visualizations. By assessing the consistency of results across methods and references, researchers can subsequently choose the most suitable setup to address their specific research questions.

2 Methods

2.1 Data collection for benchmarking

We collected single-cell data from various sources to generate references and simulate bulk data for in-silico experiments. 12 single-cell datasets were collected from human and mouse brain tissues, including dental gyrus, anterolateral motor cortex (ALM), primary visual cortex (ViSP), cerebral cortex and pons, retinas, hippocampus CA1, somatosensory cortex S1, primary motor cortex, cortex areas and glioblastoma cancer tissue. We also collected two single-cell datasets specifically focused on human pancreatic islets to investigate sample effects in the deconvolution process. To understand the impact of different platforms/single-cell technologies in deconvolution we curated three pairs of 10X and Smart-seq2 datasets, with each pair originating from the same study. For this task we covered three tissues: placenta, brain and lung. A comprehensive breakdown of the datasets is provided in [Supplementary Table S1](#), offering details on the number of samples, cells, genes and cell-types per dataset, covered conditions, and corresponding GEO accession numbers.

2.2 Single-cell data (re)analysis

Raw counts and metadata (minimum metadata required: cell-type annotations, cellIDs, sampleIDs) were obtained for each single-cell dataset used for the in-silico benchmarking experiments. Single-cell datasets used for reference generation were re-analysed utilizing the SCANPY toolkit (Wolf *et al.* 2018). Cells with less than 200 genes and genes expressed in less than 3 cells were filtered out. Cells with more than >5% of mitochondrial genes or too many total counts (dataset-specific threshold) were

also removed. To inspect the cell annotation quality all datasets were normalized (total counts = 10e4) and log-transformed, the raw count dimension was kept intact to be used as a deconvolution input in the pipeline. Next, highly variable genes were identified and we performed principal component analysis (PCA), computed the neighborhood graph as well as performed leiden clustering using different resolution parameters (0.25, 0.5, 1). After clustering, we plotted the data on UMAP space using as color key: cellType and leiden clusters in order to evaluate the author's annotations. Cell types annotated from the authors in the original papers as 'others', 'unknown', 'low quality', 'unclassified' were removed from the data. Additionally, cell type nomenclature was harmonized for consistency across datasets. This step is particularly useful for the cross-reference and real bulk deconvolution tasks. To validate cell type annotations, cell types were examined using known marker genes. This step involved visualizing the expression of genes known to be specific to particular cell types on UMAP space.

2.3 Real data validation

To validate the results from benchmarking on real data, we collected real bulk samples from various tissues and designed validation experiments using diverse ground truth data. Collection, curation and design of validation experiments are described below. More details can be found in [Supplementary Table S4](#).

2.3.1 Validation on human PBMC samples

We collected bulk RNA-seq data from 9 PBMC human samples. For these samples we have also coupled flow cytometry measurements for the cell type proportions (GSE107572) (Finotello *et al.* 2019). To construct the references for deconvolution we utilized a publicly available single-cell dataset from 14 PBMC samples (GSE150728) (Wilk *et al.* 2020). We harmonized the cell type annotation between the single-cell and the bulk data, ensuring seamless deconvolution analysis and evaluation ([Supplementary Fig. S8](#)).

2.3.2 Validations on GTEx stomach data

We acquired bulk expression data (TPM-normalized) from GTEx portal v8 and focused on the stomach tissue data. Among the 939 stomach samples, 339 were accompanied by histology slides, enabling the calculation of ground truth proportions through QuPath (Bankhead *et al.* 2017) analysis (see Methods 2.12). To use a reliable reference for this task we used the stomach subset of the human cell landscapes single-cell dataset, downloaded from cellxgene portal (CZI Single-Cell Biology Program 2023). For a standardized reduction of cell type label granularity we employed an automated ontology-based method using the scOntoMatch (v 0.1.0) package (Song *et al.* 2023).

The adapted code for this task is available at: https://github.com/nadjano/reference_preparation_for_GTEX_deconvolution.

2.3.3 Deconvolving purified flow cytometry data

To study the spillover effect previously documented in deconvolution we collected and deconvolved bulk samples of purified cell types from two tissue sources, human fetal lung (E-MTAB-9372) and mouse brain (E-GEOD-52564). The data were downloaded from Expression Atlas (Moreno *et al.* 2022). As a lung specific deconvolution reference, we used the subset from the GTEx v7 single-cell dataset. Cell type labels granularity was reduced as described in the previous section. For the deconvolution of mouse brain data, Tabula Muris dataset (The Tabula Sapiens Consortium* 2022) (E-

ENAD-15) was utilized. To expedite the deconvolution computation we downsampled the above single-cell references. Specifically, for cell types that possessed more than 300 cells, we employed a random sampling approach, selecting 300 cells per cell type. The detailed code implementation for this procedure can be found here: (https://github.com/nadjano/reference_preparation_for_GTEX_deconvolution).

2.4 Overview of CATD Snakemake pipeline

The CATD pipeline is a benchmarking pipeline meant to facilitate the assessment of cell-type deconvolution methods (currently 31) across different simulation scenarios in a standardized way. It also allows the deconvolution of real bulk samples with various input parameters allowing users to deconvolute their own in-house data following our proposed guidelines. The pipeline includes:

- Pseudo-bulk generation methods (from single-cell RNA-seq data) that allow to create diverse pseudo-bulk samples and compare deconvolution methods across different scenarios.
- 17 normalization methods implemented in the pipeline for the normalization of the input single-cell reference and the (pseudo) bulk samples.
- 4 transformation methods.
- 9 Differential Expression tests for the selection of marker genes from single-cell reference data (Seurat, FindMarkers).
- 31 deconvolution methods.
- 9 metrics to assess the results when we test deconvolution methods on pseudo-bulks or when ground truth proportions from real data are available.
- A Consensus approach, combining three well-performing deconvolution methods.

Moreover the pipeline provides:

- Visualization of the performance and scalability metrics across methods.
- Visualization of the consensus deconvolution results (bar plots with cell type proportions per sample).
- Visualization of the correlation of the results across methods for a given dataset when ground truth is not available.

2.5 Generation of pseudo-bulk (synthetic) samples

To benchmark deconvolution methods, we designed various pseudo-bulk simulation data using different sampling methodologies. Each simulation starts with splitting the single cell in half (50% testing data and 50% training data). Testing data will be then used to generate the pseudo-bulks. All the simulations require as an input to set the number of samples (n) to be generated and the number of cells to sample from the single cell for each pseudo-bulk (c). Two modes of pseudo-bulk generation have been designed for the purpose of our benchmark. For our simplest simulation (mode 1), the sampling will randomly pick c number of cells with replacement from the single-cell data to generate the first sample, this will be repeated n times until we obtain the final pseudo-bulk dataset. For mode 2, we have designed a more realistic pseudo-bulk generation method which takes into account variability of cell type proportions and cell type heterogeneity for each sample. Hence, for mode 2 the generation of pseudo-bulks will require two additional parameters, namely

‘proportional variance’ (p) and second, a logical parameter called ‘sampleCT’ that defines whether or not some cell types will be left out during generation. When including all cell types in our pseudo-bulks, if p is set to a negative number, then random proportions ranging between 0.01 and 0.99 and summing to 1, will be generated using a uniform distribution. If p is positive, then integers between 1 and $1 + 2p$ will be generated using a uniform distribution and then normalized to sum to 1. This enables more precise control of the noise, with the former (1–99 case) being the edge case of the latter (positive p) when p tends to infinity. If cell types are sampled, then a bimodal normal distribution with uniform mode probability is used to achieve the effect and the p instead becomes the standard distribution of the kept cell types. In this case, the provided p should be positive and ideally large. The absolute value of generated numbers is considered to ensure non-negative proportions and then proportions are normalized to sum to one.

The generated proportions in mode 2 are then scaled by c to determine how many cells will be chosen with replacement per cell type from the reference data to generate one sample, which is again repeated n times.

2.6 Building reference data from single-cell data

To build references used by the deconvolution methods we start by obtaining a single cell RNA-sequencing expression matrix. Different methods require different types of reference input.

- First, we have ‘single-cell’ methods that require a single-cell matrix as input. Within the pipeline, this matrix is encoded as *C0*. *C0* is a generic single-cell matrix, with rows being genes and columns being cellIDs, this matrix is used as an input in single-cell deconvolution methods (e.g. MuSiC, DWLS, Bisque and others). Single-cell methods also need an additional metadata object encoded as *phenData* in our pipeline which is associated with *C0* and provides additional information (sampleIDs, cellIDs, cell type labels).
- Next, we have reference-based deconvolution methods that require input expression data from purified cell types (e.g. FACS data). For these methods we generate the *C1 matrix* which is the averaged (arithmetic mean) expression matrix by cell types, with rows being genes and columns being cell types.
- For marker-based methods, we create the *C2 reference*. *C2* is a list of marker genes obtained from differential gene expression analysis on the single-cell data through the Seurat package. Specifically a combination of *NormalizeData()* and *FindAllMarkers()* functions are being used. Methods that use *C2* reference as input include: EPIC, DSA, ssFrobenious, ssKL, deCAMmarker. By default, we use the ‘wilcox’ test (Wilcoxon Rank Sum test) to identify marker genes. Moreover, we have implemented in the pipeline an additional functionality to use two tests each time and obtain the overlap of marker genes from the two tests. However, this could result in issues while running the pipeline in cases where the overlap of genes is very low.
- Finally, *refVar* is the row standard deviations matrix by cell types, similar in structure to *C1*. *refVar* is used only by EPIC deconvolution method due to its unique design.

2.7 Preprocessing methods

To comprehensively examine the impact of variable preprocessing on both (pseudo)bulk and single-cell matrices, we employed an array of normalization and transformation methods. While many of the normalization techniques applied in this study are suitable for both data types, we also incorporated four methods tailored specifically for single-cell data: SCTransform, scran, scatter, and Linnorm. A detailed summary of all the methods, including their characteristics, can be found in [Supplementary Table S2](#).

Furthermore, we investigated the data transformation effects. To explore this aspect, we implemented four distinct transformations: logarithmic (log), square (sqrt), variance-stabilizing transformation (vst), or no transformation (linear scale). Finally, we systematically examined the impact of the order of preprocessing steps, considering whether transformation or normalization should be conducted first. Detailed insights into the implementation of all preprocessing steps and methods are available within the pipeline.

2.8 Selection of deconvolution methods

In our benchmarking study, we curated deconvolution methods from existing literature and previous benchmark studies ([Avila Cobos et al. 2020](#), [Jaakkola and Elo 2021](#), [Jin and Liu 2021](#), [Sutton et al. 2022](#)). The selected methods represent diverse categories, including 27 supervised methods, which further break down into 10 single-cell-based methods, 12 reference-based methods, and 4 marker-based methods. Additionally, we incorporated 2 semi-supervised methods (ssKL and ssFrobenious) and 2 unsupervised methods (deconf and CDSeq).

A comprehensive summary of the enlisted methods is available in [Supplementary Table S3](#), providing details on their characteristics, algorithmic principles, implementation specifics, tool version numbers, and original publications.

2.9 Metrics for performance evaluation

To assess the efficacy of deconvolution methods, we conducted a thorough evaluation by comparing the proportions derived by these methods against known ground truth proportions. Depending on the specific task, the ground truth can either be experimentally calculated proportions or proportions derived from simulations. Our pipeline utilizes key metrics to evaluate performance, including Pearson correlation coefficients, Spearman correlation coefficients, root-mean-square error (RMSE), and weighted RMSE. For weighted RMSE each cell type is assigned a weight based on its proportion value in a way where rare cell types are given bigger weight. Additionally, we incorporated [supplementary metrics](#) such as mean absolute error/mean absolute deviation, Euclidean distance, distance correlation, cosine similarity (cos), and R-squared. The notation and equations for calculating these metrics are detailed in the [Supplementary material](#). All metrics were computed in R, and the corresponding functions can be accessed at: https://github.com/Functional-Genomics/CATD_snake_make/blob/main/Modules/Res_explore/Res_explore.R.

2.10 Metrics for scalability evaluation

In the context of deconvolution, assessing scalability is crucial, particularly when considering the application of these methods to large-scale datasets such as extensive databases with historical bulk data. Evaluating scalability ensures that the computational efficiency of deconvolution can meet the

demands of processing substantial volumes of data. Every step of the pipeline is benchmarked automatically by Snakemake's built-in 'benchmark' directive, and the results from these benchmarks are also visualized at the end of the pipeline. Four of the parameters are visualized for initial scalability assessment purposes, namely the time (s), memory (maximum resident set size), mean load (CPU usage divided by total time), and CPU time.

2.11 Consensus approach

For the consensus approach we selected three deconvolution methods: DWLS, EpiDISH, FARDEEP. These methods were performing well across our study. The main criteria for inclusion were:

- 1) well-performing methods overall,
- 2) open-source code,
- 3) reasonable time of running,
- 4) stable and they produce results (sometimes methods will not converge and produce NAs and
- 5) flexible (e.g. they can take as input different references every time to answer biological questions).

The displayed cell-type proportions of GTEx stomach, mouse purified brain, and human fetal lung cell types are computed as the averages of deconvolution outcomes obtained through DWLS, EpiDISH and FARDEEP, per sample and cell type. We also include the mean correlation between the results from the three methods across all samples to show the concordance of results across the methods. Standard deviation (SD) per cell type for each sample across methods is also computed. For the consensus visualization of deconvolution proportions in GTEx stomach samples, we employed the `CiberBarFrazer()` function. Code available at https://github.com/mkrdonovan/gtex_deconvolution, R version 4.3.1). The resulting proportions were arranged in descending order of the most abundant cell types.

2.12 QuPath analysis

To validate our deconvolution results for selected stomach samples, we leveraged histology images downloaded from the GTEx Portal Histology Viewer (<https://gtexportal.org/home/histologyPage>). Histology files for 339 samples were downloaded, and for analysis, we employed QuPath software (v 0.4.1 on Mac OS X - 11.6.2). A pixel classifier was trained on six slides, allowing us to estimate the areas of background, muscle, mucosa, and submucosal layers. Subsequently, using the assigned areas, we calculated the relative proportions of the three tissue layers on all slides.

For the mapping of relevant cell types to stomach layers, we drew upon insights from the literature on the organizational structure of stomach tissues. Specifically, the mapping was done as follows:

- Mucosa: Endocrine cells, enterocytes, goblet cells, peptic cells (parietal cells), and epithelial cells.
- Submucosa: Dendritic cells, endothelial cells, erythroid lineage cells, fibroblasts, macrophages, plasma cells, and stromal cells.
- Muscle: Mesenchymal stem cells, neurons, and smooth muscle cells.

performance of the 31 collected deconvolution methods (Supplementary Table S3). To compare the performance of deconvolution methods we implemented various evaluation metrics in the CATD pipeline. Pearson correlation coefficient (r /pcor), root mean square error (RMSE), average cosine similarity (avgcoss), weighted RMSE and Spearman's correlation are some of the key metrics implemented and used across this study to evaluate the performance of the methods in different scenarios (Methods). We select the best performing methods taking into consideration the methods that have both low RMSE and high R values consistently across different simulation scenarios, datasets and tissues. All the steps of the benchmarking study have been modularized in a snakemake pipeline for systematic exploration and reproducibility of the benchmarking results. The pipeline also allows the evaluation of the scalability of each benchmarking module using additional metrics (Methods). Finally, CATD provides a practical, user-friendly framework to accurately deconvolve real bulk samples with or without ground truth by suggesting a new approach to obtain a fair consensus of the best-performing methods and providing visualization of the deconvolution results.

3.2 Realistic pseudo-bulk simulations for systematic evaluation and method development

3.2.1 Objective

Efforts to benchmark and evaluate the performance of various methods in deconvolution often necessitate the availability of sufficient ground truth data for comparison. Nevertheless, obtaining real bulk proportions experimentally can be challenging and constantly. Moreover, existing ground truth data tends to be limited in size and primarily centered around non-solid tissues. To facilitate method comparisons and understand their limitations, the generation of pseudo-bulk data proved to be a valuable approach. It should be noted that creating realistic simulations in deconvolution studies is essential for applying them confidently afterwards to real data. A common method to simulate bulk profiles involves utilizing publicly available single-cell data and generating pseudo-bulk by averaging the expression of single cells. Simulation strategies employed in previous studies range from random sampling to more sophisticated approaches (Dietrich *et al.* 2022, Hu and Chikina 2023). Nevertheless, it is non-trivial to understand how the different parameters used in the simulations can affect the deconvolution results. In response, we have designed and compared various simulation approaches, introducing standardized parameters for simulating bulk expression profiles that can be used for the benchmarking of existing, and development of new methods.

3.2.2 Random sampling

First, we implemented a classic random sampling pseudo-bulk approach (simulation #1). In this simulation approach the parameters explored are: (a) the number of individual cells sampled for each pseudo-bulk (n) and the total number of generated pseudo-bulk samples (pool size). We systematically varied the pool size across three categories: 100 cells, 10 000 cells and 100 000 cells. At the same time, we explored the impact of the number of samples, conducting simulations with 100, 500 and 100 samples.

To explore the impact of the above parameters we performed random sampling with replacement using a single cell 10X dataset from mouse brain (Hrvatin *et al.* 2018) composed of $\sim 48\,266$ cells and 28 samples (Supplementary Fig.

S1a–c). We observed that in this simulation, the number of cells sampled per sample affects the overall similarity of the distributions more than the number of samples (s). Moreover, we show that in the random sampling scenario the pseudo-bulks created are highly correlated and become almost identical as the number of sampled cells approaches the number of cells in the single-cell (Supplementary Fig. S1d and e). However, it is crucial to note that real-world scenarios deviate from such uniformity. Bulk samples within a dataset might have extremely different proportions from sample to sample. In addition, various samples may lack certain cell types, introducing a complexity that is not fully captured in the random sampling simulations. Moreover, when the number of cells sampled from a single-cell dataset is significantly smaller than the total cell count (e.g. 100 versus 22 000 cells), the simulation may only incorporate a sparse representation of each cell type. Moreover when we sample such a small percentage of cells per cell type, we might not capture rare cell types in our pseudo-bulk. It is essential to recognize the balance in random sampling—too much can lead to highly similar samples, while too little although it can yield diverse samples, they might fail to capture cell types effectively in terms of their signatures. It is crucial also to keep in mind when performing simulations that, in real bulk, samples typically comprise thousands to millions of cells.

3.2.3 Designing and comparing simulation methods

In our realistic simulations, we address limitations by introducing controlled noise. This noise originates from variances in proportions or the absence of cell types in pseudo-bulks. Simulation #2 involves random sampling using a uniform distribution (1–99) to select proportions, summed to 1. All cell types are sampled with replacement, resulting in a fixed range with large variance in proportions. To explore variability further, we introduce propVar in simulation #3, allowing us to adjust the uniform distribution range, producing arrays with more homogeneous to variable proportions. This parameter enables us to investigate the impact of proportion variability on deconvolution results. Additionally, we examine deconvolution methods in the presence of missing cell types in the bulk for simulation #3. A bimodal distribution with equal probability assigns cell types to mode 1 or 2. Mode 1 has a distribution with mean=0 and small SD=0.0001, while mode 2 has mean=1 with variable SD. This scenario mimics real-world situations where a tissue signature is known, but the presence of all cell types in the bulk is uncertain (Fig. 2a). We finally compare the above simulation methods with random sampling as well as simulations proposed in other studies (Chu *et al.* 2022, Hu and Chikina 2023).

4 Results and conclusions

To understand the differences across simulations we compare pseudo-bulk gene expression with 10 real bulk mouse samples derived from the same study. One way to compare different bulk profiles is to calculate the mean against the variance of gene expression in logarithmic scale or the coefficient of variation (CV). Real bulk and generated pseudo-bulk samples exhibit a negative binomial distribution, as expected, for expression data (Robinson and Smyth 2007) however, random sampling pseudo-bulk deviates the most from real bulk, with less overdispersion observed (Fig. 2b). Moreover, when we compare the Coefficient of Variation (CV) in the simulated

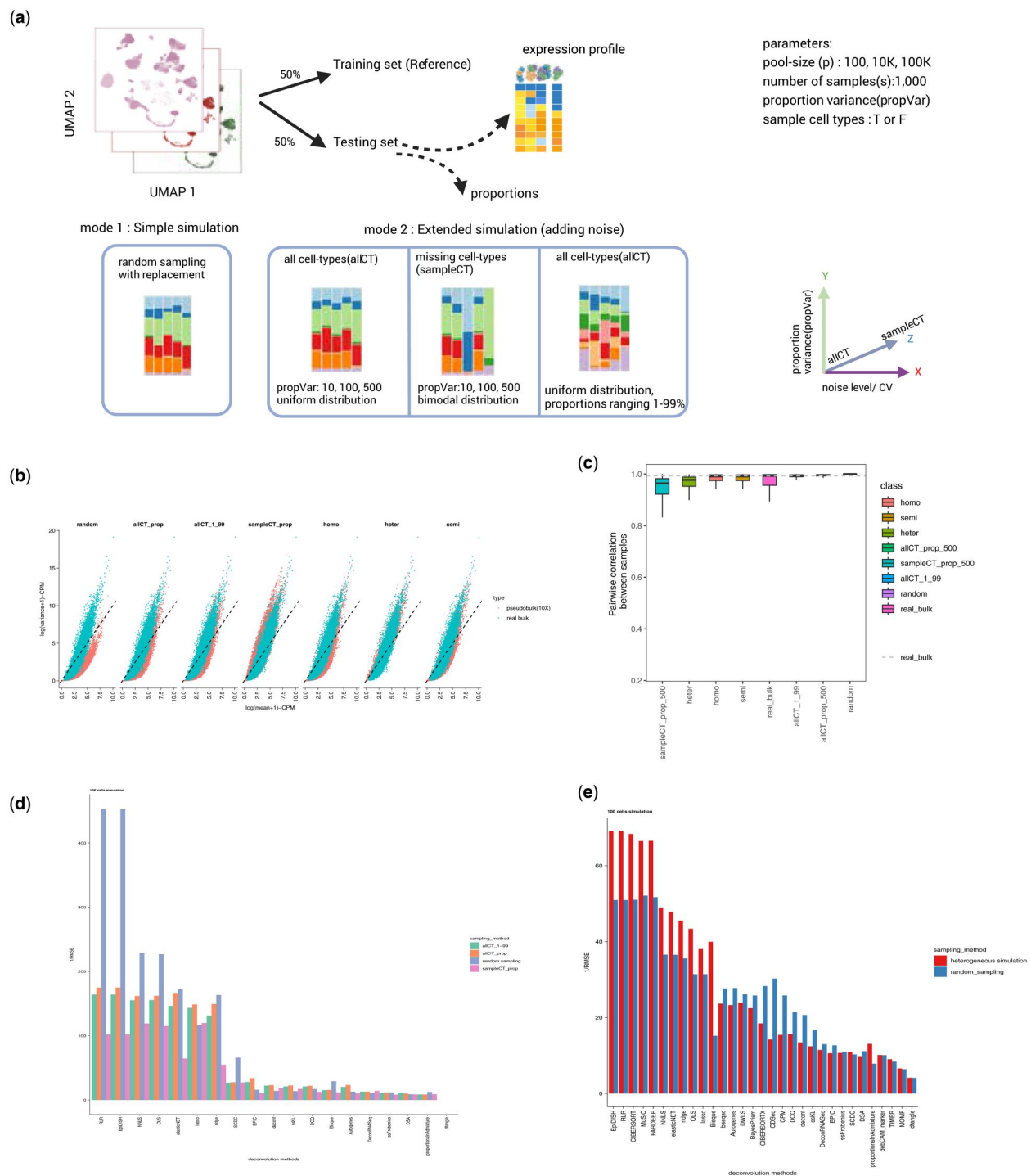


Figure 2. Evaluating pseudo-bulk simulations for benchmarking and new method development. (A) Simulation workflow: each single-cell matrix is split into half to generate a training set (reference) and a testing set (pseudo-bulk matrix). 4 different sampling techniques were used to generate pseudo-bulks and ground truth proportions. We set 4 different parameters for the simulation: the number of cells to be sampled from the single-cell (pool size (p)), the number of samples in each pseudo-bulk matrix (s)-the default is 1,000 samples, the variance of the proportions (pVar) for the cell-types and finally a logical vector to decide if we use all the cell-type or not (sampleCT, True or False)(graphic made in Biorender). (B) Scatter plot of mean and variance of expression data from 7 different simulation scenarios and matched real bulk data. All the simulated pseudo-bulk data derive from the same mouse 10X single-cell dataset and the real bulk RNA-seq derives from the same study as well (Hrvatin *et al.* 2018). (C) Boxplot plot showing the pairwise correlation of gene expression for the samples of each simulation and for the real bulk. (D) root mean square values (RMSE) of the 4 different simulation scenarios across the tested deconvolution method that were feasible to obtain results. (E) Root mean square values (RMSE) of random sampling and ‘heterogeneous’ simulation across the tested deconvolution method that were feasible to obtain results.

bulk against the CV in the real bulk, we observe that random, allCT_1_99 and allCT_propvar_500 are a lot less variant, sampleCT_prop_500, homogeneous and semi demonstrate higher CV, while heterogeneous simulation seemed to be

closer to real bulk expression CV (Supplementary Fig. S2a). At the same time we compared the methods based on the pairwise similarity of gene expression similar to here (Hu and Chikina 2023) (Supplementary Fig. S2b), showing that

heterogeneous and random simulation yield show lower gene correlation. We also compared the similarity of gene expression between samples in each pseudo-bulk strategy. A random sampling of 10k cells, and sampling all the cell types with variable proportions with replacement generates extremely similar pseudo-bulks while the other strategies produce samples with lower correlation, meaning that the simulated datasets using the latter approaches are more diverse and are closer to the bulk (Fig. 2c). Our results suggest that random sampling from single cells fails to generate pseudo-bulks that mimic real bulk expression and therefore should be avoided when developing new methods, as it oversimplifies the deconvolution problem. This has also been discussed before and aligns with results from previous studies. On the other hand, sampling cell types creates pseudo-bulks that can better mimic real bulk expression variability and can challenge deconvolution methods to deal with greater noise in the bulk. Using all the cell types from single cell with a medium variance of proportions ($\text{propVar} = 500$, allCT_prop) provides a good starting point for evaluation and method development; while sampling cell-types (sampleCT_prop), semi and heterogeneous simulation represents the noisiest scenarios that can challenge deconvolution methods as demonstrated by the performance of deconvolution methods (RMSE values, Fig. 2d). Many studies have talked about decrease in performance of methods when cell types are missing from the reference, here we show that this decrease is also happening when there are missing cell types in the bulk (sampleCT approach). Finally, we compared deconvolution performance in a mini-benchmark test between the random sampling and the heterogeneous method. We simulated 100 cells per sample for both methods and results demonstrate better performance of the heterogeneous method compared to the random sampling. This could result from the high variability depended on the small pool size (Fig. 2e). More extensive benchmarks across the different simulations would be needed in the future. Overall, we conclude that pseudo-bulk simulation parameters affect deconvolution severely and developed methods should be tested for performance in a range of scenarios before applying deconvolution to real data.

4.1 Effect of pre-processing methods in deconvolution results

Previous studies have reported that pre-processing of input data in deconvolution impacts the deconvolution performance (Avila Cobos *et al.* 2020, Jew *et al.* 2020), but it is yet unclear which methods should be used for each deconvolution method selected. Standard preprocessing of bulk and single-cell data includes the transformation and normalization of the expression matrices. Transformation methods account for the mean-variance dependencies as well as the extreme count values present in a dataset. This is a crucial pre-processing step for most downstream analyses of both single-cell and bulk RNA-seq data. For instance, both the identification of DEGs and PCA analysis requires the transformation of raw count values. The same principle applies to single-cell data for performing dimensionality reduction and identification of marker genes. Gene expression normalization is also an essential step for the above analyses. Normalization methods correct gene expression matrices for factors that affect the number of reads that will be mapped to a gene (gene length, GC content and sequencing depth) (Evans *et al.* 2017). In general, the goal of normalization is to enable the comparison of gene counts

within a sample and across samples. As mentioned before, deconvolution requires reference data that provides cell-type specific information in the form of single-cells, summed single-cells or marker genes derived from single-cell and bulk RNA-seq data to be deconvolved. Here, we test how correcting data with transformation and normalization techniques affect the deconvolution process. We implemented 2 transformation methods and 11 different normalization methods in our pipeline (Supplementary Table S2). Previous studies have applied transformation of the data first and then normalization. Here we show that normalizing the matrices first and then applying transformation yields better deconvolution results, on average in 8 brain tissues (Fig. 3a and b, Supplementary Fig. S3). We then examine the effect of different normalization and transformation methods. Results from deconvolution in brain tissue show that the majority of methods output better results when the input data are not transformed. However, bseq-sc and BisqueRNA methods benefit from the logarithmic transformation of the input matrices (Fig. 3c). Moreover, in many cases, sqrt-transformation seems to perform worse than no-transformation with the exception of the penalized linear regression methods lasso, elastic net and ridge where it performs the same as no transformation of the matrices. Lastly, we tested 11 different normalization methods across 3 different brain datasets that have been produced from different single-cell technologies (10X and SmartSeq2). In all three datasets, row normalization seems to perform worse across all the deconvolution methods (lower Pearson Correlation values), while the other normalization methods appear to perform differently depending on the dataset. The majority of normalization methods yield good results for the first dataset. On the other hand, for the other two 10X datasets, none, global min-max, column-z score and mean normalization methods appear to yield better results (Fig. 3d-f).

4.2 Batch effects between (pseudo) bulk and single-cell data affect deconvolution

The first methods developed for deconvolution were utilizing flow cytometry data from sorted cells, marker gene signatures or no references at all (unsupervised methods) to deconstruct bulk RNA-seq and microarray samples (Abbas *et al.* 2009, Gaujoux and Seoigie 2013, Chen *et al.* 2018). Nowadays, with the advancement of single-cell protocols and the extensive availability of single-cell data from different species, a great opportunity has appeared to utilize them to facilitate deconvolution.

4.2.1 Bulk versus single-cell RNA-seq

Newly developed deconvolution methods are making use of single-cell data to deconvolve bulk profiles by selecting appropriate features (genes) such as MuSiC, DWLS (Tsoucas *et al.* 2019) and many others. Yet, the large differences between single-cell and bulk quantification methods of expression data can impact the outcome of deconvolution. Here, we compared the expression of each gene between single-cell data and real bulk data from the same tissue. We observed significant expression differences between summed single-cell data (pseudo-bulk) and real bulk data originating from the same tissue, here pancreas (Fadista *et al.* 2014, Baron *et al.* 2016) (Supplementary Fig. S4a and b). In contrast, when comparing real bulk mixtures with each other, or pseudo-bulks with each other we observe very little differences (Supplementary Fig. S4c and d). These discrepancies,

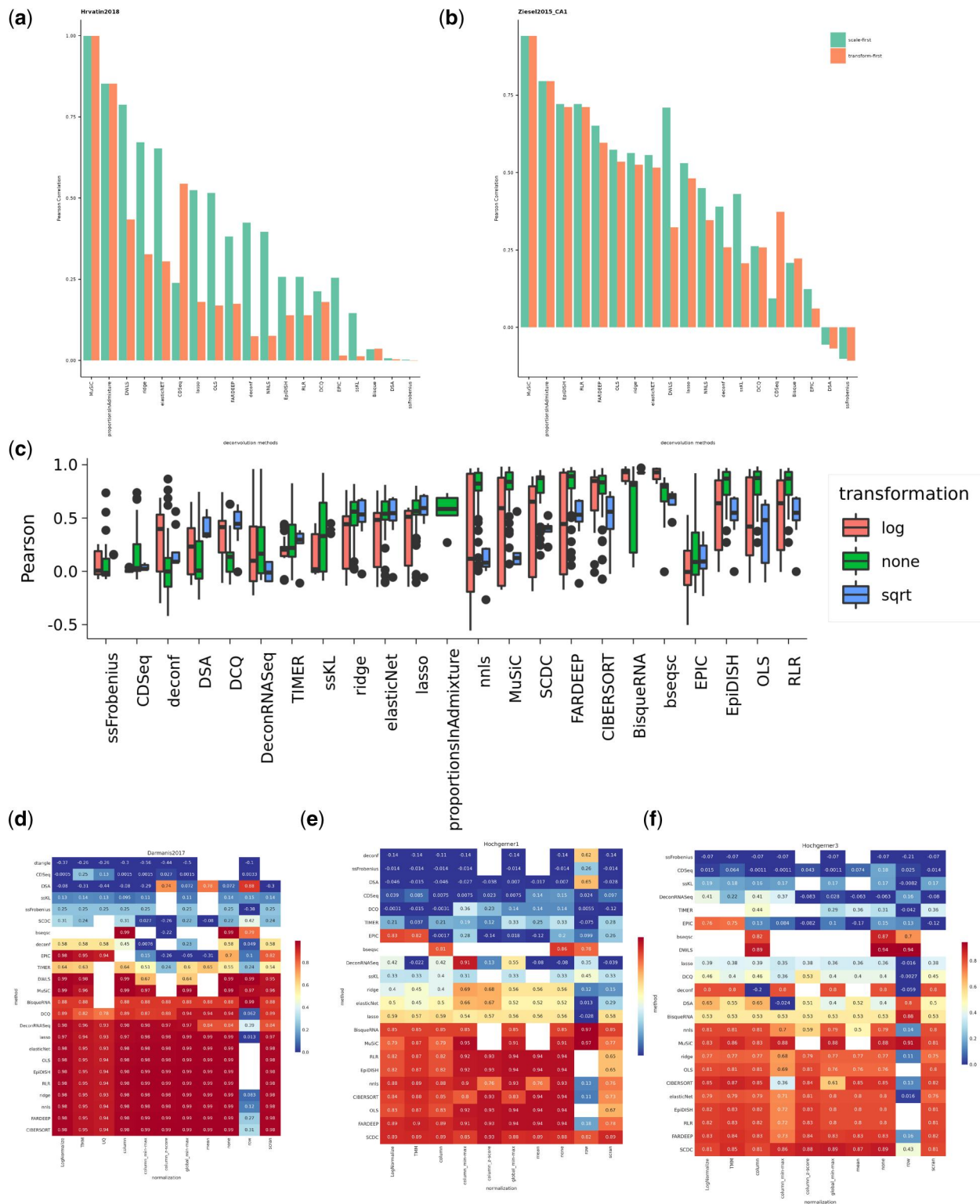


Figure 3. Pre-processing steps affect downstream decomposition of samples. Barplot showing Pearson Correlation values of a self-reference deconvolution task in which the effect of scaling first or transforming first the input matrices is tested. (A) Hrvatin *et al.* (2018)-brain tissue. (B) Ziesel 2015-brain tissue. (C) Boxplots showing Pearson correlation values estimated from the deconvolution of simulated data (from brain tissue) across 2 different transformation methods (log and sqrt) as well as no transformation (linear scale). Each data point in the boxplot represents a combination of a transformation and a normalisation method utilised in the deconvolution task (e.g log+TMM, log+column and similar for all the possible combinations). (D-F) Results from deconvolution of 3 different simulated data from brain tissues (Darmanis 2017, Hochgerner1, Hochgerner3) across 11 normalisation methods and 24 tested deconvolution methods.

observed between bulk and single-cell, are most likely due to the batch effects deriving from the different nature of the two sequencing techniques and dissimilar quantification methods.

PCA of the two data types showed that PC1 could explain 74.5% of the variance between the bulk and the single cell (Supplementary Fig. S5a-c).

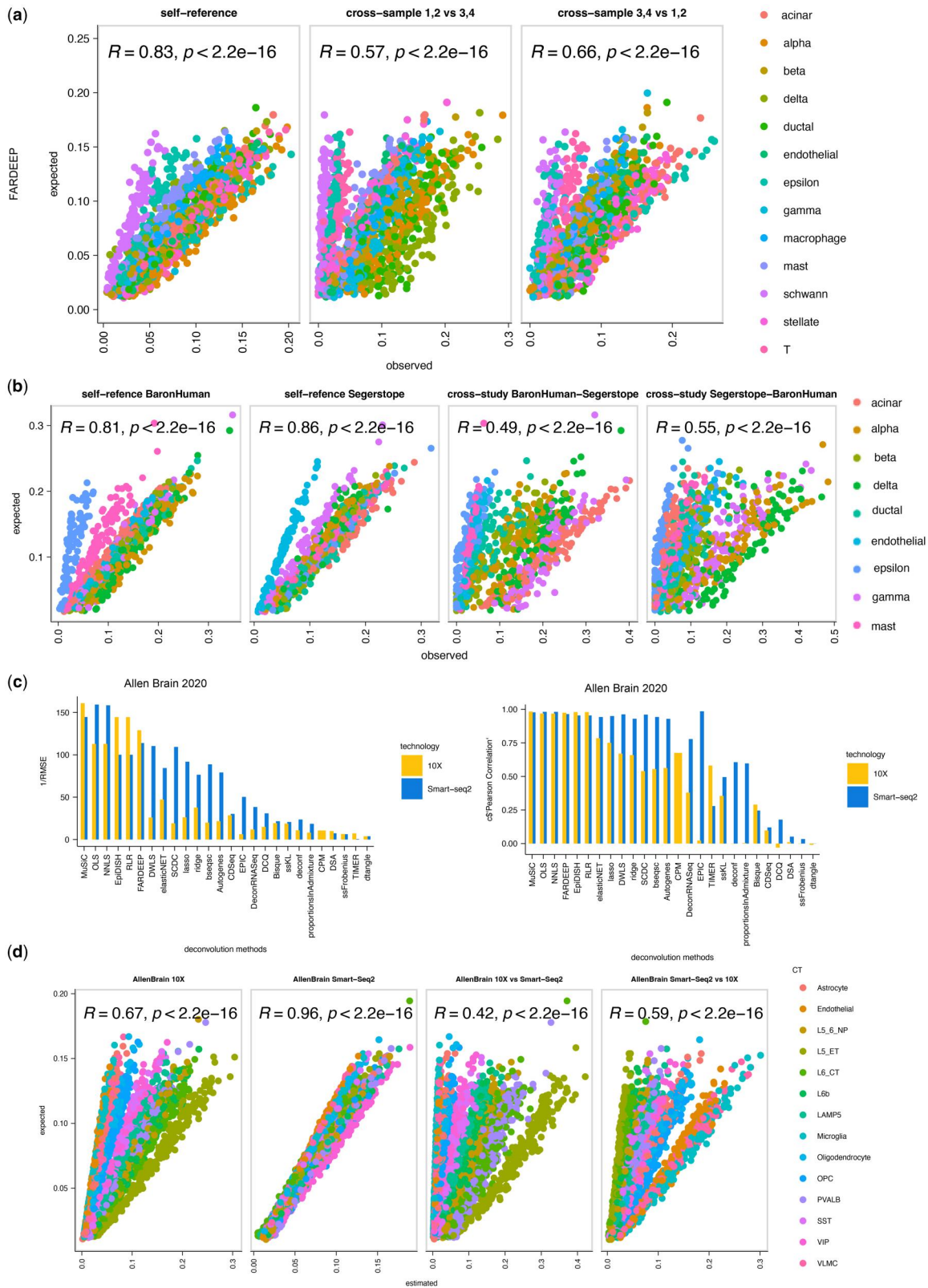


Figure 4. Impact of batch effects in deconvolution. (A) Scatter plots showing results from self-reference and 2 examples of cross-sample deconvolution from the method FARDEEP utilising a single-cell from human pancreas (13 cell-types). In the self-reference task all 4 samples are used while cross-sample tasks use two samples as reference and the other two represent the pseudo-bulk at each example. Pearson correlation (R) and p -value are reported for each task. (B) Scatter plots showing the results from two self-reference tasks from two pancreatic studies (9 cell-types) as well as two cross-study deconvolution tasks. (C) Bar plots showing $1/RMSE$ and Pearson correlation values from the 26 feasible deconvolution runs from two human brain datasets that come from the same study but different technologies (10X and Smart-seq2). (D) Scatter plots demonstrating results from the DWLS method from the two self-reference tasks (10X and Smart-seq2 of human brain tissues, 14 cell-types) and two cross-technology tasks. Pearson correlation (R) and p -values have been calculated for each task and shown on the plots.

4.2.2 Self-reference versus cross-reference deconvolution tasks

Next, we generated pseudo-bulk and single-cell references by randomly splitting the same dataset in half and then performing deconvolution afterwards. This task is a very simple case of deconvolution since both the bulk and the reference contain the same gene expression information. We name this ‘self-reference deconvolution’. While self-reference gives the opportunity to study deconvolution using simulations, it oversimplifies the problem. For this reason, we introduced cross-reference deconvolution tasks to more systematically study the impact of batch effects in cell-type deconvolution. For this task, we split the same pancreatic single-cell dataset in half and we selected cells from two individuals to generate the reference while the other two biological samples were used to build the bulk profiles. In this way, we artificially introduced batch effects in cellular decomposition. Deconvolution results from one “bulk” and one “single-cell” algorithm (FARDEEP and DWLS) show a significant drop in deconvolution performance in the cross-sample task compared to the equivalent self-reference task (Fig. 4a and Supplementary Fig. S6a). These results suggest that cell-type deconvolution is severely affected by batch effects caused by expression differences across samples. Furthermore, since the deconvolution of real bulk samples typically entails expression data from completely different studies and the use of different technologies, we next studied how cross-study and cross-technology tasks can impact the performance of deconvolution algorithms. For the cross-study deconvolution, we selected two human pancreatic single-cell datasets (Baron *et al.* 2016,

Segerstolpe *et al.* 2016) (Supplementary Table S1) which were generated from two different laboratories. We observed that deconvolution across studies can impact deconvolution performance severely when compared to self-reference tasks (Fig. 4b and Supplementary Fig. S6b). Next, we compared the performance of deconvolution when datasets from different technologies are involved. For this task, we focused on 3 independent dataset pairs from the placenta, brain and lung tissue (Supplementary Figs S7 and S8). Each dataset pair was generated from the same laboratory but using different single-cell technologies (10X and Smart-Seq2). Given that 10X and Smart-Seq2 are the most common single-cell protocols available, this test will help reveal which references are suitable in deconvolution. Results show that self-reference deconvolution with 10X sc-RNASeq data yields better results overall (higher Pearson correlation, lower RMSE values Fig. 4c, Supplementary Fig. S6c and d) across most deconvolution methods compared to Smart-seq2 self-reference. Moreover, results from the top method in this task (DWLS) demonstrate that cross-technology deconvolution results in less accurate predictions, with very low Pearson Correlation values observed (Fig. 4d). All the above highlight the importance of the reference selection and the barrier that batch effects oppose to deconvolution.

4.3 Evaluation and consensus-based approach for deconvolution in real biological samples

4.3.1 Deconvolution with flow cytometry as ground truth

In contrast to the simulation tasks performed above, here we decompose real bulk profiles from nine PBMCs (Finotello

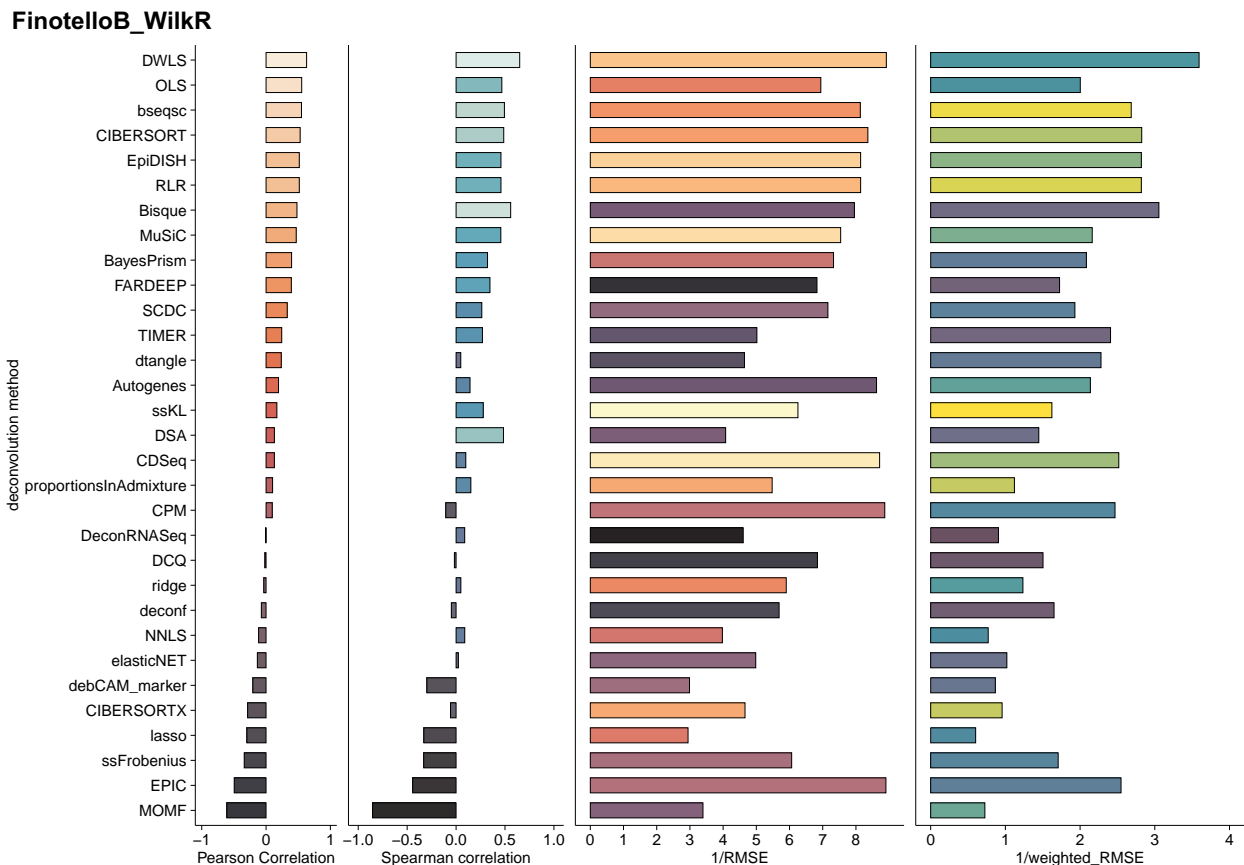


Figure 5. Deconvolution performance of real bulk RNA-seq samples. Results from 4 main metrics (Pearson Correlation, Spearman Correlation, RMSE and weighted RMSE for 31 deconvolution methods on real bulk PBMC data with ground truth using a single-cell PBMC dataset as reference (Wilk *et al.* 2020).

et al. 2019) samples, utilizing our deconvolution pipeline to assess the results and generate an honest consensus of predictions of cell-type proportions. Apart from the publicly available RNA-seq dataset from the above study, there is also flow cytometry data available for the nine individuals which can be used as ground truth in our deconvolution pipeline. The researchers have measured the proportions of 8 immune cell types using specific markers for each cell population (CD4+ T cells, CD8+ T cells, T-regulatory cells, B cells, NK cells, myeloid dendritic cells, monocytes and neutrophils (Supplementary Fig. S9a and d). For this deconvolution task, we selected a recent single-cell PBMC dataset (Wilk *et al.* 2020) to use as a reference. We re-analysed the single-cell data and matched the cell types between the single cell and the ground truth data from flow cytometry so that our reference has the same cell types that had been measured in the bulk (Supplementary Fig. S10a–e). We also summarized the proportions of CD4+positive and T-regulatory cells measured by flow cytometry since our single-cell did not contain T-regulatory cells (Supplementary Figs S9b and c and S10e). Before we apply deconvolution on the 9 real bulk samples we first performed a self-reference deconvolution to evaluate how the reference performs in the pseudo-bulk deconvolution task. Nine methods yield good results with a Pearson Correlation of more than 0.86 (Supplementary Fig. 10f). Next, we applied the all 31 deconvolution methods included in the pipeline on the real bulk PMN/PBMC data and obtained the cell-type predictions. The methods were evaluated across nine evaluation metrics. Pearson Correlation, Spearman Correlation, RMSE and weighted RMSE values are the main metrics we focused on in this benchmark task (Fig. 5). Results from the additional metrics implemented in the pipeline are shown in Supplementary (Supplementary Fig. S11). DWLS, OLS and bseq-sc showed the highest Pearson Correlation values while DWLS and bseq-sc also maintained good RMSE and Weighted RMSE values (Fig. 5). Looking closer at the DWLS cell-type predictions we observe differences across samples (Supplementary Fig. S12a and b). Moreover, we examined the deconvolution predictions per cell type and noticed that CD4+ cells were consistently overestimated while CD8+ cells were underestimated in all nine donors by the DWLS deconvolution method. The same pattern can be validated with the bseq-sc method which was ranked second in this task based on Pearson Correlation values (Supplementary Fig. S12). It has been reported before (Aliee and Theis 2021) that it is difficult to deconvolve closely related cell types due to collinearity, a problem which some single-cell methods attempt to solve with sophisticated feature selection. Since the large benchmarking across all the deconvolution methods has been performed using all the genes as initial input (except for the marker specific methods (e.g. DSA, deCAMmarker, dtangle) we explored how different signatures affect deconvolution results of marker-based and reference-based methods. For this analysis, we compared deconvolution results using all the genes, different marker selection algorithms (MAST, Wilcox test, *t*-test) and LM22 immune signature. Results show that many methods perform well when LM22 signature is used but their performance decreases when other sets of signatures are selected. Nevertheless methods such as OLS, EpiDISH, RLR, CIBERSORT, and other seem to be more robust in the selection of different signatures (Supplementary Fig. S13a and b).

4.3.2 Deconvolution with tissue slides as ground truth

Since there is no single method that clearly ranks first in all scenarios and it is still unclear how different methods work in data from different tissues, we believe that a fair deconvolution consensus is a reasonable approach to take. Here, we suggest a consensus method based on 3 methods DWLS, FARDEEP and EpiDISH (selection criteria on Methods) that perform well in previous pseudo-bulk scenarios and compute proportions relatively fast and reliably (Supplementary Fig. S10). Bseq-sc although it performs well, because of its iterative nature, it can be very slow (> 1 week) to compute proportions so it was excluded from the consensus calculations. Next, we applied the consensus on GTEx stomach data (339 samples) for which we had access to tissue slides that could inform the validation of the deconvolution results. For this task we used as a reference the stomach subset of the human cell landscapes sRNA-seq dataset containing 16 cell types (Fig. 6a and b). From the deconvolution results we can observe a clear difference between the samples on the left of the barplot and the right hand of the barplot. Samples on the right seem to contain mostly smooth muscle cells, fibroblasts and endothelial cells, whereas samples on the left are mainly peptic and goblet cells (Fig. 6c). Next, we explored if this difference in the composition can be observed in the tissue slides of the adjacent samples. We randomly selected a small number of tissue samples (five from the samples from the left side and five from the right side). We observed that the slides that come from samples from the right side of the barplot, were in general more enriched in muscle cells compared to the samples on the left side which were more enriched in mucosa and submucosa tissue which aligns with the presence of goblet and peptic cells in the deconvolved mixtures (Fig. 6e and f). We also quantified the mucosa, submucosa and muscle layers (μm^2) for each stomach sample using QuPath and computed the proportions of each layer in each sample. Barplot showing the proportions computed for the 339 samples in the sample ordered the same way as in Fig. 6c. Although we can observe that the samples on the right have a lot less mucosa and more muscle which aligns with the deconvolution results (the absence of peptic and goblet cells on the samples on the right and presence of smooth muscle), the measurements from the tissue slides seem to be quite different from the deconvolution results overall. This could be because of a number of reasons: (a) annotating mucosa, submucosa and muscle layers is not straightforward and is dependent on the individual who performs the annotation, (b) deconvolution predictions are objected to overfitting, (c) tissue slides might have different composition compared to the bulk sample that goes to sequencing due to sampling reasons. More detailed analysis on the comparison between composition in bulk samples and tissue slides should be made to conclude if tissue slides are a good ground truth data. To further explore the differences between deconvolution results and composition calculations from tissue slides we mapped each cell-type to the respective layer based on literature knowledge (see Section 2) and we performed Pearson correlation analysis. We observe the highest positive correlation values between the Mucosa proportions from tissue slides and the summarized proportions of mucosa cells from deconvolution and the similar Pearson correlation values between the Muscle layer and the summarized proportions from the cell types in the muscle layer. However the submucosa proportion from the tissue slide is not strongly correlated with any of the cell type proportions from deconvolution.

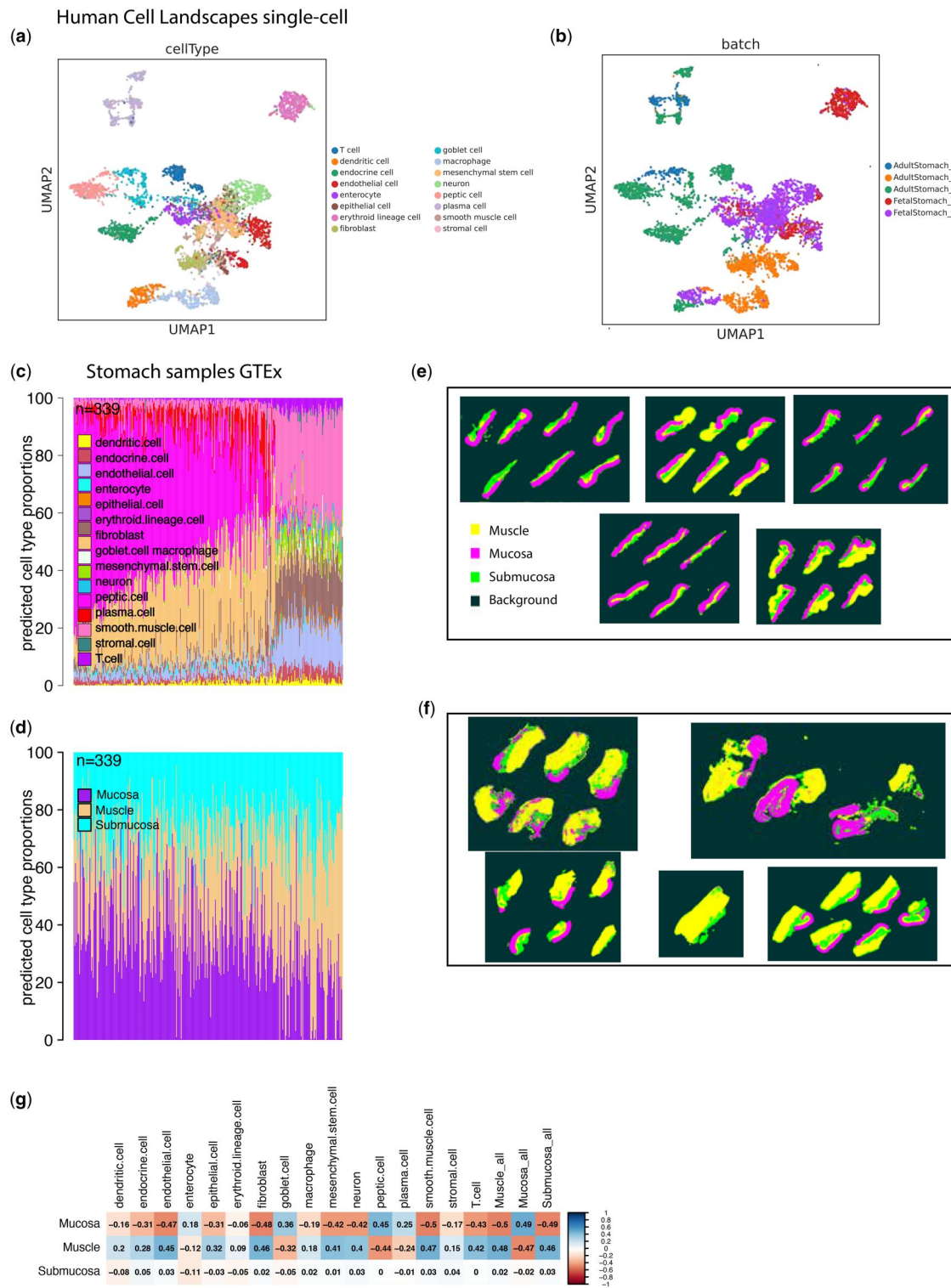


Figure 6. Application of consensus approach on stomach GTEx bulk expression data. (A,B) UMAP plots of the stomach subset of Human Cell Landscapes single-cell reference used for the stomach GTEx deconvolution showing the different annotated cell-types and the distribution of donors in the datasets. (C) Bar plot showing the prediction of cell-type proportions from the consensus approach across 339 selected stomach samples that have matched tissue slides. (D) Quantification of the different stomach layers (mucosa, submucosa, muscle) from the matched tissue slides using QuPath. Area μm^2 is measured for each layer in each tissue slide. (E) 5 randomly selected tissue slides from the left side of the barplot (C). Submucosa is shown in green, mucosa in magenta and muscle compartment in green. (F) 5 randomly selected tissue slides from samples originating from the right side of barplot (C). (G) correlation plot showing the Pearson Correlation values between the proportions measured from tissue slides and the predicted proportions from deconvolution of the bulk RNA-seq samples.

4.3.3 Background prediction using experiments from sorted cells

Next, we designed an experiment in order to evaluate the background prediction of the consensus method in two datasets with purified cell types and/or mixtures with known cell types. Our data collection included mouse cerebral cortex and human lung datasets, for which we carefully selected relevant single-cell references tailored to the task (Supplementary Fig. S14a and e). Initially, we applied the consensus deconvolution method to 17 samples obtained from the mouse cerebral cortex (Supplementary Fig. S14b). Each sample in this task represented a single cell type, serving as the ground truth, and aimed to evaluate the method's ability to accurately predict the presence of that specific cell type while quantifying background prediction. To quantify background prediction, we employed logarithmic loss, a widely used metric for evaluating classification model performance. In the context of deconvolution, this metric measures the accuracy of estimated cell type proportions within a mixture, taking into account uncertainties and penalizing incorrect estimates arising from low-level gene expression that may lead to background predictions. Across the 17 samples, we observed consistently low log-loss values, ranging from approximately 0.004–0.104 (Supplementary Fig. S14c). Notably, among all samples, two endothelial flow cytometry samples and one myelinating oligodendrocyte sample exhibited the lowest log-loss values, signifying precise predictions of the purified cell type within each sample. Moving forward, we extended our analysis to encompass 26 samples from human fetal lung. Of these, the initial 16 samples were predominantly composed of endothelial cells, while the subsequent 20 samples contained non-endothelial cells. Our objective was to observe whether the consensus approach could effectively differentiate between these two sample categories. Notably, the endothelial samples predominantly exhibited high proportions of endothelial cells in the predictions, albeit with predictions of fibroblasts, lymphocytes, and myeloid cells. In contrast, the non-endothelial samples were primarily predicted to consist of fibroblasts, with minor proportions of endothelial cells detected. These low-level endothelial cell predictions may be attributed to background predictions generated by the consensus algorithm (Supplementary Fig. S14d).

4.3.4 Deconvolution in the absence of ground truth

In situations where ground truth data is unavailable, our pipeline incorporates a separate feature to address this. We perform a pairwise comparison of the deconvolution results across all samples, enabling us to estimate Pearson Correlation coefficients. These coefficients serve as a quantitative measure of the agreement between different methods. The results of this comparison are visually presented in the form of a heatmap, where methods that yield similar results cluster together (Supplementary Fig. S10g). This clustering provides valuable insights into which methods exhibit strong agreement, both in the presence and absence of ground truth data.

It is important to mention that while the clustering approach successfully identifies method agreement, it does not inherently imply that the clustered methods consistently perform well. Finally, to comprehensively assess the reliability and scalability of our pipeline, we conducted parallel evaluations of each module, encompassing pseudo-bulk generation, normalization, and deconvolution methods. These assessments provide valuable insights into the computational efficiency of our pipeline, reporting essential metrics such as processing time, memory

usage, mean load, and CPU time for each module (Supplementary Fig. S15a and b).

5 Discussion

In this study we developed a deconvolution framework that allows the evaluation of 31 publicly available methods across different scenarios with both synthetic and real data. We first studied the effect of different pseudo-bulk techniques by defining different ways of pseudo-bulk sampling. Our results demonstrated that pseudo-bulks generated with different methods produce highly diverse results in deconvolution, which shows that the results of deconvolution are highly sensitive to the sampling of cells and the selection of proportions. Here, we suggest the use of pseudo-bulk generation methods that can challenge deconvolution methods and reflect real bulk data, when developing new, or benchmarking existing methods. Another important aspect of deconvolution, that we explore, is the normalization and the transformation of the input matrices in deconvolution. Both bulk and single-cell matrices used in the deconvolution pipeline were tested across different normalization and transformation methods. Results suggest that normalization can be beneficial in deconvolution since it corrects for differences in the library size both in the single cell and the bulk RNA-seq data. On the other hand, logarithmic transformation of the input matrices can result in worse performance with the exception of two deconvolution methods (bisqueRNA and bseq-sc). Similar to the reports from previous studies, this suggests that keeping input data in linear scale aids in assessing the cell proportions accurately (Avila Cobos *et al.* 2020, Jin and Liu 2021).

Another key aspect of the bulk sample deconstruction is the selection of the reference matrix, which is used to extract the key features (genes) that will be used in deconvolution. Results from Smart-seq2 and 10X single-cell datasets that have been used in this study, show that 10X data is consistently performing better as a reference, this is likely because this technology captures more cells of a specific cell type and as a result, more accurate features are extracted from these data. It should also be noted that in tasks when the bulk and the reference data come from different sources (studies or technologies) —the accuracy score of the deconvolution drops significantly. This indicates that differences in the expression of genes, across datasets, affect deconvolution heavily. Methods that take into account and successfully minimize these effects (DWLS, MuSiC) seem to also perform better overall in cross-reference tasks. We apply deconvolution in real PBMC samples using a suitable matching single-cell reference and conclude that two methods, DWLS and bseqsc perform best. All the above steps that should be performed before deconvolution have been implemented as part of the pipeline with a number of parameters that can be defined by the user.

Notably, in the absence of ground truth data, our pipeline introduces a valuable feature for pairwise comparisons of deconvolution results across all samples, allowing the estimation of Pearson Correlation coefficients. This innovative approach offers a quantitative measure of method agreement, even without the presence of ground truth data, providing researchers with insights into method performance. It is imperative to understand that while this clustering technique effectively identifies method agreement, it doesn't inherently imply consistent method performance. We also suggest a

consensus method based on three robust techniques: DWLS, FARDEEP, and EpiDISH and test it across three real bulk datasets with available ground truth testing its accuracy and background prediction.

Notably, the CATD pipeline, developed using Snakemake, enables the reproducibility of our results and can be additionally utilized by users to deconvolve new bulk samples by providing a bulk dataset and a single-cell dataset from the same tissue. When ground truth is provided too, the pipeline will provide evaluation metrics for each method. On the other hand when ground truth is absent the pipeline will perform pairwise comparison of the results and report a heatmap with methods to be selected as well as calculating proportions based on the consensus method.

Moreover, the advent of new high-quality single-cell atlases with standardized annotation from both healthy and disease conditions will provide more opportunities to explore the contribution of cell proportions in pathological conditions. Finally, the development of robust, accurate tools for feature selection that enable clear distinction between cell-types will pave the way for use of computational deconvolution perhaps even in clinics or in settings where scRNA-Seq is not feasible.

Acknowledgements

The authors would like to thank the reviewers for their valuable feedback and suggestions for the manuscript. The authors would also like to thank Craig Russell for his feedback and help in reviewing and editing the manuscript.

Author Contributions

Anna Vathrakokoili Pournara (Conceptualization, Data curation, Formal Analysis, Investigation, Methodology, Software, Validation, Visualization, Writing—original draft), Zhichao Miao (Conceptualization, Data Curation, Methodology, Software, Supervision, Writing—review & editing), Ozgur Beker (Methodology, Formal Analysis, Software), Nadja Nolte (Data curation, Formal Analysis, Investigation, Methodology, Software, Validation, Visualization), Alvis Brazma (Conceptualization, Supervision, Writing—review & editing), and Irene Papatheodorou (Conceptualization, Project Administration, Supervision, Writing—review & editing)

Supplementary data

[Supplementary data](#) are available at *Bioinformatics Advances* online.

Funding

This work was supported by: the European Molecular Biology Laboratory (A.V.P, Z.M., O.B.,N.N, and A.B, I.P.); the EMBL international PhD program (A.V.P.); and the OpenTargets (N.N, C.M, I.P.) [OTAR2067].

Conflict of interest

The authors declare no competing interests.

Data availability

The data underlying this article are available at https://github.com/Papatheodorou-Group/CATD_snakemake.

References

- Abbas AR, Wolslegel K, Seshasayee D *et al.* Deconvolution of blood microarray data identifies cellular activation patterns in systemic lupus erythematosus. *PLoS One* 2009;4:e6098. <https://doi.org/10.1371/journal.pone.0006098>
- Aliee H, Theis FJ. AutoGeneS: automatic gene selection using multi-objective optimization for RNA-seq deconvolution. *Cell Syst* 2021; 12:706–15.e4. <https://doi.org/10.1016/j.cels.2021.05.006>
- Alonso-Moreda N, Berral-González A, De La Rosa E *et al.* Comparative analysis of cell mixtures deconvolution and gene signatures generated for blood, immune and cancer cells. *Int J Mol Sci* 2023;24:10765. <https://doi.org/10.3390/ijms241310765>
- Aran D, Hu Z, Butte AJ. xCell: digitally portraying the tissue cellular heterogeneity landscape. *Genome Biol* 2017;18:220. <https://doi.org/10.1186/s13059-017-1349-1>
- Avila Cobos F, Alcúrcira-Hernández J, Powell JE *et al.* Benchmarking of cell type deconvolution pipelines for transcriptomics data. *Nat Commun* 2020;11:5650. <https://doi.org/10.1038/s41467-020-19015-1>
- Bankhead P, Loughrey MB, Fernández JA *et al.* QuPath: open source software for digital pathology image analysis. *Sci Rep* 2017;7: 16878. <https://doi.org/10.1038/s41598-017-17204-5>
- Becht E, Giraldo NA, Lacroix L *et al.* Estimating the population abundance of tissue-infiltrating immune and stromal cell populations using gene expression. *Genome Biol* 2016;17:218. <https://doi.org/10.1186/s13059-016-1070-5>
- Chen B, Khodadoust MS, Liu CL *et al.* Profiling tumor infiltrating immune cells with CIBERSORT. *Methods Mol Biol* 2018;1711: 243–59. https://doi.org/10.1007/978-1-4939-7493-1_12
- Chu T, Wang Z, Pe'er D *et al.* Cell type and gene expression deconvolution with BayesPrism enables bayesian integrative analysis across bulk and single-cell RNA sequencing in oncology. *Nat Cancer* 2022;3:505–17. <https://doi.org/10.1038/s43018-022-00356-3>
- Cobos FA, Panah MJN, Epps J *et al.* Effective methods for bulk RNA-seq deconvolution using scRNA-seq transcriptomes. *Genome Biol* 2023;24:177. <https://doi.org/10.1186/s13059-023-03016-6>
- CZI Single-Cell Biology Program. CZ CELLxGENE discover: a single-cell data platform for scalable exploration, analysis and modeling of aggregated data. *bioRxiv* 2023.10.30. <https://doi.org/10.1101/2023.10.30.563174>
- Denisenko E, Guo BB, Jones M *et al.* Systematic assessment of tissue dissociation and storage biases in single-cell and single-nucleus RNA-seq workflows. *Genome Biol* 2020;21:130. <https://doi.org/10.1186/s13059-020-02048-6>
- Dietrich A, Sturm G, Merotto L *et al.* SimBu: Bias-aware simulation of bulk RNA-seq data with variable cell type composition. *Bioinformatics* 2022; 38:ii141–ii147. <https://doi.org/10.1101/2022.05.06.490889>
- Dimitrakopoulou K, Wik E, Akslen LA *et al.* Deblender: a semi-/unsupervised multi-operational computational method for complete deconvolution of expression data from heterogeneous samples. *BMC Bioinform* 2018;19:408. <https://doi.org/10.1186/s12859-018-2442-5>
- Dong M, Thennavan A, Urrutia E *et al.* SCDC: bulk gene expression deconvolution by multiple single-cell RNA sequencing references. *Brief Bioinform* 2021;22:416–27. <https://doi.org/10.1093/bib/bbz166>
- Donovan MKR, D'Antonio-Chronowska A, D'Antonio M *et al.* Cellular deconvolution of GTEx tissues powers discovery of disease and cell-type associated regulatory variants. *Nat Commun* 2020;11:955. <https://doi.org/10.1038/s41467-020-14561-0>
- Dumont N, Liu B, Defilippis RA *et al.* Breast fibroblasts modulate early dissemination, tumorigenesis, and metastasis through alteration of extracellular matrix characteristics. *Neoplasia* 2013;15:249–62. <https://doi.org/10.1593/neo.121950>

- Edgar R, Domrachev M, Lash AE. Gene expression omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res* 2002;30:207–10. <https://doi.org/10.1093/nar/30.1.207>
- Evans C, Hardin J, Stroebel DM. Selecting between-sample RNA-Seq normalization methods from the perspective of their assumptions. *Brief Bioinform* 2017;19:776–92. <https://doi.org/10.1093/bib/bbx008>
- Fadista J, Vikman P, Laakso EO *et al.* Global genomic and transcriptomic analysis of human pancreatic islets reveals novel genes influencing glucose metabolism. *Proc Natl Acad Sci USA* 2014;111:13924–9. <https://doi.org/10.1073/pnas.1402665111>
- Finotello F, Mayer C, Plattner C *et al.* Molecular and pharmacological modulators of the tumor immune contexture revealed by deconvolution of RNA-seq data. *Genome Med* 2019;11:34. Available at: <https://doi.org/10.1186/s13073-019-0638-6>
- Frishberg A, Peshes-Yaloz N, Cohn O *et al.* Cell composition analysis of bulk genomics using single-cell data. *Nat Methods* 2019;16:327–32. <https://doi.org/10.1038/s41592-019-0355-5>
- Garmire LX, Li Y, Huang Q *et al.* Challenges and perspectives in computational deconvolution of genomics data. *Nat Methods* 2024;21:391–400. <https://doi.org/10.1038/s41592-023-02166-6>
- Gaujoux R, Seoighe C. CellMix: a comprehensive toolbox for gene expression deconvolution. *Bioinformatics* 2013;29:2211–2. <https://doi.org/10.1093/bioinformatics/btt351>
- Gaujoux R, Seoighe C. Semi-supervised nonnegative matrix factorization for gene expression deconvolution: a case study. *Infect Genet Evol* 2012;12:913–21. <https://doi.org/10.1016/j.meegid.2011.08.014>
- Gierahn TM, Wadsworth MH, Hughes TK *et al.* Seq-Well: portable, low-cost RNA sequencing of single cells at high throughput. *Nat Methods* 2017;14:395–8. <https://doi.org/10.1038/nmeth.4179>
- Gong T, Hartmann N, Kohane IS *et al.* Optimal deconvolution of transcriptional profiling data using quadratic programming with application to complex clinical blood samples. *PLoS One* 2011;6:e27156. <https://doi.org/10.1371/journal.pone.0027156>
- GTEX Consortium. Genetic effects on gene expression across human tissues. *Nature* 2017;550:204–13. <https://doi.org/10.1038/nature24277>
- Jaakkola MK, Elo LL. Computational deconvolution to estimate cell type-specific gene expression from bulk data. *NAR Genomics and Bioinformatics* 2021;3:lqaa110. <https://doi.org/10.1093/nargab/lqaa110>
- Hanahan D, Coussens LM. Accessories to the crime: functions of cells recruited to the tumor microenvironment. *Cancer Cell* 2012;21:309–22. <https://doi.org/10.1016/j.ccr.2012.02.022>
- Hao Y *et al.* 2019. Fast and robust deconvolution of tumor infiltrating lymphocyte from expression profiles using least trimmed squares. In Ioshikhes I. (ed.), *PLOS Computational Biology*, 15, p. e1006976. <https://doi.org/10.1371/journal.pcbi.1006976>
- Hu M, Chikina M. Heterogeneous pseudobulk simulation enables realistic benchmarking of cell-type deconvolution methods. *bioRxiv* 2023.01.05.522919. <https://doi.org/10.1101/2023.01.05.522919>
- Hashimshony T, Wagner F, Sher N *et al.* CEL-Seq: single-Cell RNA-Seq by multiplexed linear amplification. *Cell Rep* 2012;2:666–73. <https://doi.org/10.1016/j.celrep.2012.08.003>
- Hrvatin S, Hochbaum DR, Nagy MA *et al.* Single-cell analysis of experience-dependent transcriptomic states in the mouse visual cortex. *Nat Neurosci* 2018;21:120–9. <https://doi.org/10.1038/s41593-017-0029-5>
- Hudson TJ, Anderson W, Artez A, *et al.* International network of cancer genome projects. *Nature* 2010;464:993–8. <https://doi.org/10.1038/nature08987>
- Inkeles MS, Teles RM, Pouldar D *et al.* Cell-type deconvolution with immune pathways identifies gene networks of host defense and immunopathology in leprosy. *JCI Insight* 2016;1:e88843. <https://doi.org/10.1172/jci.insight.88843>
- Jew B, Alvarez M, Rahmani E *et al.* Accurate estimation of cell composition in bulk expression through robust integration of single-cell information. *Nat Commun* 2020;11:1971. <https://doi.org/10.1038/s41467-020-15816-6>
- Jiménez-Sánchez A, Cast O, Miller ML. Comprehensive benchmarking and integration of tumor microenvironment cell estimation methods. *Cancer Res* 2019;79:6238–46. <https://doi.org/10.1158/0008-5472.CAN-18-3560>
- Jin H, Liu Z. A benchmark for RNA-seq deconvolution analysis under dynamic testing environments. *Genome Biol* 2021;22:102. <https://doi.org/10.1186/s13059-021-02290-6>
- Jorge NAN, Cruz JGV, Pretti MAM *et al.* Poor clinical outcome in metastatic melanoma is associated with a microRNA-modulated immunosuppressive tumor microenvironment. *J Transl Med* 2020;18:56. <https://doi.org/10.1186/s12967-020-02235-w>
- Kang K, Meng Q, Shats I *et al.* CDSeq: a novel complete deconvolution method for dissecting heterogeneous samples using gene expression data. *PLoS Comput Biol* 2019;15:e1007510. <https://doi.org/10.1371/journal.pcbi.1007510>
- Lopez D, Montoya D, Ambrose M *et al.* SaVanT: a web-based tool for the sample-level visualization of molecular signatures in gene expression profiles. *BMC Genomics* 2017;18:824. <https://doi.org/10.1186/s12864-017-4167-7>
- Lowe R, Rakyán VK. Correcting for cell-type composition bias in epigenome-wide association studies. *Genome Med* 2014;6:23. <https://doi.org/10.1186/gm540>
- Macosko EZ, Basu A, Satija R *et al.* Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell* 2015;161:1202–14. <https://doi.org/10.1016/j.cell.2015.05.002>
- Maden SK, Kwon SH, Huuki-Myers LA *et al.* Challenges and opportunities to computationally deconvolve heterogeneous tissue with varying cell sizes using single-cell RNA-sequencing datasets. *Genome Biol* 2023;24:288. <https://doi.org/10.1186/s13059-023-03123-4>
- Matos LL, Truffelli DC, de Matos MG, da Silva Pinhal MA. Immunohistochemistry as an important tool in biomarkers detection and clinical practice. *Biomarker Insights* 2010;5:9–20. <https://doi.org/10.4137/bmi.s2185>
- Menden K, Marouf M, Oller S *et al.* Deep learning-based cell composition analysis from tissue expression profiles. *Sci Adv* 2020;6:eaba2619. <https://doi.org/10.1126/sciadv.aba2619>
- Monaco G, Lee B, Xu W *et al.* RNA-Seq signatures normalized by mRNA abundance allow absolute deconvolution of human immune cell types. *Cell Rep* 2019;26:1627–40.e7. <https://doi.org/10.1016/j.celrep.2019.01.041>
- Moreno P, Fexova S, George N *et al.* Expression atlas update: gene and protein expression in multiple species. *Nucleic Acids Res* 2022;50:D129–D140. <https://doi.org/10.1093/nar/gkab1030>
- Nadel BB, Oliva M, Shou BL *et al.* Systematic evaluation of transcriptomics-based deconvolution methods and references using thousands of clinical samples. *Brief Bioinform* 2021;22:bbab265. <https://doi.org/10.1093/bib/bbab265>
- Newman AM, Liu CL, Green MR *et al.* Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods* 2015;12:453–7. <https://doi.org/10.1038/nmeth.3337>
- Patrick E, Taga M, Ergun A *et al.* Deconvolving the contributions of cell-type heterogeneity on cortical gene expression. *PLoS Comput Biol* 2020;16:e1008120. <https://doi.org/10.1371/journal.pcbi.1008120>
- Picelli S, Björklund ÅK, Faridani OR *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat Methods* 2013;10:1096–8. <https://doi.org/10.1038/nmeth.2639>
- Repsilber D, Kern S, Telaar A *et al.* Biomarker discovery in heterogeneous tissue samples -taking the in-silico deconvolution approach. *BMC Bioinform* 2010;11:27. <https://doi.org/10.1186/1471-2105-11-27>
- Robinson MD, Smyth GK. Moderated statistical tests for assessing differences in tag abundance. *Bioinformatics* 2007;23:2881–7. <https://doi.org/10.1093/bioinformatics/btm453>
- Segerstolpe Å, Palasantza A, Eliasson P *et al.* Single-Cell transcriptome profiling of human pancreatic islets in health and type 2 diabetes. *Cell Metab* 2016;24:593–607. <https://doi.org/10.1016/j.cmet.2016.08.020>
- Shen-Orr SS, Tibshirani R, Khatri P *et al.* Cell type-specific gene expression differences in complex tissues. *Nat Methods* 2010;7:287–9. <https://doi.org/10.1038/nmeth.1439>
- Song Y, Miao Z, Brazma A *et al.* Benchmarking strategies for cross-species integration of single-cell RNA sequencing data. *Nat Commun* 2023;14:6495. <https://doi.org/10.1038/s41467-023-41855-w>

- Sturm G, Finotello F, Petitprez F *et al.* Comprehensive evaluation of transcriptome-based cell-type quantification methods for immunoncology. *Bioinformatics* 2019;35:i436–i445. <https://doi.org/10.1093/bioinformatics/btz363>
- Sutton GJ, Poppe D, Simmons RK *et al.* Comprehensive evaluation of deconvolution methods for human brain gene expression. *Nat Commun* 2022;13:1358. <https://doi.org/10.1038/s41467-022-28655-4>
- Taube JM, Galon J, Sholl LM *et al.* Implications of the tumor immune microenvironment for staging and therapeutics. *Mod Pathol* 2018;31:214–34. <https://doi.org/10.1038/modpathol.2017.156>
- The Tabula Sapiens Consortium*. The tabula sapiens: a multiple-organ, single-cell transcriptomic atlas of humans. *Science* 2022;376:eabl4896. <https://doi.org/10.1126/science.abl4896>
- Teschendorff AE, Breeze CE, Zheng SC *et al.* A comparison of reference-based algorithms for correcting cell-type heterogeneity in Epigenome-Wide association studies. *BMC Bioinformatics* 2017;18:105. <https://doi.org/10.1186/s12859-017-1511-5>
- The International Cancer Genome Consortium. International network of cancer genome projects. *Nature* 2010;993–8.
- Tsoucas D, Dong R, Chen H *et al.* Accurate estimation of cell-type composition from gene expression data. *Nat Commun* 2019;10:2975. <https://doi.org/10.1038/s41467-019-10802-z>
- Wang X, Park J, Susztak K *et al.* Bulk tissue cell type deconvolution with multi-subject single-cell expression reference. *Nat Commun* 2019;10:380. <https://doi.org/10.1038/s41467-018-08023-x>
- White BS, de Reyniès A, Newman AM *et al.* Community assessment of methods to deconvolve cellular composition from bulk gene expression. *bioRxiv* 2022.06.03.494221. <https://doi.org/10.1101/2022.06.03.494221>
- Wilk AJ, Rustagi A, Zhao NQ *et al.* A single-cell atlas of the peripheral immune response in patients with severe COVID-19. *Nat Med* 2020;26:1070–6. <https://doi.org/10.1038/s41591-020-0944-y>
- Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol* 2018;19:15. <https://doi.org/10.1186/s13059-017-1382-0>
- Zhang Y, Cheng S, Zhang M. High-infiltration of tumor-associated macrophages predicts unfavorable clinical outcome for node-negative breast cancer. *PLoS One* 2013;8:e76147. <https://doi.org/10.1371/journal.pone.0076147>
- Zheng GXY, Terry JM, Belgrader P *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat Commun* 2017;8:14049. <https://doi.org/10.1038/ncomms14049>
- Zilionis R, Nainys J, Veres A *et al.* Single-cell barcoding and sequencing using droplet microfluidics. *Nat Protoc* 2017;12:44–73. <https://doi.org/10.1038/nprot.2016.154>