

1 Haplotype-resolved genome assembly of the tetraploid potato cultivar Désirée

2

3 Tim Godec^{*1,2}, Sebastian Beier³, Natalia Yaneth Rodriguez-Granados⁴, Rashmi Sasidharan⁴,

4 Lamis Abdelhakim⁵, Markus Teige⁶, Björn Usadel^{3,7}, Kristina Gruden¹, Marko Petek¹

5

6

7 ¹ National Institute of Biology, Department of Biotechnology and Systems Biology, Ljubljana,
8 Slovenia

9 ² Jožef Stefan International Postgraduate School, Ljubljana, Slovenia

10 ³ Institute of Bio- and Geosciences (IBG-4 Bioinformatics), Bioeconomy Science Center
11 (BioSC), CEPLAS, Forschungszentrum Jülich GmbH, Jülich, Germany

12 ⁴ Plant Stress Resilience, Institute of Environmental Biology, Utrecht University, Utrecht, The
13 Netherlands

14 ⁵ PSI (Photon Systems Instruments), Drásov, Czech Republic

15 ⁶ Molecular Systems Biology (MOSYS), Department of Functional and Evolutionary Ecology,
16 University Vienna, Vienna, Austria

17 ⁷ Faculty of Mathematics and Natural Sciences, Institute for Biological Data Science, Cluster
18 of Excellence on Plant Sciences (CEPLAS), Heinrich Heine University Düsseldorf,
19 Düsseldorf, Germany

20

21

22 * corresponding author

23 corresponding author email: tim.godec@nib.si

24 Abstract

25 Cultivar Désirée is an important model for potato functional genomics studies to assist
26 breeding strategies. Here, we present a haplotype-resolved genome assembly of Désirée,
27 achieved by assembling PacBio HiFi reads and Hi-C scaffolding, resulting in a
28 high-contiguity chromosome-level assembly. We implemented a comprehensive annotation
29 pipeline incorporating gene models and functional annotations from the *Solanum tuberosum*
30 Phureja DM reference genome alongside RNA-seq reads to provide high-quality gene and
31 transcript annotations. Additionally, we provide a genome-wide DNA methylation profile
32 using Oxford Nanopore reads, enabling insights into potato epigenetics. The assembled
33 genome, annotations, methylation and expression data are visualised in a publicly
34 accessible genome browser (<https://desiree.nib.si>), providing a valuable resource for the
35 potato research community.

36 Background & Summary

37

38 Potato (*Solanum tuberosum*) is one of the most important and widely cultivated crops
39 worldwide, with a significant role in global food security and agricultural research. Despite its
40 significance, many studies still rely on the genome of the double monoploid (DM) clone of
41 group Phureja DM1–3 516 R44^{1,2} which lacks a substantial portion of the gene repertoire
42 and variability found in cultivated tetraploid potato varieties.

43

44 The potato cultivar Désirée is a red-skinned late-season potato variety, originally bred in the
45 Netherlands in 1962 by crossing parent cultivars Urgenta and Depesche (Potato Pedigree
46 Database)³. It is still cultivated due to its favourable agronomic traits, such as predictable
47 yields and high tolerance to drought and some pathogens⁴. It has also been used in
48 breeding programs, yet a genome assembly for the Désirée cultivar has not been available.
49 In research, it has been propagated in tissue cultures, and used for genetic manipulation
50 including gene overexpression⁵, gene silencing⁶, and Crispr-Cas gene editing⁷.

51

52 Although haplotype-resolved genome assemblies are becoming common in diploid
53 organisms, the high heterozygosity rate, extensive repeat content, and the autopolyploid
54 nature of cultivated potatoes still present significant challenges for generating high-quality
55 haplotype-resolved assemblies. Currently, five haplotype-resolved genomes of autotetraploid
56 potato cultivars are publicly available^{8–12} as well as several phased diploid genomes^{13–15}. The
57 recently published haplotype-resolved tetraploid potato assemblies rely on labour-intensive
58 techniques such as single-pollen sequencing¹⁰ or the use of parental and crossing material¹¹,
59 which may not always be available.

60

61 Adding to existing publicly available genomes, we provide a reference quality (CRAQ overall
62 AQI of 97.5) haplotype-resolved genome assembly of the tetraploid cultivar Désirée,
63 assembled using solely PacBio HiFi and Illumina Hi-C data. Our assembly is accompanied
64 by a comprehensive structural and functional gene annotation reaching 99.4 % BUSCO
65 completeness for Solanaceae, accompanied by orthology to DM genes. For the potato
66 research community, we provide an online resource featuring a genome browser and
67 downloadable genomic assembly and annotation files, providing a valuable tool for studies
68 involving allele-specific expression or promoter analysis.

69 Methods

70 Sample preparation and sequencing

71 Leaves from 4-week old *S. tuberosum* cv. Désirée plants were collected and flash-frozen.
 72 High molecular weight genomic DNA (HMW gDNA) used for PacBio HiFi, Illumina and
 73 Oxford Nanopore Technologies (ONT) sequencing was extracted from the leaf tissues using
 74 a modified CTAB method¹⁶. The concentration and quality of the extracted DNA were
 75 assessed using a NanoDrop spectrophotometer.

76

77 *PacBio HiFi*

78 HMW gDNA was sent to National Genomics Infrastructure (NGI) Sweden for library
 79 preparation and sequencing on the PacBio Sequel II platform. We obtained 79.4 Gbp of raw
 80 data, consisting of 4.1 million reads.

81

82 *Illumina Hi-C*

83 Leaves from 4-week old *S. tuberosum* cv. Désirée plants were collected, flash-frozen in
 84 liquid nitrogen and ground using mortar and pestle. Hi-C library prep using the Omni-C kit
 85 (Dovetail Genomics) and sequencing were performed on an Illumina NovaSeq 6000 platform
 86 by NGI Sweden. Sequencing generated 2018.4 million paired-end (2 × 150 bp) reads.

87

88 *ONT*

89 The HMW gDNA was used for ONT DNA library prep using the SQK-LSK110 kit and
 90 sequenced on a MinION using the FLO-MIN106 flow cell. Reads were basecalled using
 91 Dorado (v0.7.2) with the model dna_r9.4.1_e8_sup@v3.3 which generated 5.8 Gbp. The
 92 reads with methylation-related tags were converted to bedMethyl format using modkit
 93 (v0.4.1).

94

95 *Illumina short reads*

96 Illumina short-read library was constructed from the HMW gDNA and sequenced on Illumina
 97 NextSeq 2000 by ELIXIR Slovenia node to generate 150 bp paired-end reads. The
 98 short-read sequencing generated approximately 138 Gbp of raw data, consisting of 460.1
 99 million paired-end (2 × 150 bp) reads.

100

101 Genome size and heterozygosity estimation

102 The genome characteristics of *S. tuberosum* cv. Désirée, including genome size,
 103 heterozygosity, and repeat content, were estimated using Illumina short-read data and a
 104 k-mer based approach. A 21-mer frequency distribution was generated with Jellyfish
 105 (v2.2.10), and the genome's key features were inferred using GenomeScope2 (v2.0). The
 106 haploid genome size was estimated at 669.6 Mbp, with a heterozygosity rate estimated at
 107 3.8–5.7%.

108

De novo genome assembly, Hi-C scaffolding and quality assessment

111

112 PacBio HiFi and Illumina Hi-C reads were initially assembled into four sets of
113 haplotype-resolved contigs using Hifiasm (v0.19.8-r603)^{17–19}. Hifiasm primary unitigs were
114 searched against DM genome assembly with blastn (v2.5.0)²⁰ and best matches were
115 visualised on Graphical Fragment Assembly with Bandage (v0.8.1, Fig. 1a)²¹. We performed
116 quality control of the contigs using Merqury (v1.3, Fig. 1b)²² k-mer spectra and BUSCO
117 completeness scores (v5.4.7, solanales_odb10 dataset)²³. The length of haplotype draft
118 assemblies ranged from 761.6 Mbp to 888.4 Mbp with contig N50 sizes ranging from
119 7.0 Mbp to 13.7 Mbp (Table 1).

120

121 Contigs identified as contaminants were removed based on blastn (v0.8.1) searches against
122 a custom-built contaminant database, which includes *Solanum* plastid and mitochondrial
123 sequences and bacterial NCBI RefSeq sequences.

124

125 Decontaminated scaffolds were anchored to chromosomes by mapping Hi-C reads to each
126 haplotype set separately following the manufacturer's recommended pipeline for Omni-C
127 data (<https://omni-c.readthedocs.io>). Briefly, Hi-C reads were mapped using BWA-MEM
128 (v0.7.17-r1188)²⁴ then the mappings were parsed with *pairtools* (v0.3.0)²⁵ followed by
129 samtools (v1.3.1)²⁶ to identify and extract valid pairs. Valid pairs were used to anchor and
130 orient scaffolds into chromosomes using YaHS (v1.2a.1)²⁷ and Juicebox Assembly Tools
131 (v2.17.00)^{28,29}.

132

133 Chromosomes 11 and 12 of haplotype 4 lacked ~20 Mbp and ~30 Mbp part of the
134 pericentromeric region, respectively, and haplotype 1 contained two additional unplaced
135 scaffolds (scaffold_22 and scaffold_23). Alignment of these scaffolds to reference genome
136 (DM v6.1) and inspection of Hi-C contacts suggested that these scaffolds are the missing
137 regions of chromosomes 11 and 12 in haplotype 4. Therefore, we remapped Hi-C reads and
138 incorporated these two scaffolds in haplotype 4 using Juicebox Assembly Tools (v2.17.00).

139

140 The final scaffolded assembly size amounts to 3.3 Gbp, with individual haplotypes ranging
141 between 762 and 888 Mb. As expected, one haplotype is highly similar to the DM haplotype,
142 whereas other haplotypes can be more dissimilar (Fig. 1c). A comparison of Merqury k-mer
143 spectra between the initial contigs and the scaffolded chromosomes (Fig. 1a) reveals that
144 many apparent duplications in the contigs are resolved during scaffolding. A small proportion
145 of sequences remains missing from the chromosomes and those can be found in the whole
146 genome FASTA.

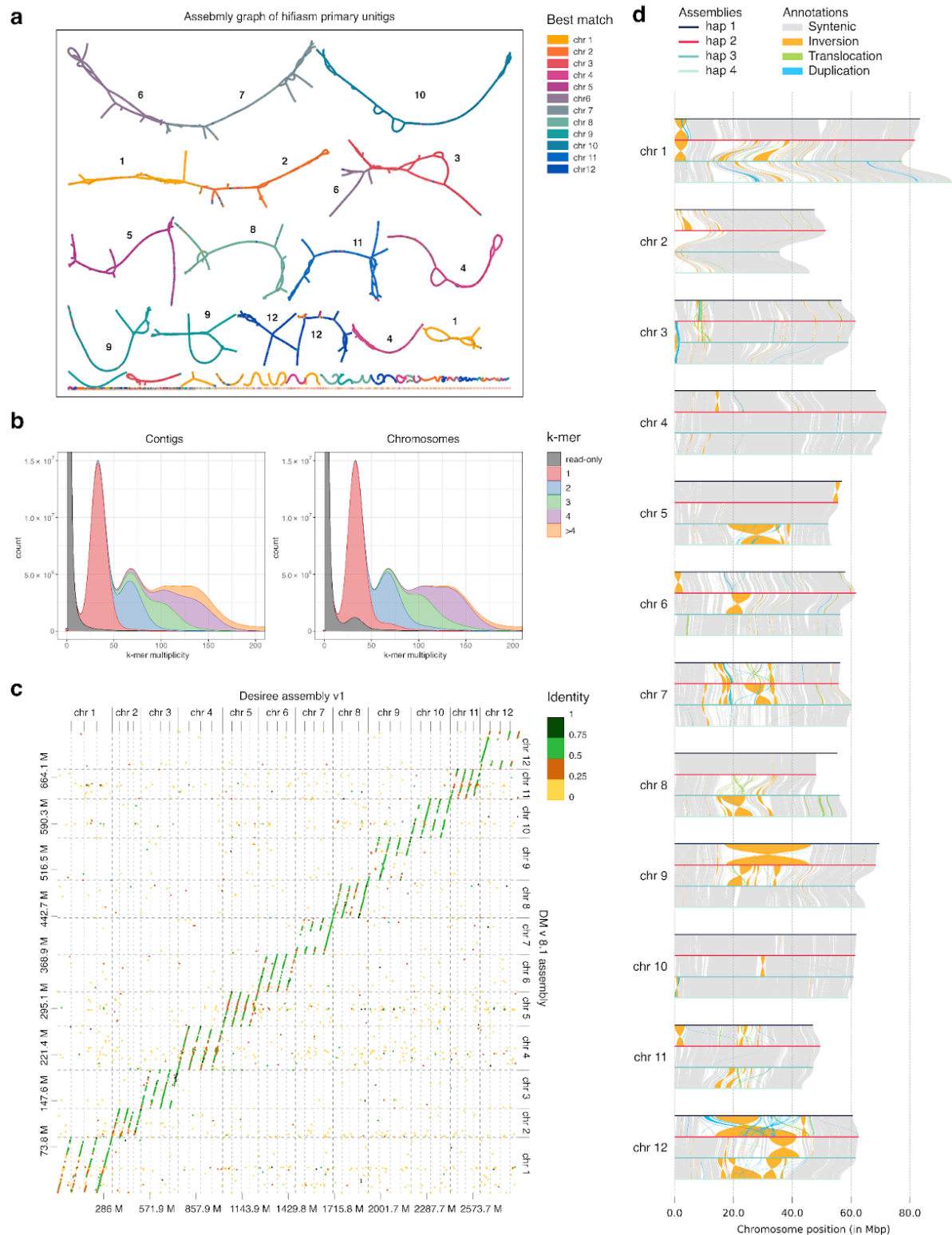
147

148 The haplotype assemblies were sequentially aligned using minimap2 (v2.28) and analyzed
149 with SyRi (1.7.0) to identify syntenic regions and structural rearrangements which were
150 visualized using plotsr (v1.1.1, Fig. 1d).

151

| | haplotype 1 | haplotype 2 | haplotype 3 | haplotype 4 | all haplotypes |
|-------------------------------|-------------|-------------|-------------|-------------|----------------|
| Genome length (Mb) | 888.4 | 862.7 | 761.6 | 858.5 | 3371.2 |
| GC content (%) | 35.31 | 35.27 | 35.12 | 35.47 | 35.3 |
| Contig N50 (Mb) | 11.5 | 13.7 | 11.7 | 7.0 | 10.8 |
| Number of contigs | 1126 | 867 | 1048 | 2695 | 5736 |
| Chromosome length (Mb) | 721.9 | 729.9 | 698.5 | 709.4 | 2859.6 |
| Scaffold N50 (Mb) | 56.9 | 61.4 | 60.1 | 57.1 | 58.0 |
| Number of scaffolds | 705 | 496 | 523 | 1350 | 3074 |
| Complete BUSCO (%) | 96.2% | 96.1% | 96.6% | 95.7% | 99.6% |
| Size of repeat sequences (Mb) | 514.2 | 534.1 | 489.3 | 503.6 | 2041.2 |
| Total gene number | 76903 | 81184 | 75816 | 75550 | 309453 |

152 **Table 1.** Summary of the four haplotypes of the Désirée genome assembly.



153

Fig. 1 General characteristics of Désirée genome assembly **a)** Assembly graph of primary
 units coloured by best match to DM chromosomes (also designated with numbers on the
 graph). **b)** Merqury k-mer spectra for initial contigs and scaffolded chromosomes. The k = 21
 was used. K-mers are categorized as read-only (grey), unique (red), and shared (blue,
 green, purple, orange). Peaks corresponding to higher multiplicities indicate the presence of
 highly repeated k-mers. **c)** Dot plot comparing cv. Désirée chromosome-anchored contigs

with DM v8.1 chromosomes. The colour designates contig identity. **d)** Genomic synteny of cv. Désirée haplotype-resolved assembly.

Genome annotation

Repeat elements in the *S. tuberosum* cv. Désirée genome were identified using the Extensive *de novo* TE Annotator (EDTA, v2.2.1)³⁰. Repetitive sequences cover 489 - 534 Mbp per haplotype, representing more than 70% of the genome (Table 2).

The prediction of protein-coding genes in the assembled *S. tuberosum* cv. Désirée was determined using five complementary approaches: *de novo*, homology-based, transcriptome-based, deep-learning, and reference-based predictions (Fig. 2).

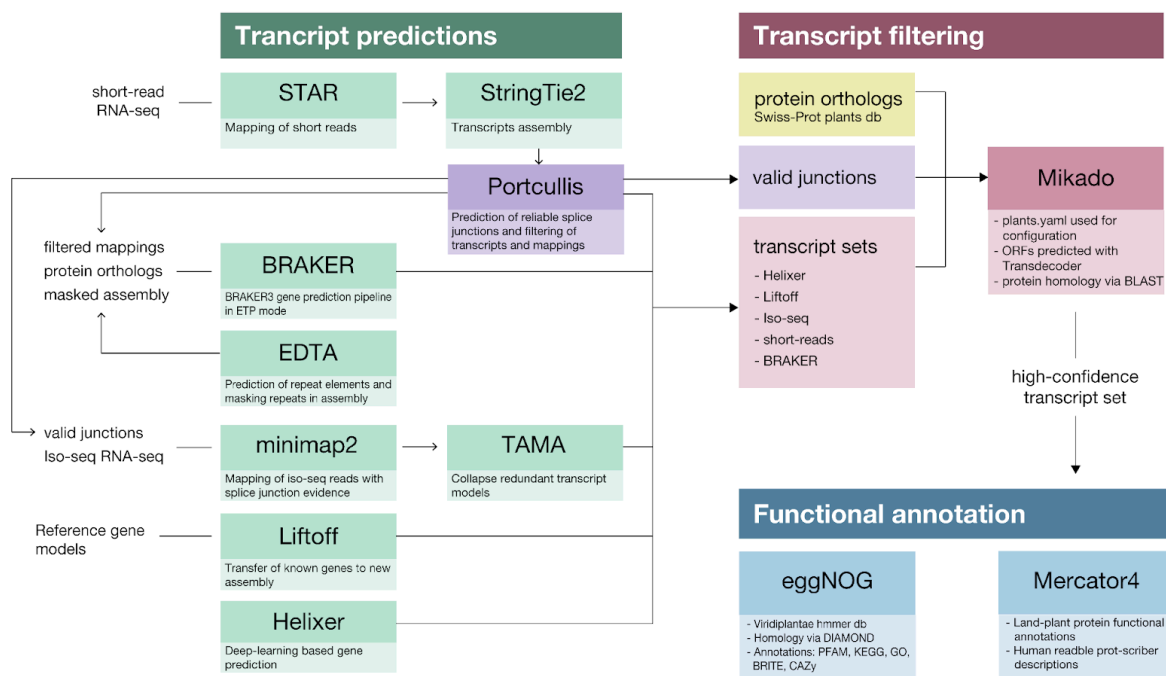


Fig. 2 Workflow overview of *S. tuberosum* cv. Désirée genome annotation.

For transcriptome-based prediction, two methods were applied for short reads and Iso-Seq reads, respectively. Short reads from multiple tissues were aligned to each haplotype using STAR (2.7.10a)³¹, and transcripts were assembled with StringTie2 (v2.2.1)³², followed by Portcullis (v1.2.4)³³ for junction validation. Iso-Seq reads from five *S. tuberosum* cultivars were mapped to both haplotypes using minimap2 (v2.28)³⁴, and transcripts were generated using IsoQuant (v3.3.1)³⁵ and TAMA Collapse (tc_version_date_2023_03_28)³⁶.

BRAKER3 (v3.0.8)³⁷ was used in ETP mode to predict gene models by integrating *de novo*, homology-based, and transcriptome-based predictions. Repeat masking of the assembly was performed with RepeatMasker (v4.1.2), using EDTA annotations. Protein sequences from OrthoDB (green plant orthologs) were provided as evidence, and short-read STAR alignments with invalid junctions removed were included.

188

189 Helixer (v0.3.3)^{38,39} was used for deep-learning-based gene prediction via its web interface
 190 (https://www.plabipd.de/helixer_main.html). Gene models from the *S. tuberosum* reference
 191 genome (DM v6.1, UniTato annotation) were transferred to the Désirée assembly using
 192 Liftoff (v1.6.3)⁴⁰. All five transcript or gene model sets were consolidated using Mikado
 193 (v2.3.4)⁴¹ to generate a non-redundant set of transcripts. Protein-coding gene completeness
 194 was assessed using BUSCO (Table 2, v5.4.7, solanales_odb10 dataset) and OMArk (v0.3.0,
 195 omamer v2.0.2)⁴².

196

197 The predicted protein-coding genes were functionally annotated using EggNOG Mapper
 198 (v2.1.11)⁴³ with the EggNOG database (version 5.0.2)⁴⁴ for the Viridiplantae subset. This
 199 included categories such as gene names, Gene Ontologies (GOs), enzyme functions (EC),
 200 and KEGG pathways, reactions, and modules, along with CAZy families, PFAM domains,
 201 and more. Additionally, functional land-plant protein annotations were predicted using
 202 Mercator4 (v7)⁴⁵ via the web platform (https://www.plabipd.de/mercator_main.html).
 203 Annotations from EggNOG and Mercator4 were combined into the final GFF3 annotation file.

204

205 Orthologous groups between haplotypes and UniTato genes were identified using
 206 OrthoFinder (v2.5.5)⁴⁶. Across haplotypes, 55.3% of orthogroups contained genes from all
 207 four haplotypes, 22.9% from three haplotypes, 19.2% from two haplotypes, and 2.7% from a
 208 single haplotype. When comparing the Désirée annotation to UniTato, 17.24% of genes were
 209 specific to the Désirée annotation.

210

| Type | haplotype 1 | haplotype 2 | haplotype 3 | haplotype 4 |
|-----------------------------------|-------------------|-------------------|-------------------|-------------------|
| Repeat elements | | | | |
| DNA | 46.6 Mbp (6.4%) | 54.1 Mbp (7.4%) | 44.6 Mbp (6.4%) | 45.9 Mbp (6.5%) |
| Helitron | 36.1 Mbp (5.0%) | 38.0 Mbp (5.2%) | 33.6 Mbp (4.8%) | 42.3 Mbp (6.0%) |
| LINE | 12.4 Mbp (1.7%) | 8.1 Mbp (1.1%) | 7.5 Mbp (1.1%) | 8.1 Mbp (1.1%) |
| LTR | 176.5 Mbp (24.4%) | 188.0 Mbp (25.8%) | 165.9 Mbp (23.7%) | 193.6 Mbp (27.3%) |
| LTR/Copia | 16.8 Mbp (2.3%) | 19.3 Mbp (2.6%) | 20.3 Mbp (2.9%) | 23.9 Mbp (3.4%) |
| LTR/Gypsy | 136.2 Mbp (18.9%) | 133.6 Mbp (18.3%) | 130.0 Mbp (18.6%) | 102.8 Mbp (14.5%) |
| MITE | 11.9 Mbp (1.6%) | 10.2 Mbp (1.4%) | 13.0 Mbp (1.9%) | 10.6 Mbp (1.5%) |
| Other | 72.8 Mbp (10.1%) | 76.2 Mbp (10.4%) | 69.7 Mbp (10.0%) | 71.3 Mbp (10.1%) |
| SINE | 5.1 Mbp (0.7%) | 6.6 Mbp (0.9%) | 4.7 Mbp (0.7%) | 4.9 Mbp (0.7%) |
| Total | 514.2 Mbp (71.2%) | 534.1 Mbp (73.2%) | 489.3 Mbp (70.1%) | 503.6 Mbp (71.0%) |
| Protein-coding genes | | | | |
| Total gene number | 76903 | 81184 | 75816 | 75550 |
| Mean gene length (bp) | 1695.85 | 1610.97 | 1687.71 | 1677.79 |
| Mean CDS length (bp) | 1062.59 | 1032.74 | 1060.23 | 1061.68 |
| Mean exon number | 5.28 | 5.04 | 5.31 | 5.28 |
| Mean intron number | 4.28 | 4.04 | 4.31 | 4.28 |
| Complete BUSCO (%) | 94.1% | 93.3% | 95.4% | 93.7% |
| Single Omark HOGs | 82.9% | 82.5% | 84.3% | 82.8% |
| Duplicated Omark HOGs | 11.6% | 11.6% | 11.5% | 11.9% |
| Missing Omark HOGs | 5.5% | 5.9% | 4.2% | 5.4% |
| Mercator4 proteins annotated (%) | 93.5% | 93.5% | 93.7% | 93.5% |
| Mercator4 proteins classified (%) | 50.5% | 46.5% | 50.7% | 50.0% |
| Mercator4 bins occupied (%) | 94.2% | 93.9% | 94.6% | 94.3% |

211 **Table 2.** Summary of genome annotations for each haplotype.

212 Data Records

213 The raw sequencing data, including Illumina Hi-C, Illumina paired-end, PacBio HiFi, and
214 ONT reads, have been deposited at the National Center for Biotechnology Information
215 (NCBI) Sequence Read Archive (SRA) under BioProject number PRJNA1185028. Plastid,

mitochondrial and bacterial sequences used for removal of contaminant contigs were downloaded from NCBI RefSeq release 218. Transcriptomic data used for gene annotation was downloaded from public repositories: SRA under accessions PRJNA1192223, PRJNA1186376, PRJNA718240, PRJNA803222, PRJNA1209787 and PRJNA1191209; the Gene Expression Omnibus (GEO) under accession GSE232028; and the National Genomics Data Center (NGDC) under accession CRA006012. Existing gene models used in the gene annotation pipeline were downloaded from <https://unitato.nib.si> and <https://spuddb.uga.edu>. The genome assemblies of the four haplotypes have been submitted to NCBI GenBank under the BioProject accessions PRJNA1196677, PRJNA1196678, PRJNA1196679 and PRJNA1196680. The assembled genome, including annotations, methylation profile and identified orthologs, is hosted in a Zenodo repository under DOI: 10.5281/zenodo.14609304 and is also accessible via an interactive genome browser at <https://desiree.nib.si>.

Technical Validation

We assessed the assembly quality and completeness using DNA sequencing read mapping, CRAQ, BUSCO analysis, and Merqury k-mer based evaluation. Illumina reads were mapped with BWA (v0.7.17), while PacBio and ONT reads were aligned using minimap2 (v2.28). Mapping rates were 99.90%, 100.00%, and 99.74% for Illumina paired-end, PacBio, and ONT reads, respectively. CRAQ (v1.0.9)⁴⁷ analysis of PacBio and Illumina mappings yielded a regional AQI of 96.3 and an overall AQI of 97.5, classifying the assembly as reference quality (AQI > 90). Assembly completeness was assessed with BUSCO (v5.4.7) using the solanales_odb10 lineage database, identifying 5930 (99.6%) of the 5950 BUSCO orthologous groups in both the whole genome and chromosome-only assemblies (Table 1). Merqury (v1.3) analysis, using a Meryl (v1.3) database constructed from Illumina reads, estimated genome completeness at 98.57% for the whole genome and 95.73% for the chromosomes. The estimated QV values were 54.30 and 58.53 for the whole genome and chromosomes, respectively.

Completeness of gene annotation was assessed using OMArk (v0.3.0, omamer v2.0.2), BUSCO (v5.4.7) and Mercator4 (v7). OMArk analysis demonstrated that our annotation captured 94.1%–94.6% of Hierarchical Orthologous Groups (HOGs) per haplotype, with duplication rates ranging from 11.5% to 11.9% (Fig. 3a). When combining genes from all haplotypes, the proportion of complete HOGs reaches 99.3%, meaning that not all conserved genes are present in all haplotypes. Similarly, BUSCO analysis reported a haplotype completeness range of 93.3%–95.4% (Table 2), while the whole genome annotation achieved 99.4% completeness. Protein classification via Mercator4 revealed that 93.9%–94.6% of Mercator bins were occupied per haplotype, increasing to 97.5% when combining all proteins (Table 2). As expected, the Mercator bin with the largest proportion of missing proteins was associated with clade-specific metabolism (Fig. 3b). Additionally, the classified proteins showed no significant deviation from the median protein length, confirming consistency in annotation quality (Fig. 3c).

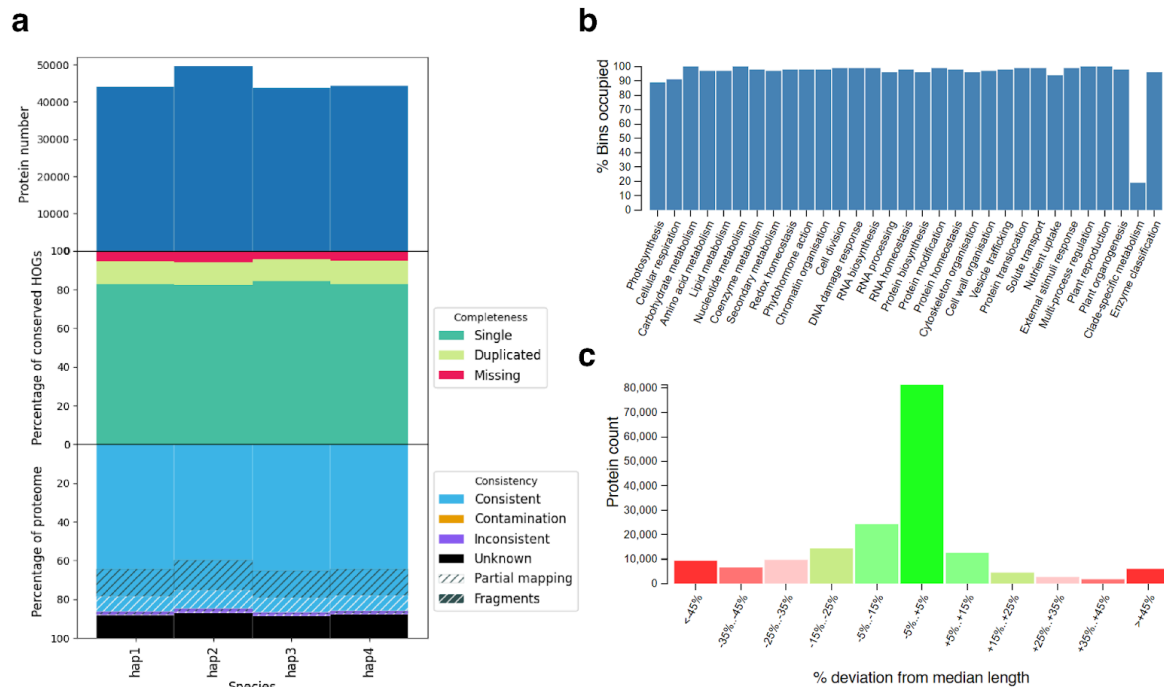


Fig. 3 Validation of gene annotation. **a)** OMArk quality assessment showing consistency, completeness and count of proteins across all four haplotypes. **b)** Histogram showing the percentage of Mercator4 functional bins occupied by the Désirée proteins. **c)** Histogram displaying the distribution of proteins grouped by their percentage deviation from the median protein length.

Usage Notes

The presented Désirée genome assembly is of high contiguity, completeness and phasing quality and presents a valuable resource for haplotype-aware transcriptomics, proteomics and epigenomics analyses. The transfer of UniTato annotations⁴⁸ provides translation of gene identifiers from the DM to the Désirée genome. The RNA-seq datasets used to supplement gene model annotation are predominantly from mature leaf and root tissue, thus genes specifically expressed in other tissue and developmental stages may not be fully captured in the current annotation.

The genome was produced from a plant propagated in tissue culture for over a decade. A recent pangenome study⁴⁹ found that *in vitro* propagated plants of the *Solanum* section Petota have greater numbers of TEs in their genomes. While this seems to hold for LTR elements and DNA transposons in the Désirée genome, overall TE expansion is not evident. Examining the DNA methylation profile available in the Désirée genome browser might provide more insight into specific transposable element expansion in this cultivar.

Recently, efforts were made to generate potato pangenomes^{9,49}. However, the number of included phased tetraploid genomes is still limited. Including Désirée and more phased tetraploid genomes will improve the completeness of potato pangenome. This will bridge knowledge gaps in potato genomics and give potato breeders a powerful toolkit for developing more resilient and productive cultivars.

284 Code Availability

285 The code, scripts and command-line tool commands used for genome assembly, annotation
286 and quality control are freely available in the GitHub repository
287 <https://github.com/NIB-SI/desiree-genome>.

288

289 Acknowledgement

290 This work benefits from resources and services provided by ELIXIR, a distributed
291 infrastructure for life science data, funded by national governments and the European
292 Commission, particularly the Elixir-SI node for performing Illumina paired-end sequencing.

293

294 Funding for this work was provided by the European Union's Horizon 2020 research and
295 innovation programme project ADAPT (grant agreement No GA 2020 862-858), Slovenian
296 Research and Innovation Agency (ARIS) project grants P4-0165, P4-0431, and J4-3089. SB
297 and BU are supported by the German Federal Ministry of Education and Research (BMBF)
298 in the frame of the German Network for Bioinformatics Infrastructure (de.NBI).

299

300 Author contributions

301 **TG**: Methodology, Data curation, Investigation, Visualization, Writing - Original Draft. **SB**:
302 Investigation, Writing - Review & Editing. **BU**: Writing - Review & Editing. **NYRG**: Resources,
303 Writing - Review & Editing. **RS**: Resources, Writing - Review & Editing. **LA**: Resources,
304 Writing - Review & Editing. **MT**: Funding acquisition, Writing - Review & Editing. **KG**: Funding
305 acquisition, Conceptualization, Writing - Review & Editing. **MP**: Conceptualization, Validation,
306 Resources, Supervision, Project administration, Writing - Review & Editing.

307

308 Competing interests

309 The author(s) declare no competing interests.

References

1. Yang, X. *et al.* The gap-free potato genome assembly reveals large tandem gene clusters of agronomical importance in highly repeated genomic regions. *Molecular Plant* **16**, 314–317 (2023).
2. Pham, G. M. *et al.* Construction of a chromosome-scale long-read reference genome assembly for potato. *GigaScience* **9**, giaa100 (2020).
3. van Berloo, R., Hutten, R. C. B., van Eck, H. J. & Visser, R. G. F. An Online Potato Pedigree Database Resource. *Potato Res.* **50**, 45–57 (2007).
4. The European Cultivated Potato Database.
<https://www.europotato.org/varieties/view/Desiree-E>.
5. Tomaž, Š. *et al.* A mini-TGA protein modulates gene expression through heterogeneous association with transcription factors. *Plant Physiology* **191**, 1934–1952 (2023).
6. Halim, V. A. *et al.* PAMP-induced defense responses in potato require both salicylic acid and jasmonic acid. *The Plant Journal* **57**, 230–242 (2009).
7. Lukan, T. *et al.* CRISPR/Cas9-mediated fine-tuning of miRNA expression in tetraploid potato. *Horticulture Research* **9**, uhac147 (2022).
8. Bao, Z. *et al.* Genome architecture and tetrasomic inheritance of autotetraploid potato. *Molecular Plant* **15**, 1211–1226 (2022).
9. Hoopes, G. *et al.* Phased, chromosome-scale genome assemblies of tetraploid potato reveal a complex genome, transcriptome, and predicted proteome landscape underpinning genetic diversity. *Molecular Plant* **15**, 520–536 (2022).
10. Sun, H. *et al.* Chromosome-scale and haplotype-resolved genome assembly of a tetraploid potato cultivar. *Nat Genet* **54**, 342–348 (2022).
11. Serra Mari, R. *et al.* Haplotype-resolved assembly of a tetraploid potato genome using long reads and low-depth offspring data. *Genome Biology* **25**, 26 (2024).
12. Reyes-Herrera, P. H. *et al.* Chromosome-scale genome assembly and annotation of the tetraploid potato cultivar Diacol Capiro adapted to the Andean region. *G3 Genes|Genomes|Genetics* **14**, jkae139 (2024).
13. Freire, R. *et al.* Chromosome-scale reference genome assembly of a diploid potato clone derived from an elite variety. *G3 Genes|Genomes|Genetics* **11**, jkab330 (2021).
14. van Lieshout, N. *et al.* Solyntus, the New Highly Contiguous Reference Genome for Potato (*Solanum tuberosum*). *G3 Genes|Genomes|Genetics* **10**, 3489–3495 (2020).
15. Zhou, Q. *et al.* Haplotype-resolved genome analyses of a heterozygous diploid potato. *Nat Genet* **52**, 1018–1023 (2020).
16. Doyle, J. DNA extraction by using DTAB-CTAB procedures. *Phytochemical Bulletin* **19**, 11–17 (1987).
17. Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat Methods* **18**, 170–175 (2021).
18. Cheng, H. *et al.* Haplotype-resolved assembly of diploid genomes without parental data. *Nat Biotechnol* **40**, 1332–1335 (2022).
19. Cheng, H., Asri, M., Lucas, J., Koren, S. & Li, H. Scalable telomere-to-telomere assembly for diploid and polyploid genomes with double graph. *Nat Methods* **21**, 967–970 (2024).
20. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).
21. Wick, R. R., Schultz, M. B., Zobel, J. & Holt, K. E. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics* **31**, 3350–3352 (2015).
22. Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies. *Genome Biology* **21**, 245 (2020).
23. Manni, M., Berkeley, M. R., Seppey, M., Simão, F. A. & Zdobnov, E. M. BUSCO Update: Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic

- 363 Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. *Molecular Biology*
364 *and Evolution* **38**, 4647–4654 (2021).
- 365 24. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
366 *arXiv:1303.3997 [q-bio]* (2013).
- 367 25. Open2C *et al.* Pairtools: From sequencing data to chromosome contacts. *PLOS*
368 *Computational Biology* **20**, e1012164 (2024).
- 369 26. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *GigaScience* **10**, giab008
370 (2021).
- 371 27. Zhou, C., McCarthy, S. A. & Durbin, R. YaHS: yet another Hi-C scaffolding tool.
372 *Bioinformatics* **39**, btac808 (2023).
- 373 28. Dudchenko, O. *et al.* The Juicebox Assembly Tools module facilitates de novo assembly
374 of mammalian genomes with chromosome-length scaffolds for under \$1000. 254797
375 Preprint at <https://doi.org/10.1101/254797> (2018).
- 376 29. Juicebox Provides a Visualization System for Hi-C Contact Maps with Unlimited Zoom.
377 *Cell Systems* **3**, 99–101 (2016).
- 378 30. Ou, S. *et al.* Benchmarking transposable element annotation methods for creation of a
379 streamlined, comprehensive pipeline. *Genome Biology* **20**, 275 (2019).
- 380 31. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21
381 (2013).
- 382 32. Shumate, A., Wong, B., Pertea, G. & Pertea, M. Improved transcriptome assembly using
383 a hybrid of long and short reads with StringTie. *PLOS Computational Biology* **18**,
384 e1009730 (2022).
- 385 33. Mapleson, D., Venturini, L. & Swarbreck, D. El-CoreBioinformatics/portcullis.
386 El-CoreBioinformatics (2024).
- 387 34. Li, H. New strategies to improve minimap2 alignment accuracy. *Bioinformatics* **37**,
388 4572–4574 (2021).
- 389 35. Pribelski, A. D. *et al.* Accurate isoform discovery with IsoQuant using long reads. *Nat*
390 *Biotechnol* **41**, 915–918 (2023).
- 391 36. Kuo, R. I. *et al.* Illuminating the dark side of the human transcriptome with long read
392 transcript sequencing. *BMC Genomics* **21**, 751 (2020).
- 393 37. Gabriel, L. *et al.* BRAKER3: Fully automated genome annotation using RNA-seq and
394 protein evidence with GeneMark-ETP, AUGUSTUS, and TSEBRA. *Genome Res.* **34**,
395 769–777 (2024).
- 396 38. Holst, F. *et al.* Helixer–de novo Prediction of Primary Eukaryotic Gene Models
397 Combining Deep Learning and a Hidden Markov Model. 2023.02.06.527280 Preprint at
398 <https://doi.org/10.1101/2023.02.06.527280> (2023).
- 399 39. Stiehler, F. *et al.* Helixer: cross-species gene annotation of large eukaryotic genomes
400 using deep learning. *Bioinformatics* **36**, 5291–5298 (2021).
- 401 40. Shumate, A. & Salzberg, S. L. Liftoff: accurate mapping of gene annotations.
402 *Bioinformatics* **37**, 1639–1643 (2021).
- 403 41. Venturini, L., Caim, S., Kaithakottil, G. G., Mapleson, D. L. & Swarbreck, D. Leveraging
404 multiple transcriptome assembly methods for improved gene structure annotation.
405 *GigaScience* **7**, giy093 (2018).
- 406 42. Nevers, Y. *et al.* Quality assessment of gene repertoire annotations with OMArk. *Nat*
407 *Biotechnol* 1–10 (2024) doi:10.1038/s41587-024-02147-w.
- 408 43. Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J.
409 eggNOG-mapper v2: Functional Annotation, Orthology Assignments, and Domain
410 Prediction at the Metagenomic Scale. *Molecular Biology and Evolution* **38**, 5825–5829
411 (2021).
- 412 44. Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically
413 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids*
414 *Research* **47**, D309–D314 (2019).
- 415 45. MapMan4: A Refined Protein Classification and Annotation Framework Applicable to
416 Multi-Omics Data Analysis. *Molecular Plant* **12**, 879–892 (2019).
- 417 46. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative

- 418 genomics. *Genome Biology* **20**, 238 (2019).
- 419 47. Li, K., Xu, P., Wang, J., Yi, X. & Jiao, Y. Identification of errors in draft genome
420 assemblies at single-nucleotide resolution for quality assessment and improvement. *Nat*
421 *Commun* **14**, 6556 (2023).
- 422 48. Zagorščak, M. *et al.* Evidence-based unification of potato gene models with the UniTato
423 collaborative genome browser. *Front. Plant Sci.* **15**, (2024).
- 424 49. Bozan, I. *et al.* Pangenome analyses reveal impact of transposable elements and ploidy
425 on the evolution of potato species. *Proceedings of the National Academy of Sciences*
426 **120**, e2211117120 (2023).