



**HAL**  
open science

## **Assembly and annotation of the 'Golden Delicious' Doubled-Haploid GDDH18 apple genome**

Marine Salson, Christelle Lalanne, Maryline Cournol, Sylvain Hanteville, Damien Hinsinger, Patricia Faivre-Rampant, Aurélie Bérard, Arnaud Bellec, Stéphane Cauet, Nathalie Choisne, et al.

### ► **To cite this version:**

Marine Salson, Christelle Lalanne, Maryline Cournol, Sylvain Hanteville, Damien Hinsinger, et al.. Assembly and annotation of the 'Golden Delicious' Doubled-Haploid GDDH18 apple genome. G3, 2026, <10.1093/g3journal/jkag104/8662901>. <hal-05606946>

**HAL Id: hal-05606946**

**<https://hal.science/hal-05606946v1>**

Submitted on 29 Apr 2026

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License

1 **Title:** Assembly and annotation of the ‘Golden Delicious’ Doubled-Haploid GDDH18  
2 apple genome

3  
4 **Authors**

5 Marine Salson<sup>1,\*</sup>, Christelle Lalanne<sup>1</sup>, Maryline Cournol<sup>1</sup>, Sylvain Hanteville<sup>1</sup>, Damien  
6 Hinsinger<sup>2</sup>, Patricia Faivre-Rampant<sup>2</sup>, Aurélie Bérard<sup>2</sup>, Arnaud Bellec<sup>3</sup>, Stéphane  
7 Cauet<sup>3</sup>, Nathalie Choisne<sup>4</sup>, Sébastien Aubourg<sup>1</sup>, Anne-Laure Fanciullino<sup>1</sup>, Jean-Marc  
8 Celton<sup>1,\*†</sup>, Sandrine Balzergue<sup>1,\*†</sup>

9  
10 † These authors contributed equally.

11 \*Corresponding authors: [sandrine.balzergue@inrae.fr](mailto:sandrine.balzergue@inrae.fr),  
12 [jean-marc.celton@inrae.fr](mailto:jean-marc.celton@inrae.fr), [marine.salson@inrae.fr](mailto:marine.salson@inrae.fr)

13  
14 **Affiliation**

15 <sup>1</sup> Université Angers, Institut Agro, INRAE, IRHS, SFR QUASAV, F-49000 Angers,  
16 France

17  
18 <sup>2</sup> Université Paris-Saclay, Centre INRAE Île-de-France Versailles-Saclay, EPGV, Evry,  
19 91057, France

20  
21 <sup>3</sup> INRAE, CNRGV French Plant Genomic Resource Center, F-31320, Castanet Tolosan,  
22 France

23  
24 <sup>4</sup> IJPB, INRAE, Université Paris-Saclay, Versailles, F-78026, France

25  
26 **Keywords:** genome assembly, *Malus domestica*, *de novo* gene annotation,  
27 transposable elements annotation

## 28 **Abstract**

29 Apple (*Malus domestica* Borkh.) is an important fruit crop cultivated worldwide in  
30 temperate regions. Due to their frequent high content in repetitive DNA sequences  
31 and large size, plant genomes have long been considered challenging to assemble. A  
32 first reference genome of a 'Golden Delicious' doubled-haploid apple tree, GDDH13,  
33 was previously produced in 2017. The chromosomes of GDDH13 present however a  
34 high percentage of N stretches, around 12%, and above 8% of the assembly was not  
35 anchored on chromosomes. Here, we provide a chromosome-level assembly of  
36 another 'Golden Delicious' doubled-haploid apple tree, GDDH18. Interestingly, this  
37 tree presents phenotypic differences with GDDH13, notably associated with fruits  
38 size. In this assembly, 99.4% of the sequences were assembled into 17 contigs  
39 corresponding to the 17 expected chromosomes of apple, and all but one of these  
40 17 contigs present telomeric repeats at both their extremities. A high complete  
41 BUSCO (Benchmarking Universal Single-Copy Orthologs) score and a high LTR  
42 Assembly Index, respectively of 99.5% and above 22, attest to the high quality of  
43 the assembly and to the completeness of both the gene and the repetitive space. A  
44 *de novo* annotation of the assembly was produced, allowing to predict 51,892  
45 protein-coding genes, of which 93% are functionally annotated. A total of 6,239  
46 additional genes were functionally annotated in comparison with the GDDH13  
47 assembly. This high-quality assembly of GDDH18 along with a *de novo* annotation  
48 will both help to better understand the phenotypic differences between the trees  
49 GDDH13 and GDDH18, and will serve as an enhanced reference genome for  
50 advancing the study of apple genomics.

## 51 **Introduction**

52 Apple (*Malus domestica* Borkh.) belongs to the *Rosaceae* family and is an important  
53 fruit crop cultivated worldwide in temperate regions. It has a diploid genome of  
54 around 650 Mb with 17 chromosomes ranging from 31 Mb to 59 Mb ( $2n=2x=34$ ),  
55 and is composed of approximately 57% of transposable elements (TE) repeats  
56 (Daccord *et al.* 2017). A *Rosaceae* ancestor from which the apple originates notably  
57 underwent a whole genome duplication (WGD) around 27 Mya (Lallemand *et al.*  
58 2023). The domestication of apple likely began in the Tian Shan Mountains in  
59 Central Asia between 4,000 and 10,000 years ago, and the wild Central Asian  
60 species *Malus sieversii* has been shown to be the main contributor to the genome of  
61 the domesticated apple *Malus domestica* (Harris *et al.* 2002; Velasco *et al.* 2010).  
62 The spread of domesticated apples to Europe through the Silk Road was associated  
63 with additional contributions of other local wild species, including *Malus sylvestris*  
64 and *Malus orientalis* (Cornille *et al.* 2014; Cornille *et al.* 2019). Hybridization with  
65 wild relatives played thus an important role in the evolution of the cultivated apples  
66 (Cornille *et al.* 2012; Cornille *et al.* 2014; Duan *et al.* 2017; Sun *et al.* 2020). A  
67 doubled-haploid and highly homozygous 'Golden Delicious' line called Golden  
68 Delicious Doubled Haploid 13 (GDDH13) was developed thanks to breeding effort  
69 (Lespinasse *et al.* 1998; Daccord *et al.* 2017). This highly homozygous line allowed  
70 to facilitate the assembly of a reference genome (Daccord *et al.* 2017).

71

72 Plant genomes have long been considered challenging to assemble, notably because  
73 they are often characterized by high repetitiveness, large genome size, or high  
74 ploidy levels (Belser *et al.* 2018; Kong *et al.* 2023). A number of projects have  
75 however permitted the production of several high-quality plant genomes these last  
76 few years, notably thanks to the development and improvement of long reads  
77 sequencing technologies (Belser *et al.* 2018; Istace *et al.* 2021; Belser *et al.* 2021;

78 Aury *et al.* 2022; Salson *et al.* 2023). A chromosome-level assembly of GDDH13 was  
79 produced in 2017 (Daccord *et al.* 2017). However, this genome contains 12% of  
80 ambiguous bases (N), telomeric repeats are missing at the ends of most of the  
81 chromosomes, and 52.7 Mb of unanchored contigs were grouped onto an artificial  
82 pseudochromosome 0.

83

84 Here, we provide a high-quality assembly of the apple line Golden Delicious  
85 Doubled Haploid 18 (GDDH18). This tree originates from the same haploid seedling  
86 used to generate GDDH13 (Daccord *et al.* 2017), but differs by producing  
87 significantly smaller fruits and a higher number of seeds, which facilitates  
88 laboratory experimentation. This GDDH18 assembly showed improved continuity  
89 and better quality metrics in comparison with the previous GDDH13 assembly  
90 (Daccord *et al.* 2017). Additionally, we also provided a *de novo* annotation of this  
91 GDDH18 assembly. A high-quality assembly of GDDH18, combined with a *de novo*  
92 annotation, will not only facilitate the identification of genes or alleles associated  
93 with fruit development and clarify the phenotypic differences between GDDH13  
94 and GDDH18, but more importantly, it will provide an enhanced reference genome.  
95 This improved genomic resource represents a cornerstone for advancing the study  
96 and understanding of apple genomics, enabling more precise investigations and  
97 fostering future applications in fruit biology and breeding.

## 98 **Materials & Methods**

99

### 100 **Plant materials and sequencing**

101 Young leaves of a 'Golden Delicious' Doubled-Haploid tree GDDH18 (Daccord *et al.*  
102 2017) from UE Horti orchard  
103 (<https://doi.org/10.15454/1.5573931618268674E12>) were sampled. High DNA  
104 Molecular Weight was extracted with NucleoBond HMW DNA kit (Macherey-Nagel,  
105 Düren, Germany) according to the manufacturer. DNA quality was assessed with  
106 FemtoPulse (Agilent, CA, USA) after SRE-XL kit purification (Pacific Biosciences,  
107 Menlo Park, CA, USA) to eliminate fragments smaller than 30Kb and quantity was  
108 checked with a Qubit 4 fluorometer (Thermo Fisher Scientific, MA, USA).

109 A sequencing library was built with the ligation Oxford Nanopore sequencing kit  
110 V14 (SQK-LSK114) using 1000ng as input, with a ligation time of 30mn and using  
111 the LFB buffer. Sequencing was performed by using a R10.4.1 flowcell on a  
112 PromethION 24 instrument for 96 hours.

113 After the first 24 hours, the flowcell was washed with a nuclease-flush according to  
114 the manufacturer's instructions, then reloaded with a fresh library built as above.

115 Raw signal was live-base called with Dorado v7.8.3+f64462b6f and the SUP model  
116 (super accurate), as implemented in minKNOW v6.4.8. Taxonomic affiliation of the  
117 basecalled reads was assessed with centrifuge 1.0.3 (Kim *et al.* 2016) to detect  
118 contamination during the library and sequencing processes.

119

120

121

### 122 **Long reads assembly**

123 Long reads were filtered using Fitlong (v. 0.2.0, --length\_weight 10  
124 <https://github.com/rrwick/Fitlong>). Several read coverages were assessed for the

125 assembly (ranging from 30X to 60X, Supplementary Table 1). The two assemblers  
126 Flye (v. 2.9.2, Kolmogorov *et al.* 2019, and Hifiasm (v. 0.25.0, Cheng *et al.* 2021,  
127 2022, 2024) were tested. Telomeres were sought on the different assemblies using  
128 quarTeT (v. 1.2.5, Lin *et al.* 2023). We selected the assembly with the minimum  
129 number of contigs and for which telomeres were found at both extremities for a  
130 maximum number of chromosomes. We used both FCS-adaptor and FCS-GX  
131 (Astashyn *et al.* 2024; <https://github.com/ncbi/fcs>) to detect and remove adaptors  
132 and contaminations from foreign organisms in the selected assembly. We used  
133 publicly available paired-end short reads of the GDDH18 tree (run SRR5351715,  
134 BioProject PRJNA379390) for polishing of the assembly. The short reads were  
135 aligned to the assembly using bwa-mem2 (v.2.2.1, mem option, Li & Durbin 2010;  
136 Vasimuddin *et al.* 2019), and we only kept properly paired reads using samtools (v.  
137 1.9, -f 0x02 option, Danecek *et al.* 2021). Hapo-G was then used to polish the  
138 assembly using paired-end short reads (v. 2.31, default parameters, Aury & Istace  
139 2021).

140

### 141 **Quality assessment and comparison with previous *Malus* genomes**

142 The completeness of the gene space assembly was assessed using BUSCO (v. 6.0.0,  
143 Manni *et al.* 2021) and three databases: embryophyta\_odb12 (2,026 BUSCO genes),  
144 viridiplantae\_odb12 (822 BUSCO genes) and rosaceae\_odb12 (10,071 BUSCO  
145 genes). We computed the LTR Assembly Index (LAI) score (Ou *et al.* 2018, 2019)  
146 using LTR\_retriever (v 3.0.4) to assess the continuity of the assembly based on the  
147 repeat space following the protocol described here:  
148 [https://github.com/oushujun/LTR\\_retriever](https://github.com/oushujun/LTR_retriever). We used LTRharvest and  
149 LTR\_FINDER\_parallel to identify LTR. LTR\_retriever was then used to compute the  
150 LAI score. In order to assess the impact of the polishing step with the short reads,

151 we computed both the LAI and the BUSCO scores of the GDDH18 assembly before  
152 and after the polishing step.

153 For comparison purposes, the BUSCO and the LAI scores were also computed for  
154 the GDDH13 assembly (Daccord *et al.* 2017) and a recent high-quality ‘Golden  
155 Delicious’ diploid phased genome, GDT2T (Su *et al.* 2024). We used the same  
156 databases and the same procedures as previously described to compute the BUSCO  
157 and the LAI scores.

158 We used D-GENIES to compare the chromosomes of the GDDH18 assembly with the  
159 chromosomes of both GDDH13 and GDT2T (v. 1.4, hide noise option enabled and  
160 small matches filtered, Cabanettes & Klopp 2018).

161

### 162 **Telomeric and centromeric repeats identification**

163 We used quarTeT (v. 1.2.5, default parameters, Lin *et al.* 2023) to identify telomeric  
164 repeats on the chromosomes of the three ‘Golden Delicious’ assemblies GDDH18,  
165 GDDH13 and GDT2T.

166 We mapped the 9,716 bp HODOR (High cOpY goldEn deliciOus Repeat, accession  
167 KX869746, Daccord *et al.* 2017) sequence to the GDDH18 chromosomes using  
168 Minimap2 (Li 2018, -N 1000 -p 0.5 parameters), and computed the number of  
169 HODOR alignments found on non-overlapping genomic windows of 1 Mb for each  
170 chromosome. We only kept alignments longer than 3,000 bp. We used these  
171 alignments to estimate the putative centromere position on the GDDH18 assembly.  
172 We performed the same analysis on the GDDH13 genome for comparison.

173

### 174 **Mitochondrial and chloroplast sequences identification**

175 We used Oatk (Zhou *et al.* 2025 , v.1.0, default parameters) with the angiosperm  
176 database for assembling the mitochondrial and chloroplastic genomes using all the  
177 Nanopore long reads. We also mapped the mitochondrion (NCBI accession  
178 NC\_018554.1, Goremykin *et al.* 2012) and the chloroplast (NCBI accession

179 NC\_061549.1, Li *et al.* 2022) sequences to the GDDH18 assembly obtained with  
180 Hifiasm, using Minimap2 (v. 2.22, default parameters, Li 2018). This allowed to  
181 identify the contigs corresponding to mitochondrial and chloroplast sequences and  
182 removed them from the final assembly.

183

## 184 **Gene annotation**

185 Structural annotation of the *Malus domestica* GDDH18 genotype was performed  
186 using the EuGene pipeline (Sallet *et al.* 2019; Carrere & Gouzy 2023). This  
187 integrative gene prediction workflow combined multiple sources of evidence,  
188 including transcript alignments and protein homology searches. Transcriptome  
189 assemblies (Supplementary Table 2) were used to guide gene model prediction,  
190 while protein sequence homology was assessed against several reference  
191 databases, including Swiss-Prot, and the UniProt/TrEMBL plant subset. Detection  
192 and masking of transposable elements (TEs) were carried out using the Dfam  
193 consensus sequence database (<https://dfam.org>) and the TransposonPSI protein  
194 database (<http://transposonpsi.sourceforge.net>).

195 Functional annotation of predicted protein-coding genes was performed using a  
196 combination of eggNOG-mapper v2.1.12 (Cantalapiedra *et al.* 2021), BLAST+  
197 v2.15.0 (Altschul *et al.* 1990), InterProScan v5.64-96.0 (Jones *et al.* 2014), and  
198 KofamScan v1.3.0 (Aramaki *et al.* 2020). Outputs from these tools were integrated  
199 to produce final annotation files in GFF3 and TSV formats, prioritizing eggNOG  
200 annotations when overlaps occurred between sources. We applied the default  
201 parameter values for all analyses, unless otherwise specified.

202 We estimated the proportions of the genome in protein-coding genes using the  
203 outputs of this pipeline.

204

## 205 **Identification of syntenic blocks**

206 Identification of syntenic genes along the chromosomes of the GDDH18 assembly  
207 was performed using i-ADHoRe (Simillion *et al.* 2008, Proost *et al.* 2012) and an  
208 approach previously used for GDDH13 (Lallemand *et al.* 2023). For the  
209 identification of homologous genes, we first used blastP (e-value  $10^{-5}$ , max target  
210 seq 5, Camacho *et al.* 2009) and aligned all the proteins identified in GDDH18  
211 versus all. The resulting homology information was then used as input with i-  
212 ADHoRe (v.3.0, Proost *et al.* 2012) for the construction of syntenic blocks. We used  
213 the same parameters for i-ADHoRe (v.3.0, Proost *et al.* 2012) as previously used for  
214 GDDH13 (Lallemand *et al.* 2023): cluster type = colinear; tandem gap = 15; gap size  
215 = 30; cluster gap = 30; q value = 0.75; prob cutoff = 0.01, and anchor points = 5. The  
216 identification of homologous genes between the chromosomes allowed the  
217 reconstruction of syntenic fragments that were displayed using Circos (Krzywinski  
218 *et al.* 2009). To estimate the number of ohnologous genes we retained the pairs  
219 from the last WGD 27 Mya (Lallemand *et al.* 2023) belonging to the ohnologous  
220 chromosome pairs: 01-07, 01-15, 02-07, 02-15, 03-11, 04-06, 04-12, 05-10, 06-14,  
221 08-15, 09-17, 12-14 and 13-16.

222

## 223 **TE annotation**

224 The TEannot pipeline from the REPET package V3.0 (Quesneville *et al.* 2005, Flutre  
225 *et al.* 2011, Hoede *et al.* 2014,  
226 <https://urgi.versailles.inrae.fr/Tools/REPET/TEannot-tuto>) was used with default  
227 parameters to annotate TE copies in the GDDH18 genome. We used the consensus  
228 library of TE previously constructed for GDDH13 (Daccord *et al.* 2017,  
229 <https://urgi.versailles.inrae.fr/repetdb/report.do?id=49000001>), to identify  
230 repetitive elements in the GDDH18 assembly. A total of 2,456 TE consensus  
231 sequences were used, ranging from 358 to 17,023 bp, with a mean size of 3,150 bp.

232 We estimated the percentage of the different classes of TE copies in the GDDH18  
233 genome, using the GFF3 file generated with Teannot.

## 234 **Results & Discussion**

235

### 236 **Long reads sequencing**

237 We generated 4.4 M reads 'PASS' with a quality score Q larger than 10 (85.5% of the  
238 total sequencing data) with a N50 of 31.8 kb and a mean Qvalue of 20.48, totalling  
239 82.45 Gb of sequencing data. Of these reads, 91.64% were complete with Nanopore  
240 adapters kept at both ends, indicative both of the quality of the sequencing and of  
241 the library. A taxonomical analysis of contaminants was performed with  
242 centrifuge1.0.3, using a 10,000 reads subsample (with options -k 15 --min-hitlen 26  
243 --host-taxids 3750 --exclude-taxids 32644,28384,12908,77133). This analysis  
244 showed no significant contamination, with only 0.30% and 0.08% of bacterial and  
245 fungi reads, respectively (Supplementary Tables 3-5).

246

### 247 **Assembly of the long reads**

248 We tested long read assembly with both Hifiasm (Cheng *et al.* 2021, 2022, 2024)  
249 and Flye (Kolmogorov *et al.* 2019), and evaluated different depths ranging from 30X  
250 to 60X (Supplementary Table 1). A depth of 45X yielded the least fragmented  
251 assembly using Hifiasm, composed of a maximum number of contigs with telomeric  
252 repeats at both extremities (Supplementary Table 1). This assembly was obtained  
253 using 508,796 filtered reads, with a N50 and a mean length of 57 kb  
254 (Supplementary Table 6). A total of 26 contigs were obtained with a N50 of 36 Mb  
255 (Table 1). The N50 of the contigs of the GDDH18 assembly was 58 fold larger than  
256 GDDH13 (Supplementary Table 7). The total size of the assembly was 658 Mb  
257 (Table 1), corresponding to the expected size of the apple 'Golden Delicious' haploid  
258 genome (Daccord *et al.* 2017; Su *et al.* 2024).

259

### 260 **Identification of telomeric repeats and of the chromosomes**

261 Telomeric repeats were identified at both extremities of 16 contigs, likely  
262 corresponding to 16 chromosomes (Supplementary Figure 1, and Supplementary  
263 Table 8). In comparison, no chromosome of the GDDH13 assembly (Daccord *et al.*  
264 2017) showed telomeric repeats at both ends (Supplementary Table 7). Telomeric  
265 repeats were also found at one single extremity of two other contigs  
266 (Supplementary Fig. 4). One of these contig, ptg00006l, had a size of around 47 Mb,  
267 close to the expected size of chromosome 5 of GDDH13 (Supplementary Table 8).  
268 The other smallest contig was 498 kb in size (ptg000023l, Supplementary Table 8).  
269 While this telomeric region might represent one extremity of the chromosome 5,  
270 the contig mapped ambiguously to multiple positions of the GDDH13 and the  
271 GDT2T genomes. It was thus included in the Chr00 (see the corresponding section).  
272 Alignment of the 17 largest contigs of the GDDH18 assembly to the GDDH13  
273 chromosomes confirmed that these contigs corresponded to the 17 chromosomes  
274 of the apple genome (Figure 1). The remaining nine small contigs that were not  
275 included in the nuclear genome assembly represented 0.6% percent of the  
276 assembled sequences, highlighting the high contiguity and completeness of the  
277 contig sequences.

278

### 279 **Identification of mitochondrial and chloroplast sequences**

280 The nine contigs shorter than 1 Mb in size were assembled separately and were not  
281 integrated into the putative chromosomes by Hifiasm. Three contigs (ptg000024l,  
282 ptg000021c and ptg000022c) were mapped to the published mitochondria and  
283 chloroplast genomes of apple (NCBI accessions NC\_018554.1, Goremykin *et al.*  
284 2012, and NC\_061549.1 respectively, Li *et al.* 2022) with sequence identities  
285 between 75 and 100% (Supplementary Figure 2). The contig ptg000022c notably  
286 seemed to correspond to multiple adjacent chloroplast sequences, based on  
287 visualisation of the alignments (Supplementary Figure 2). These contigs that

288 aligned to the mitochondria and the chloroplast genomes were discarded from the  
289 final assembly. OatK (Zhou *et al.* 2025) allowed the assembly of a 162 kb sequence  
290 that aligned all along a previously published chloroplast apple sequence  
291 (NC\_061549.1 NCBI accession, Li *et al.* 2022) with high identity (Supplementary  
292 Figure 3, Supplementary Table 9 for annotation). However, no sequence was  
293 assembled with OatK for the mitochondrial genome. A previous study showed that  
294 the apple mitochondrial genome displays few similarity with other angiosperm  
295 mitochondrial genomes (Goremykin *et al.* 2012), which could make it difficult to  
296 reconstruct the sequences (Ni *et al.* 2025), notably when using databases.

297

### 298 **Construction of a Chr00 and of the final assembly**

299 We constructed a Chr00 (named MdG1800) totalling 2.2 Mb and comprising three  
300 contigs (ptg000013l, ptg000019l and ptg000023l of 587 kb, 782 kb and 498 kb,  
301 respectively), composed of repetitive sequences (identified by FCS-GX,  
302 Supplementary Table 10), and two additional contigs (ptg000025l and ptg000026l  
303 of 253 kb and 97 kb, respectively). The Chr00 gathers the sequences that have not  
304 been anchored with enough confidence to a chromosome. The Chr00 of the  
305 GDDH18 assembly was 24 fold smaller than the Chr00 of the GDDH13 assembly  
306 (Daccord *et al.* 2017), demonstrating a greater completeness of the assembly of the  
307 chromosomes (Supplementary Table 7).

308 Using FCS-GX, one contig, ptg000020c, was found to be similar to a *Cannabis sativa*  
309 sequence (Supplementary Table 10) and was thus removed from the assembly. The  
310 final assembly is composed of 17 chromosomes (MdG1801-MdG1817, Figure 1,  
311 Supplementary Table 10), including 16 chromosomes with telomeric repeats at  
312 both ends (Supplementary Table 8, Supplementary Figure 1), along with the Chr00,  
313 the mitochondrial and the chloroplast sequences. Three short adaptors sequences  
314 were found and removed from the final assembly (Supplementary Table 11).

315

## 316 **Quality metrics**

317 The complete BUSCO scores obtained for the whole genome assembly of GDDH18  
318 were high, between 97.4% and 99.8% depending on the database used (Table 1 and  
319 Supplementary Table 7), and were indicative of the high completeness of the gene  
320 assembly. These BUSCO scores were notably slightly higher than those found on  
321 GDDH13, between 96.0% and 98.1% using the same databases (Supplementary  
322 Table 7). Duplicated BUSCO score of the GDDH18 assembly was high 55.0%  
323 (*viridiplantae\_odb12* (n=822), Table 1). The ancestor of apple underwent a recent  
324 WGD around 27 Mya (Lallemand *et al.* 2023), likely explaining this high fraction of  
325 duplicated BUSCOs in the assembly.

326 The LAI score imputed for the assembly was 22.16, higher than the score of the  
327 GDDH13 assembly of 20.16 (Supplementary Table 7). An LAI score above 20 is  
328 characteristic of high-quality (“gold standard”) assemblies (Ou *et al.* 2018).  
329 Accordingly, a high LAI reflects a good continuity of repetitive genomic regions.

330 The short-reads polishing step did not enhance the BUSCO scores, and resulted in a  
331 slight drop in the LAI score (Supplementary Table 12). This could be explained by  
332 the low mean depth of the short reads available for the GDDH18 tree, around 17X  
333 after mapping and filtering of the properly paired reads. Therefore, we propose a  
334 final assembly that does not include results from the polishing step. It has recently  
335 been suggested that assemblies using Oxford Nanopore R10.4 and the latest Kit V14  
336 technology may not require a short-reads polishing step (Belinchon-Moreno *et al.*  
337 2025; Sereika *et al.* 2022).

338

## 339 **Gene annotation**

340 *De novo* annotation with Eugene allowed the identification of 65,170 predicted  
341 genes, of which 51,892 corresponded to protein-coding genes and 13,278 to non-

342 coding RNA genes (Supplementary Table 13). Protein-coding genes had a mean  
343 gene length of 3,130 bp, with the shortest gene spanning 150 bp and the longest  
344 79,620 bp (Supplementary Table 13). About 80% of protein-coding genes  
345 contained introns. The protein-coding genes comprised on average 4.79 exons  
346 (mean exon length of 315 bp) and 3.79 introns (mean intron length of 427 bp). The  
347 average coding sequence (CDS) length was 1,075 bp, while mean 5' and 3' UTR  
348 lengths were 255 bp and 421 bp, respectively. Non-coding RNA genes including  
349 snoRNA, rRNA, tRNA, miRNA, snRNA and other undefined ncRNA, displayed an  
350 average length of 1,074 bp.

351 Completeness assessment of the predicted proteome was evaluated using BUSCO  
352 v6.0.0 (viridiplantae\_odb12 dataset, 822 genes), yielding completeness scores of  
353 96.4%, with less than 0.4% missing BUSCOs. Among the 51,892 predicted proteins,  
354 48,379 (93%) were functionally annotated by at least one tool (Supplementary  
355 Table 14). Specifically, 43,063 proteins were annotated by eggNOG, 46,060 by  
356 BLAST, 38,520 by InterProScan, and 40,641 by KofamScan (Supplementary Table  
357 14). A total of 41,686 proteins were associated with at least one Gene Ontology  
358 (GO) term, and 42,625 with a KEGG Orthology (KO) term (Supplementary Table  
359 14). To assess the reliability of predicted proteins, the PSAURON machine-learning  
360 model (Sommer *et al.* 2024) was applied, assigning a probability score for each  
361 protein based on sequence features, conservation patterns, and length  
362 distributions. Overall, 89.3% of predicted proteins exhibited a high PSAURON  
363 confidence score (higher than 0.8). Among the 4,316 unannotated proteins, 43.8%  
364 received a significant PSAURON score indicative of genuine protein-coding  
365 potential. The majority of unannotated proteins (3,019 i.e. 70.7%) were short, with  
366 lengths below 100 amino acids.

367 The GDDH18 assembly presented 6,239, 4,004 and 4,075 more functionally  
368 annotated protein-coding genes in comparison with GDDH13 (Daccord *et al.* 2017),

369 and with the two haplotypes of the GDT2T genome (Su *et al.* 2024), respectively  
370 (Supplementary Table 7). The protein-coding and annotated genes corresponded to  
371 23.4% of the genome (Figure 2), in accordance with the percentage found in the  
372 GDDH13 genome (Daccord *et al.* 2017). A lower density of genes was observed  
373 around the putative centromeric regions of GDDH18 (Supplementary Figure 4).

374

### 375 **Identification of syntenic blocks**

376 A total of 700 syntenic blocks were identified between the chromosomes of the  
377 GDDH18 genome with i-ADHoRe (Supplementary Tables 15-17, Supplementary  
378 Figure 5). These syntenic blocks were identified between chromosomes similar to  
379 GDDH13 (Lallemand *et al.* 2023, Supplementary Figure 5). These syntenic blocks  
380 had a mean size of 1.5 Mb and encompassed 106 genes on average, with a median  
381 number of 31 genes. A smaller number of syntenic blocks were identified in the  
382 GDDH18 assembly in comparison with the GDD13 assembly (865 syntenic blocks in  
383 GDDH13, Lallemand *et al.* 2023), and they consisted in a larger median number of  
384 genes, 31 in GDDH18 and nine in GDDH13 (Lallemand *et al.* 2023), which could be  
385 further indicative that the GDDH18 assembly was less fragmented. We estimated a  
386 total of 12,006 pairs of ohnologous genes likely originating from the last WGD 27  
387 Mya (Supplementary Table 18), in accordance with a previous study in GDDH13  
388 (Lallemand *et al.* 2023).

389

### 390 **TE annotation**

391 All of the 2,456 consensus had at least one copy on the GDDH18 genome. We found  
392 that approximately 63.2% of the GDDH18 assembly was composed of repeats, in  
393 accordance with previous results on 'Golden Delicious' genomes (Daccord *et al.*  
394 2017; Su *et al.* 2024). The class I elements were the most prevalent, representing  
395 46.6% of the genome, including 41.9% LTR, 4.1% LINE and 0.4% SINE (Figure 2).

396 DNA transposon or class II elements represented 13.2% of the genome, including  
397 10.2% TIR, 1.0% MITE, 0.3% HELITRON, and 1.8% class II elements no further  
398 classified (Figure 2). A total of 4.1% of copies represented unclassified TEs (Figure  
399 2). The proportions of class I and class II elements were similar as the GDDH13  
400 genome (Daccord *et al.* 2017). *Gypsy* and *Copia* LTR elements were present in  
401 proportions close to those reported for the GDT2T genome (Su *et al.* 2024),  
402 representing approximately 24.8% and 12.6%, respectively.

403

#### 404 **HODOR repeats**

405 We found HODOR (High cOpy goldEN deliciOus Repeat, accession KX869746,  
406 Daccord *et al.* 2017) repeats gathered on specific genomic regions of each  
407 chromosome (Supplementary Figure 6). These genomic regions enriched in HODOR  
408 repeats correspond to the putative centromeres of the apple chromosomes (Su *et*  
409 *al.* 2024) and were found in similar regions in GDDH18 and GDDH13  
410 (Supplementary Figure 6). The high repetitiveness of the HODOR sequence might  
411 explain partially the fragmentation of the GDDH13 assembly obtained with shorter  
412 reads (Daccord *et al.* 2017).

413

#### 414 **Inversions found between the GDDH13 and GDDH18 assemblies**

415 Alignments of the GDDH18 and the GDDH13 chromosomes showed an overall  
416 continuity between homologous chromosomes (Figure 1).

417 Some DNA sequence inversions were identified between GDDH13 and GDDH18 on  
418 chromosomes 1, 8 and 16 (Supplementary Figure 7). These putative inversions  
419 were also found when comparing GDDH13 and GDT2T (Su *et al.* 2024), but they  
420 were not found when comparing GDDH18 to the two haplotypes of GDT2T  
421 (Supplementary Figure 7), suggesting that sequences at these particular loci might  
422 have been misassembled in GDDH13.

423 Other putative inversions were found when comparing the chromosomes 2, 3, 6  
424 and 12 of GDDH18 to one of the haplotypes of GDT2T (Supplementary Figure 8).  
425 These putative inversions were also previously observed when comparing GDDH13  
426 to GDT2T (Su *et al.* 2024). These putative inversions are close to the centromeres  
427 (Su *et al.* 2024). While these inverted sequences may originate from a mis-  
428 assembly, it has already been proposed and observed that centromeric regions  
429 might be naturally prone to inversions (Harringmeyer & Hoesktra 2022, Salson *et*  
430 *al.* 2025). These putative inversions have to be validated with independent method  
431 such as Hi-C data.

## 432 **Data and code availability**

433 The GDDH18 ONT long reads, the genome assembly and the annotation of the  
434 genes are available in the European Nucleotide Archive (ENA), under the project  
435 PRJEB97937: <https://www.ebi.ac.uk/ena/browser/view/PRJEB97937>. The raw  
436 long reads can be found here:  
437 <https://www.ebi.ac.uk/ena/browser/view/ERR15761846>, the assembly here:  
438 [https://www.ebi.ac.uk/ena/browser/view/GCA\\_977014495](https://www.ebi.ac.uk/ena/browser/view/GCA_977014495), and the annotation  
439 file with the functionally annotated genes here:  
440 <https://www.ebi.ac.uk/ena/browser/view/ERZ28547962>. The plastid sequence,  
441 the GFF3 of all the genes annotation file and of the transposable elements  
442 annotation file have been deposited via the GSA figshare portal  
443 (<https://doi.org/10.25387/g3.32012313>), and can be shared under request . The  
444 scripts used for the GDDH18 assembly and the analysis of the genome can be found  
445 on the GitLab repository: [https://forge.inrae.fr/OPTIMAE/gddh18\\_assembly](https://forge.inrae.fr/OPTIMAE/gddh18_assembly).

446

## 447 **Author contributions**

448 S.H and M.C. collected the leaves of GDDH18 tree. M.C. performed the DNA  
449 extraction and first QC. D.H, P.F-V and A.B. performed the sequencing. M.S.  
450 performed the assembly, and the quality assessments of the GDDH18 genome. A.B.  
451 and S.C performed the *de novo* genes annotation. N.C. performed the transposable  
452 elements annotation. M.S., C.L., D.H., S.A., A.B., S.C., N.C., S.A., A-L.F. J-M.C. and S.B.  
453 contributed to the interpretations of the results. M.S., J-M.C and S.B. took the lead of  
454 writing of the manuscript. All authors read, contributed to the writing, provided  
455 feedback and approved the final manuscript.

456

## 457 **Funding**

458 This work has been funded by PlantAlliance consortium, INRAE, 75338, Paris,  
459 France.

460

### 461 **Acknowledgements**

462 We would like to thank the Experimental Unit Horti at Beaucozé, France (EU Horti,  
463 DOI: <https://doi.org/10.15454/1.5573931618268674E12>). The authors  
464 acknowledge PTM ANAN (Balzergue *et al.* 2024) for DNA quality check notably. We  
465 thank Claudine Landès and Andréa Bouanich for their help in the syntenic blocks  
466 analysis.

467

468 **Literature cited**

469

470 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment  
471 search tool. *J Mol Biol.* ;215(3):403–410. [https://doi.org/10.1016/S0022-  
472 2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2)

473

474 Aramaki, Blanc-Mathieu, Endo, Ohkubo, Kanehisa et al. 2020. KofamKOALA: KEGG  
475 Ortholog assignment based on profile HMM and adaptive score threshold,  
476 *Bioinformatics*, Volume 36, Issue 7, Pages 2251–2252,  
477 <https://doi.org/10.1093/bioinformatics/btz859>

478

479 Astashyn, A., Tvedte, E.S., Sweeney, D., Sapojnikov V., Bouk N. et al. 2024. Rapid and  
480 sensitive detection of genome contamination at scale with FCS-GX. *Genome Biol* 25,  
481 60. <https://doi.org/10.1186/s13059-024-03198-7>

482

483 Aury, JM & Istace, B. 2021. Hapo-G, haplotype-aware polishing of genome  
484 assemblies with accurate reads, *NAR Genomics and Bioinformatics*, Volume 3, Issue  
485 2, lqab034, <https://doi.org/10.1093/nargab/lqab034>

486

487 Aury, JM, Engelen, S, Istace, B, Monat, C, Lasserre-Zuber et al. 2022. Long-read and  
488 chromosome-scale assembly of the hexaploid wheat genome achieves high  
489 resolution for research and breeding. *Gigascience*, 11.  
490 <https://doi.org/10.1093/gigascience/giac034>

491

492 Belinchon-Moreno, Berard, Canaguier, Le-Clainche, Rittener-Ruff et al. 2025.  
493 Nuclear and organelle genome assemblies of 5 *Cucumis melo* L. accessions, Ananas,  
494 Canton, PI 414723, Vedrantaïs, and Zhimali, belonging to diverse botanical groups,  
495 *G3 Genes/Genomes/Genetics*, Volume 15, Issue 7, jkaf098,  
496 <https://doi.org/10.1093/g3journal/jkaf098>

497

498 Belser, C, Istace, B, Denis, E, Dubarry, M, Baurens, FC et al. 2018. Chromosome-scale  
499 assemblies of plant genomes using nanopore long reads and optical maps. *Nat*  
500 *Plants*, 4, 11:879-887. <https://doi.org/10.1038/s41477-018-0289-4>

501 Belser, C, Baurens, FC, Noel, B, Martin, G, Cruaud, C et al. 2021. Telomere-to-  
502 telomere gapless chromosomes of banana using nanopore sequencing. *Commun*  
503 *Biol*, 4, 1:1047. <https://doi.org/10.1038/s42003-021-02559-3>

504 Cabanettes F & Klopp C. 2018. D-GENIES: dot plot large genomes in an interactive,  
505 efficient and simple way. *PeerJ* 6:e4958. <https://doi.org/10.7717/peerj.4958>

- 506 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J et al. 2009. BLAST+:  
507 architecture and applications. BMC Bioinformatics. 15;10:421.  
508 <https://doi.org/10.1186/1471-2105-10-421>
- 509 Cantalapiedra, Hernandez-Plaza, Letunic, Bork & Huerta-Cepas. 2021. eggNOG-  
510 mapper v2: functional annotation, orthology assignments, and domain prediction  
511 at the metagenomic scale. Molecular Biology and Evolution, msab293,  
512 <https://doi.org/10.1093/molbev/msab293>
- 513 Carrere, Sébastien, & Gouzy, Jérôme. 2023. Eukaryote EuGene pipeline Version 2  
514 (2.0.0). Zenodo <https://doi.org/10.5281/zenodo.7648710>
- 515 Cheng, H., Concepcion, G.T., Feng, X., Zhang, H. & Li H. 2021. Haplotype-resolved de  
516 novo assembly using phased assembly graphs with hifiasm. *Nat Methods*, 18:170-  
517 175. <https://doi.org/10.1038/s41592-020-01056-5>  
518
- 519 Cheng, H., Jarvis, E.D., Fedrigo, O., Koepfli, K.P., Urban, L. et al. 2022. Haplotype-  
520 resolved assembly of diploid genomes without parental data. *Nature Biotechnology*,  
521 40:1332–1335. <https://doi.org/10.1038/s41587-022-01261-x>  
522
- 523 Cheng, H., Asri, M., Lucas, J., Koren, S. & Li, H. 2024. Scalable telomere-to-telomere  
524 assembly for diploid and polyploid genomes with double graph. *Nat Methods*,  
525 21:967-970. <https://doi.org/10.1038/s41592-024-02269-8>
- 526 Cornille A, Gladieux P, Smulders MJM, Roldán-Ruiz I, Laurens F et al. 2012. New  
527 Insight into the History of Domesticated Apple: Secondary Contribution of the  
528 European Wild Apple to the Genome of Cultivated Varieties. *PLoS Genet* 8(5):  
529 e1002703. <https://doi.org/10.1371/journal.pgen.1002703>
- 530 Cornille A, Giraud T, Smulders MJ, Roldán-Ruiz I & Gladieux P. 2014. The  
531 domestication and evolutionary ecology of apples. *Trends Genet.* ;30(2):57-65.  
532 <https://doi.org/10.1016/j.tig.2013.10.002>
- 533 Cornille A, Antolín F, Garcia E, Vernesi C, Fietta A et al. 2019. A Multifaceted  
534 Overview of Apple Tree Domestication. *Trends Plant Sci.* Aug;24(8):770-782.  
535 <https://doi.org/10.1016/j.tplants.2019.05.007>
- 536 Daccord, N., Celton, JM., Linsmith, G., Becker C., Choisne N. et al. 2017. High-quality  
537 *de novo* assembly of the apple genome and methylome dynamics of early fruit  
538 development. *Nat Genet* 49, 1099–1106. <https://doi.org/10.1038/ng.3886>  
539
- 540 Danecek, Bonfield, Liddle, Marshall, Ohan et al. 2021. Twelve years of SAMtools and  
541 BCFtools, *GigaScience*, Volume 10, Issue 2, February giab008,  
542 <https://doi.org/10.1093/gigascience/giab008>

543

544 Duan, N., Bai, Y., Sun, H., Wang N., Ma Y. et al. 2017. Genome re-sequencing reveals  
545 the history of apple and supports a two-stage model for fruit enlargement. *Nat*  
546 *Commun* 8, 249. <https://doi.org/10.1038/s41467-017-00336-7>

547

548 Flutre T, Duprat E, Feuillet C, Quesneville H. 2011.  
549 Considering transposable element diversification in de novo annotation  
550 approaches. *PLoS ONE* 6(1): e16526.  
551 <https://doi.org/10.1371/journal.pone.0016526>

552

553 Goremykin VV, Lockhart PJ, Viola R, Velasco R. 2012. The mitochondrial genome of  
554 *Malus domestica* and the import-driven hypothesis of mitochondrial genome  
555 expansion in seed plants. *Plant J.* Aug;71(4):615-26.  
556 <https://doi.org/10.1111/j.1365-313X.2012.05014.x>

557

558 Grabherr, M., Haas, B., Yassour, M. Levin J., Thompson D. et al. 2011. Full-length  
559 transcriptome assembly from RNA-Seq data without a reference genome. *Nat*  
560 *Biotechnol* 29, 644–652. <https://doi.org/10.1038/nbt.1883>

561

562 Harringmeyer, O.S., Hoekstra, H.E. 2022. Chromosomal inversion polymorphisms  
563 shape the genomic landscape of deer mice. *Nat Ecol Evol* 6, 1965–1979.  
564 <https://doi.org/10.1038/s41559-022-01890-0>

565

566 Harris SA, Robinson JP, Juniper BE. 2002. Genetic clues to the origin of the apple.  
567 *Trends Genet.*;18(8):426-30. [https://doi.org/10.1016/s0168-9525\(02\)02689-6](https://doi.org/10.1016/s0168-9525(02)02689-6)

568

569 Hoede C, , Arnoux S, Moissette M, Chaumier T, Inizan O et al. 2014. PASTEC: An  
570 Automatic Transposable Element Classification Tool. *PLoS One.* 2014 May  
571 2;9(5):e91929. <https://doi.org/10.1371/journal.pone.0091929>

572

573 Istace, B, Belser, C, Falentin, C, Labadie, K, Boideau, F et al. 2021. Sequencing and  
574 Chromosome-Scale Assembly of Plant Genomes, *Brassica rapa* as a Use Case.  
575 *Biology (Basel)*, 10,732. <https://doi.org/10.3390/biology10080732>

576

577 Jones, Binns, Chang, Fraser, Li et al. 2014. InterProScan 5: genome-scale protein  
578 function classification, *Bioinformatics*, Volume 30, Issue 9, Pages 1236–1240,  
579 <https://doi.org/10.1093/bioinformatics/btu031>

580

581 Kim, D., Song, L., Breitwieser, F. P., & Salzberg, S. L. 2016. Centrifuge: rapid and  
582 sensitive classification of metagenomic sequences. *Genome research*, 26(12), 1721-  
583 1729, <https://doi.org/10.1101/gr.210641.116>

584

585 Kolmogorov, M, Yuan, J, Lin, Y, Pevzner, PA. 2019. Assembly of long, error-prone  
586 reads using repeat graphs. *Nat Biotechnol*, 37, 5:540-546,  
587 <https://doi.org/10.1038/s41587-019-0072-8>

588 Kong W, Wang Y, Zhang S, Yu J & Zhang X. 2023. Recent Advances in Assembly of  
589 Complex Plant Genomes, *Genomics, Proteomics & Bioinformatics*, Volume 21, Issue  
590 3, Pages 427–439, <https://doi.org/10.1016/j.gpb.2023.04.004>

591 Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R et al. 2009. Circos: an  
592 information aesthetic for comparative genomics. *Genome Res*. Sep;19(9):1639-45.  
593 <https://doi.org/10.1101/gr.092759.109>

594

595 Lallemand, Leduc, Desmazières, Aubourg, Rizzon et al. 2023. Insights into the  
596 Evolution of Ohnologous Sequences and Their Epigenetic Marks Post-WGD in *Malus*  
597 *Domestica*, *Genome Biology and Evolution*, Volume 15, Issue 10, evad178,  
598 <https://doi.org/10.1093/gbe/evad178>

599

600 Lespinasse, Y., Bouvier, L., Djulbic, M. & Chevreau, E. 1998. Haploidy in apple and  
601 pear. *Acta Hort.* 484, 49-54 <https://doi.org/10.17660/ActaHortic.1998.484.4>

602 Li H & Durbin R. 2010. Fast and accurate long-read alignment with Burrows-  
603 Wheeler transform. *Bioinformatics*, 26, 589-595,  
604 <https://doi.org/10.1093/bioinformatics/btp698>

605 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences, *Bioinformatics*,  
606 Volume 34, Issue 18, Pages 3094–3100,  
607 <https://doi.org/10.1093/bioinformatics/bty191>

608 Li X, Ding Z, Miao H, Bao J & Tian X. 2022. Complete chloroplast genome studies of  
609 different apple varieties indicated the origin of modern cultivated apples from  
610 *Malus sieversii* and *Malus sylvestris*. *PeerJ*. 18;10:e13107.  
611 <https://doi.org/10.7717/peerj.13107>

612 Lin, Ye, Li, Chen, Wu et al. quarTeT: a telomere-to-telomere toolkit for gap-free  
613 genome assembly and centromeric repeat identification, *Horticulture Research*,  
614 Volume 10, Issue 8, uhad127, <https://doi.org/10.1093/hr/uhad127>

615 Manni M, Berkeley MR, Seppey M, Zdobnov EM. 2023. BUSCO: assessing genomic  
616 data quality and beyond. *Curr Protoc*. 2021;1(12):e323.  
617 <https://doi.org/10.1002/cpz1.323>

618 Ni, Y., Li, J., Tan, Y., Shen, G. & Liu, C. 2025. Advance in the assembly of the plant  
619 mitochondrial genomes using high-throughput DNA sequencing data of total

- 620 cellular DNAs. *Plant Biotechnol. J.* 23: 4944-4965.  
621 <https://doi.org/10.1111/pbi.70249>
- 622 Ou S & Jiang N. 2018. LTR\_retriever: A Highly Accurate and Sensitive Program for  
623 Identification of Long Terminal Repeat Retrotransposons , *Plant Physiology*,  
624 Volume 176, Issue 2, Pages 1410–1422, <https://doi.org/10.1104/pp.17.01310>
- 625 Ou S, Chen J & Jiang N. 2018. Assessing genome assembly quality using the LTR  
626 Assembly Index (LAI), *Nucleic Acids Research*, Volume 46, Issue 21, 30, Page e126,  
627 <https://doi.org/10.1093/nar/gky730>
- 628 Proost S, Fostier J, De Witte D, Dhoedt B, Demeester P et al. 2012. i-ADHoRe 3.0--  
629 fast and sensitive detection of genomic homology in extremely large data sets.  
630 *Nucleic Acids Res.* ;40(2):e11. <https://doi.org/10.1093/nar/gkr955>
- 631 Quesneville H, Bergman C, Andrieu O, Autard D, Nouaud D et al. 2005. Combined  
632 evidence annotation of transposable elements in genome sequences. *PLoS Comput*  
633 *Biol* 1(2): e22. <https://doi.org/10.1371/journal.pcbi.001>
- 634 Sallet, E., Gouzy, J. & Schiex, T. 2019. EuGene: An Automated Integrative Gene  
635 Finder for Eukaryotes and Prokaryotes. In: Kollmar, M. (eds) *Gene Prediction.*  
636 *Methods in Molecular Biology*, vol 1962. Humana, New York, NY.  
637 [https://doi.org/10.1007/978-1-4939-9173-0\\_6](https://doi.org/10.1007/978-1-4939-9173-0_6)
- 638 Salson, Orjuela, Mariac, Zekraoui, Couderc et al. 2023. An improved assembly of the  
639 pearl millet reference genome using Oxford Nanopore long reads and optical  
640 mapping, *G3 Genes/Genomes/Genetics*, Volume 13, Issue 5, jkad051,  
641 <https://doi.org/10.1093/g3journal/jkad051>
- 642 Salson, M., Duranton, M., Huynh, S., Mariac C., Tranchant-Dubreuil C. et al. 2025.  
643 Interplay between large low-recombining regions and pseudo-overdominance in a  
644 plant genome. *Nat Commun* 16, 6458. <https://doi.org/10.1038/s41467-025-61529-z>
- 646 Sereika, M., Kirkegaard, R.H., Karst, S.M., Michaelsen, T.Y., Sorensen, E.A. et al. 2022.  
647 Oxford Nanopore R10.4 long-read sequencing enables the generation of near-  
648 finished bacterial genomes from pure cultures and metagenomes without short-  
649 read or reference polishing. *Nat Methods* 19, 823–826.  
650 <https://doi.org/10.1038/s41592-022-01539-7>
- 651 Simillion, Janssens, Sterck & Van de Peer. 2008. i-ADHoRe 2.0: an improved tool to  
652 detect degenerated genomic homology using genomic profiles, *Bioinformatics*,  
653 Volume 24, Issue 1, Pages 127–128,  
654 <https://doi.org/10.1093/bioinformatics/btm449>

- 655 Sommer, Zimin & Salzberg. 2025. PSAURON: a tool for assessing protein annotation  
656 across a broad range of species, *NAR Genomics and Bioinformatics*, Volume 7, Issue  
657 1, lqae189, <https://doi.org/10.1093/nargab/lqae189>
- 658 Su, Yang, Wang, Li, Long et al. 2024. Phased telomere-to-telomere reference  
659 genome and pangenome reveal an expansion of resistance genes during apple  
660 domestication, *Plant Physiology*, Volume 195, Issue 4, Pages 2799–2814,  
661 <https://doi.org/10.1093/plphys/kiac258>
- 662 Sun, X., Jiao, C., Schwaninger, H., Chao, C.T, Ma, Y. et al. 2020. Phased diploid genome  
663 assemblies and pan-genomes provide insights into the genetic history of apple  
664 domestication. *Nat Genet* 52, 1423–1432. [https://doi.org/10.1038/s41588-020-](https://doi.org/10.1038/s41588-020-00723-9)  
665 00723-9
- 666 M. Vasimuddin, S. Misra, H. Li & S. Aluru. 2019. Efficient Architecture-Aware  
667 Acceleration of BWA-MEM for Multicore Systems. *2019 IEEE International Parallel  
668 and Distributed Processing Symposium (IPDPS)*, Rio de Janeiro, Braz, pp. 314-324,  
669 <https://doi.org/10.1109/IPDPS.2019.00041>
- 670 Velasco, R., Zharkikh, A., Affourtit, J., Dhingra, A., Cestaro, A. et al. 2010. The genome  
671 of the domesticated apple (*Malus × domestica* Borkh.). *Nat Genet* 42, 833–839.  
672 <https://doi.org/10.1038/ng.654>
- 673 Zhou, C., Brown, M., Blaxter, M. Darwin Tree of Life Project Consortium, McCarthy,  
674 S.A. et al. 2025. Oatk: a de novo assembly tool for complex plant organelle genomes.  
675 *Genome Biol* 26, 235. <https://doi.org/10.1186/s13059-025-03676-6>

676

## 677 **Figures legends**

678 **Figure 1.** Chromosome-to-chromosome alignments between GDDH13 and  
679 GDDH18.

680 The 17 longest contigs of the GDDH18 assembly identified as the 17 putative  
681 chromosomes of the apple genome are aligned to the 17 chromosomes of GDDH13  
682 (Daccord *et al.* 2017). On the x-axis are the GDDH18 chromosomes and on the y-  
683 axis are the GDDH13 chromosomes. Yellow, orange, light green and dark green  
684 colors correspond to alignments with identities between 0% and 25%, 25% and  
685 50%, between 50% and 75%, and above 75%, respectively.

686 **Figure 2.** Percentage of the GDDH18 genome represented by protein-coding genes  
687 and TEs regions.

688 Percentage of the GDDH18 genome that represents the protein-coding genes and  
689 the TEs regions. Light green represents the functionally annotated protein-coding  
690 genes and dark green the unannotated protein-coding genes. The class I TEs are in  
691 three different shades of blue and include in the order LTR, LINE and SINE  
692 elements. The class II TEs are in orange, brown, coral and red, and include TIR,  
693 unclassified class II, MITE, and HELITRON elements respectively. Unclassified TEs  
694 are in grey.

695

<b>Total assembly</b>	
Total length	658,096,990 bp
GC content	38.9%
Complete BUSCOs viridiplantae_odb12 (n=822)	99.5%
Complete and duplicated BUSCOs	55.0%
Fragmented BUSCOs	0.1%
Missing BUSCOs	0.4%
<b>Chromosomes</b>	
Number of chromosomes	17
Total length of chromosomes	655 Mb
Percentage of Ns	No gap
<b>Contigs</b>	
Number of contigs	28
Longest contig	59 Mb
N50 (contigs)	36 Mb
<b>Gene annotation</b>	
Total number of genes	65,170
Number of protein-coding genes	51,182
Number of protein-coding and functionally annotated genes	48,379

Figure 1

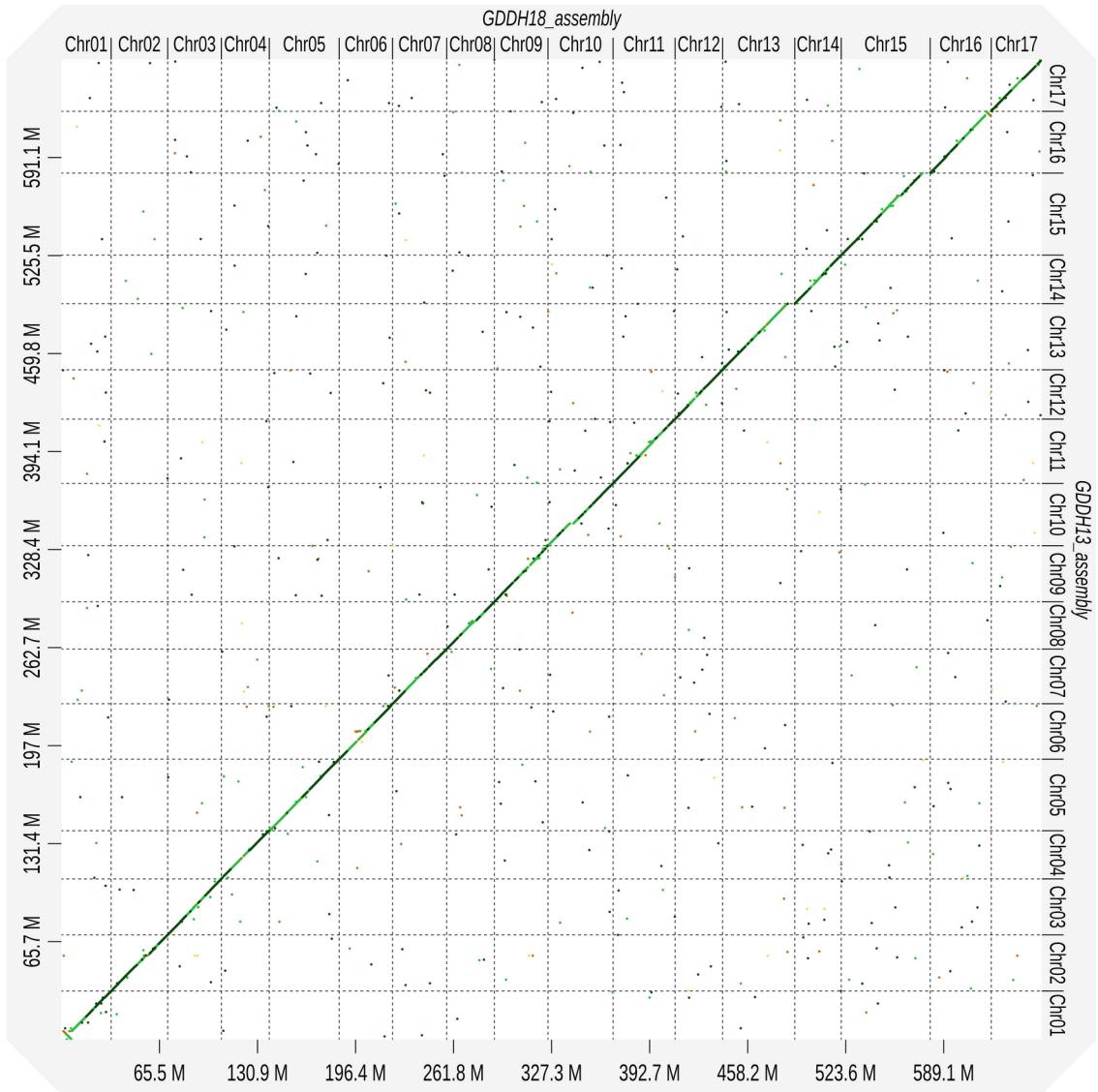


Figure 2

