1 **The complete reference genome for grapevine (*Vitis vinifera* L.)**

2 **genetics and breeding**

3

4 Xiaoya Shi[1, 2, #,] Shuo Cao[2, 3, #], Xu Wang[2, 6, #], Siyang Huang[2, 4], Yue Wang[2, 5],

5 Zhongjie Liu[2], Wenwen Liu[2], Xiangpeng Leng[1], Yanling Peng[2], Nan Wang[2], Yiwen

6 Wang[2], Zhiyao Ma[2], Xiaodong Xu[2], Fan Zhang[2], Hui Xue[2], Haixia Zhong[7], Yi Wang[8],

7 Kekun Zhang[9], Amandine Velt[10], Komlan Avia[10], Daniela Holtgräwe[11], Jérôme

8 Grimplet[12], José Tomás Matus[13], Doreen Ware[14,15]; Xinyu Wu[7], Haibo Wang[16],

9 Chonghuai Liu[17], Yuling Fang[9], Camille Rustenholz[10, *], Zongming Cheng[18, *], Hua

10 Xiao[2, 7, *], Yongfeng Zhou[2, 19, *]

11

12

13

14

15

16

17

18

19

20

21

22

28    1 College of Horticulture, Qingdao Agricultural University, 266109, Qingdao, China

29    2 State Key Laboratory of Tropical Crop Breeding, Shenzhen Branch, Guangdong Laboratory of

30    Lingnan Modern Agriculture, Key Laboratory of Synthetic Biology, Ministry of Agriculture and

31    Rural Affairs, Agricultural Genomics Institute at Shenzhen, Chinese Academy of Agricultural

32    Sciences, Shenzhen, China

33    3 Key Laboratory of Horticultural Plant Biology Ministry of Education, Huazhong Agricultural

34    University, Wuhan, People's Republic of China

35    4 Guangxi Key Lab for Sugarcane Biology, Guangxi University, Nanning, Guangxi, 530005

36    China

37    5 State Key Laboratory of Silkworm Genome Biology, Southwest University, Chongqing, China

38    6 School of Agriculture and Food Science, University College Dublin, Belfield, Dublin 4, Ireland

39    7 Institute of Horticulture Crops, Xinjiang Academy of Agricultural Sciences, Urumqi, China

40    8 Beijing Key Laboratory of Grape Science and Enology, Institute of Botany, Chinese Academy

41    of Sciences, Xiangshan, Beijing 100093, China

42    9 College of Enology, Northwest A&F University, Yangling 712100, China

43    10 SVQV, INRAE - University of Strasbourg, 68000 Colmar, France

44    11 Genetics and Genomics of Plants, CeBiTec & Faculty of Biology, Bielefeld University, 33615

45    Bielefeld, Germany

46    12 Unidad de Hortofruticultura, Centro de Investigación y Tecnología Agroalimentaria de Aragón

47    (CITA), 50059 Zaragoza, Spain

48    13 Institute for Integrative Systems Biology (I2SysBio), Universitat de València-CSIC, Paterna,

49    46908, Valencia, Spain

50    14 Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 11724, USA; ware@cshl.edu (D.W.)

51    15 USDA ARS NEA Robert W. Holley Center for Agriculture and Health, Agricultural Research

52    Service, Ithaca, NY 14853, USA

53    16 Fruit Research Institute,Chinese Academy of Agricultural Sciences/Key Laboratory of Biology

54    and Genetic Improvement of Horticultural Crops (Germplasm Resources Utilization), Ministry of

55  Agriculture/Key Laboratory of Mineral Nutrition and Fertilizers Efficient Utilization of Deciduous

56  Fruit Tree, Liaoning Province, Xingcheng, China

57  17 Zhengzhou Fruit Research Institute, Chinese Academy of Agricultural Sciences, Zhengzhou,

58  China

59  18 College of Horticulture, Nanjing Agricultural University, Nanjing 210095, China

60  19 State Key Laboratory of Tropical Crop Breeding, Tropical Crops Genetic Resources Institute,

61  Chinese Academy of Tropical Agricultural Sciences, Haikou, China

62

63  #, these authors contributed equally to this study

64  * Corresponding authors: zhouyongfeng@caas.cn; xiaohua01@caas.cn; zmc@njau.edu.cn;
65  camille.rustenholz@inrae.fr
66

## Abstract

Grapevine is one of the most economically important crops worldwide. However, the previous versions of the grapevine reference genome consisted of thousands of fragments with missing centromeres and telomeres, which limited the accessibility of the repetitive sequences, the centromeric and telomeric regions, and the inheritance of important agronomic traits in these regions. Here, we assembled a telomere-to-telomere (T2T) gap-free reference genome for the pinot noir cultivar (PN40024) using the PacBio HiFi long reads. The T2T reference genome (PN_T2T) was 69 Mb longer with 9026 more genes identified than the 12X.v2 version (Canaguier et al., 2017). We annotated 67% repetitive sequences, 19 centromeres and 36 telomeres, and incorporated gene annotations of previous versions into the PN_T2T. We detected a total of 377 gene clusters, which showed associations with complex traits, such as aroma and disease resistance. Even though the PN40024 sample had been selfed for nine generations, we still found nine genomic hotspots of heterozygous sites associated with biological processes, such as the oxidation-reduction process and protein phosphorylation. The fully annotated complete reference genome, therefore, provides important resources for grapevine genetics and breeding.

87 **Introduction**

88 Since the first human genome was published in 2000, lots of reference genomes have

89 been assembled successively in a variety of species (Lander et al., 2001; Venter et al.,

90 2001; Rice and Green, 2019). The reference genome is essential for biological and

91 genetic studies. Thus, acquiring a high-quality genome has persistently been pursued.

92 However, there are many missing segments due to highly repetitive sequences

93 clustered across the genome, especially three representative regions: telomere,

94 centromere and ribosome DNA (rDNA) (Rice and Green, 2019; Giani et al., 2020;

95 Nurk et al., 2022). Centromere, which hosts CENPA/CENH3-variant nucleosomes

96 and where the kinetochore forms and attaches to spindle microtubules, play an

97 essential role during cell division. It consists of alpha satellites, a highly repetitive

98 DNA sequence. The alpha satellite is composed of monomeric DNA repeats known as

99 Higher Order Repeats (HOR), which contains arranged monomers that ranged from

100 100 to 200 bp (Talbert and Henikoff, 2020; Naish et al., 2021; Sundararajan and

101 Straight, 2022). Despite their conserved function across species, their structure and

102 sequence can change rapidly within and between species, and diverse organizations

103 were observed from one species to another. However, centromere shows concerted

104 evolution within the genome (Liao et al., 2018, Rudd et al., 2006; Melters et al., 2013;

105 Naish et al., 2021). Currently, the centromere remains mostly unknown for

106 researchers. A similar situation also exists in telomeres, which are composed of

107 tandem repeats of relatively conserved microsatellite sequences located at the end of a

108 chromosome in eukaryotes (Fajkus et al., 2005; Podlevsky and Chen, 2016).

109 Telomeres are important for protecting chromosome terminal sequences during cell

110 division (Turner et al., 2019; Coulon and Vaurs, 2020; Yuan et al., 2020; Engin and

111 Engin, 2021). The rDNA is one of the most abundant repetitive elements in the

112 genome that are essential for ribosome formation and play an important role in

113 driving cell growth and cell proliferation (Kobayashi, 2011; Xu et al., 2020; Sasaki

114 and Kobayashi, 2021). Because of the missing information on previously assembled

115  genomes, the investigation of these regions has been extremely limited in the past two

116  decades.

117  Fortunately, benefiting from the improvement of sequencing technology and

118  computational algorithms, the genome assembly ushered in a new era: the telomere-

119  to-telomere genome (T2T genome, Kille et al., 2022). Compared with the fragmented

120  genome, the T2T genome has fewer or no gap on it based on the third-generation

121  sequencing using PacBio high-fidelity long reads (HiFi), ultra-long Oxford Nanopore

122  Technologies (ONT) and Hi-C data. Moreover, the T2T genome includes nearly

123  complete information of telomere, centromere and rDNA regions (Logsdon et al.,

124  2020; Miga and Sullivan, 2021). Promisingly, the T2T genome allows us to access

125  these regions, opening a window into understanding the structure of these regions and

126  the function of genes in these regions. Since the first complete human X chromosome

127  was published in 2020, T2T assembly quickly become a research hotspot (Miga et al.,

128  2020; Logsdon et al., 2021). In plants, the first T2T genome was reported in

129  *Arabidopsis thaliana* in 2021 (Naish et al., 2021; Wang et al., 2022a). Presently, T2T

130  genome assembly has been obtained in several species, such as rice, banana and

131  watermelon, which fascinated the research on genomic structure and function, and

132  crop breeding (Belser et al., 2021; Deng et al., 2022; Zhang et al., 2022; Yue et al.,

133  2022).

134  The grapevine (*Vitis vinifera ssp. vinifera*), a Near East originated fruit tree, is one of

135  the most widely cultivated and economically valuable crops worldwide (Grassi and

136  De Lorenzis, 2021). Domesticated grapes often have highly heterozygous genomes

137  (Zhou et al., 2019), which greatly impeded the acquirement of a high-quality genome.

138  For instance, ~15% of genes were hemizygous in the Chardonnay genome (Zhou et

139  al., 2019). Fortunately, PN40024, a highly homozygous Pinot Noir genotype that

140  originated from successive selfings, was sequenced and the reference genome (8X) of

141  grapevine was first obtained in 2007, which was also the first one for fruit crops

142  (Jaillon et al., 2007). Subsequently, several updated versions have been released: the

ORIGINAL UNEDITED MANUSCRIPT

143 12X.v2 version and its upgraded annotation VCost.v3 in 2017, and the PN40024.v4.1

144 version in 2021 (Canaguier et al., 2017; Navarro-Payá et al., 2021). In addition,

145 fragmented genome assemblies of various grape cultivars have been produced in

146 recent years such as Black Corinth (Massonnet et al., 2020), Cabernet Franc (Vondras

147 et al., 2021; Minio et al., 2022), Cabernet Sauvignon (Chin et al., 2016; Minio et al.,

148 2019b; Minio et al., 2022), Carménère (Minio et al., 2019a), Chardonnay (Roach et al.,

149 2018; Zhou et al., 2019), Merlot (Massonnet et al., 2020), and Nebbiolo (Maestri et al.,

150 2022). As a representative dicotyledonous plant in fruit trees, the high-quality genome

151 will greatly facilitate the research on gene function, genetic structure and evolution of

152 *Vitis* and eudicots species.

153 However, the previous incomplete genome assembly of grapes makes it difficult to

154 access the highly repetitive regions on the genome. Here, we generated a T2T-level

155 gap-free grape reference genome using the PN40024 material and aimed to address

156 the following four questions. First, with the application of third-generation sequencing

157 and assembly technologies, high-fidelity long reads have contributed to gap-free

158 genome assemblies (Cheng et al., 2021; Mascher et al., 2021). Can we complete the

159 reference genome of grapes using these new sequencing and assembly approaches?

160 Second, the studies on centromere, telomere and rDNA have long been neglected. We

161 analyzed the feature, structure, and distribution of these regions based on the

162 assembled gapless grape genome. Third, the annotation of TE and genes in highly

163 repetitive regions were improved based on the T2T genome, which could further

164 improve our understanding of their biological functions, especially the gene clusters.

165 Finally, the PN40024 genome was almost fully homozygous (Jaillon et al., 2007), but

166 some sites remained heterozygous after nine generations of selfing. It is worthwhile to

167 investigate the genomic distribution and the genetic effects of such heterozygous sites.

168 **Results**

169 **A telomere to telomere gap-free reference genome for PN_T2T**

170  PN40024, a highly homozygous Pinot Noir inbred line (Jaillon et al., 2007), was used
171  for T2T genome assembly. In total, 21 Gb (21,024,461,524 bp, ~42X coverage) HiFi
172  reads were generated by the PacBio platform. For the preliminary assembly, HiFiasm
173  was used to assemble the HiFi reads. Using the published grapevine genomes as the
174  reference, we then used Mummer to order the 38 contigs into 19 chromosomes
175  (Figure 1). Only one gap left after initial assembly into contigs (Figure S1). After
176  filling the gap with HiFi reads, a gap-free PN_T2T genome was finally generated
177  (494.87 Mb), which is 69 Mb longer than 12X.v2 (426.18 Mb, Table 1) using the
178  same statistical method. The K-mer was used to evaluate genomic homozygosity,
179  estimated at 99.8% (Figure S2A-D). The BUSCO was used to evaluate genomic
180  completeness, 98.5% of the core conserved plant genes were found complete in the
181  genome assembly (Figure S2E), which is 4.8% more than the 12X.v2 (93.7%, Table
182  1).

183  Compared with the 12X.v2 genome, a substantial improvement was observed in our
184  PN_T2T assembly. The contig N50 length of PN_T2T was ~ 250 times higher than
185  that of 12X.v2 (26.89 Mb VS 102 kb), and all the 9429 gaps in 12X.v2 were filled in
186  PN_T2T genome (Table 1, Table S1, Figure 1A). As shown in Figure 1C, 28 gaps in
187  12X.v2 were filled in PN_T2T with the largest gap being 16,951 bp in the 1Mb
188  syntenic region on chromosome 18 (Figure 1C). Many orientation errors in 12X.v2
189  were also corrected such as inversions and translocations compared to PN_T2T
190  (Figure 1A, Figure S3). For example, two large inversions, which were located
191  surrounding the centromere of chromosome 3 and at the ends of chromosome 5, with
192  the length of 0.9 M and 4.9 M were observed between two versions of assemblies,
193  respectively (Figure 1A, Figure 1B). Moreover, 19 centromeres and 36 out of the 38
194  telomeres were detected on the PN_T2T genome assembly, except one telomere on
195  chromosome 15 and one telomere on chromosome 17, which were missing in all
196  previous grape genome assemblies. A total of 37,534 genes and 41,064 transcripts
197  were annotated, among which 24526 (86.01%), 27696 (78.83%), 27717 (78.75%)
198  were shared with older versions PN40024.V2.1 (https://phytozome-

199 next.jgi.doe.gov/info/Vvinifera_v2_1), PN40024.v4, and PN40024.v4.1

200 (https://integrape.eu/resources/genes-genomes/genome-accessions/), respectively

201 (Table S2). A total of 5472 (14.58%) genes were not found to correspond in any of

202 the three versions. A total of 97.9% single-copy genes completely assembled was

203 assessed by the BUSCO analysis, and structural domains were detected in 35508

204 sequences out of 40307 unique sequences (88.1%) while PN40024.V4.1 has 38364

205 unique sequences, and 29688 sequences were detected with structural domains

206 (77.4%, Table S2).

207 Based on the species-specific Pan-TE database constructed by RepeatModeler2, the

208 repeats were detected with a pipeline shown in Figure 2A. Finally, 66.47% of our

209 gap-free grape genome was marked as repetitive sequences (Figure 1D). As a

210 comparison, 62.47% of the repetitive sequences were identified in 12X.v2 genome

211 using the same pipeline (Table S3). Among the repeats predicted in PN_T2T genome,

212 the largest portion is transposable elements (TEs, 63.90%) with a total length of 316

213 Mb (59.96% and 292 Mb in 12X.v2). The TEs mainly consisted of the long terminal

214 repeat (LTR) type (47.54%), predominantly Gypsy (20.22%) and Copia (19.67%)

215 elements. In total, we detected 276 rDNA sequences, representing 0.019% of the

216 genome.

217 **The telomeres and centromeres**

218 To access the telomeric and centromeric regions in PN_T2T, we identified the

219 telomeres and centromeres using the pipeline in Figure 2A. For telomeres, we

220 checked the 150 kb sequences at both ends of each chromosome, and the length of the

221 telomere repeat unit was set to range from 5 to 12 bp. Finally, the telomere repeat unit

222 (TTTAGGG/CCCTAAA) was detected, which was the most abundant in the genome

223 and carried by all chromosomes. The same telomere repeat unit was reported in

224 grapes by Melters et al. 2013 and Castro et al. 2021. We further predicted the

225 telomeres in 36 out of 38 telomeres in the PN T2T genome, except the short arms of

226 chromosome 15 and chromosome 17 (Figure 1A, Figure 2B and Table S4). Among

227    them, the longest telomere (31 kb) was in the short arm of chromosome 8 with 4,479

228    repeats while the shortest telomere (1,260 bp) was in the long arm of chromosome 7

229    with only 180 repeats.

230    To detect the centromeric region, we scanned candidate repeats from 30 to 500 bp

231    along the genome. The Tandem Repeats Finder (TRF) found 470 different repeat

232    units in the PN_T2T genome, of which the 107 bp repeats were the most abundant

233    unit in the whole genome, which had 182,620.5 (copies ≥ 2) repetitions accounted for

234    about 3.95% of the genome, followed by 321 bp (2.45%), 214 bp (1.94%), and 135 bp

235    (1.05%) (Figure 3A). Interestingly, we found the sequences of 214 bp and 321 bp

236    repeat units consisted of two and three copies of the 107 bp repeat unit, respectively.

237    The TE analyses also support the centromeric feature of the 107bp repetitive region

238    (Figure 2). Thus, the centromeres were recognized mainly based on 107 bp repeat

239    units, and localized on all 19 chromosomes (Figure 1A, Figure 2B, and Table S5). As

240    shown in Figure 3B, the total length of 107 bp repeats varied from 1.4 kb to 3.8 Mb,

241    but the sequences of the 107 bp repeats were highly conserved among chromosomes

242    (Figure 3C). The 107 bp repeats were the most abundant in all chromosomes, except

243    chromosomes 3, 14 and 18 (Figure 3D-H, Table S6). We found that the 187 bp was

244    the main repeat unit in chromosome 14, it was scattered throughout the whole

245    chromosome, and that 51 bp, 56 bp, 105 bp and 107 bp repeat units were highly

246    overlapped and enriched in the centromere, which showed a core region of

247    chromosome. The centromeric repeat unit in chromosome 3 was the 135 bp repeats

248    and its integer multiples (270 bp and 405 bp). As for chromosome 18, 66 bp and its

249    integer multiple 132 bp were the main repeat units (Figure S4).

250    To locate the centromeric repeats, we further examined the relationship between TE

251    and centromeres. LTR retrotransposons or centromeric retrotransposons (CRs), were

252    usually mixed with tandem repeats and enriched in plant centromeric regions (Guo et

253    al., 2016; Fernandes et al., 2019). We found (Figure 4A) that the genes and TE repeats,

254    such as LTR (Gypsy and Copia), DNA TE(MULE-MuDR) and RC (Helitron), have a

ORIGINAL UNEDITED MANUSCRIPT

255  low density in the special region where the enormous centromeric tandem repeats

256  enriched in the chromosome were viewed in IGV (Figure 2 and Figure S4). We then

257  inferred the region with centromeric repeats and low TE density as the centromeres

258  after zooming one by one (Figure S4 and Table S5). The pattern of 107 bp was the

259  target, which was highly linked with the centromeric region in grapes. However, there

260  were likely different repeat units and patterns that appeared on chromosomes 3, 14,

261  and 18 (Figure 3F-H). The scattering of transposons and the distribution of the

262  centromere showed that specific sequence-defined repeat superfamilies are correlated

263  or anticorrelated to various levels with centromeric proximity (Figure 2B, Figure 4A),

264  forming density gradients that are the main chromosome-scale repeat-associated

265  features, presumably reflecting overall chromatin structure (Figure S4).

266  To detect the captured genes, we then screened all genes in these regions in the highly

267  linked centromeric region. Interestingly, we found 343 genes (Table S7) captured in

268  the centromeres, which included 179 genes with Uniprot ID through blastp. Through

269  GO (Gene Ontology) functional annotation, 12 genes were enriched in protein binding

270  (molecular function, MF), such as *VvAMP1* (*Vitis01g01298* and *Vitis13g01021*;

271  Uniprot ID: Q9M1S8) involved in ethylene (ETH), gibberellin (GA), and abscisic

272  acid (ABA) signaling pathways (Saibo et al., 2007; Shi et al., 2013). In addition, we

273  found 10 genes enriched in the cellular component (CC) of the cytosol, mitochondrion

274  and cytoplasm respectively, including auxin transport protein *VvBIG* (*Vitis02g01141*;

275  Uniprot ID: Q9SRU2), which influences general growth and development in plants

276  (Gil et al., 2001); fumarate hydratase 1 *VvFUM1* (*Vitis02g01128*; Uniprot ID: P93033)

277  catalyzes the active of mitochondrial Krebs cycle-associated enzyme (Zubimendi et

278  al., 2018); 6-phosphogluconate dehydrogenase, decarboxylating 2 *VvPGD2*

279  (*Vitis02g01123*, Uniprot ID: Q9FWA3) plays a key role in the development of the

280  male gametophytes and interaction between the pollen tube and the ovule (Hölscher et

281  al., 2016). Moreover, RNA modification, protein autophosphorylation, DNA

282  integration, DNA recombination and photomorphogenesis were enriched in the

283  biological process (BP) (Figure 4C).

284 **Gene clusters in the grapevine reference genome**

285 To infer the gene clusters in the grapevine genome, protein-to-protein alignments

286 among the Pinot protein coding gene exposed a rich panoply of duplication structures

287 in terms of genomic positions and functions. Prominent and complex tandem-like

288 blocks of high-similarity genes can be seen via visualizations of all–vs.–all alignments

289 (Figure S5). We found a total of 377 gene clusters in the grapevine reference genome

290 (Table S8). These duplications often involve local rearrangements and can extend into

291 megabases with dozens to hundreds of genes involved (Figure 5). On chromosome 16

292 (23-27 Mb), there were 599 genes enriched domains mainly including WAKs (Wall

293 associated receptor kinase galacturonan binding), PPR repeat, Leucine-rich, ABC

294 transporter, Intergrase domain, Peptidase family, Protein kinase and Reverse

295 transcriptase (Figure 5A). And on chromosome 18 (25~36 Mb), there were 1237

296 genes enriched domains mainly including Intergrase domain, C JID domain, NB ARC

297 domain, Leucine rich repeat, Multicopper oxidase, Reverse transcriptase, Terpene

298 synthase and TIR. The results showed that many of the strongly enriched structural

299 domains were part of the structural domains of plant disease resistance genes (R

300 genes), including NB-ARC, TIR and structures identified by the Colis database. We

301 analyzed the domain architecture of our 41,064 PN_T2T PCGs and identified 3,381

302 possible R genes. Collectively, these R genes and gene clusters in grapes indicated a

303 tremendous opportunity for exploring plant defense mechanisms.

304 **The genetic heterozygosity after the ninth generation of selfing**

305 We are interested in the genomic changes associated with the inbreeding process.

306 Based on the reference genome of PN_T2T, the resequencing data of four PN40024

307 clones were downloaded from NCBI and analyzed (Jaillon et al., 2007, Magris et al.,

308 2019). A total of 244,215 SNPs were detected, among which 208,330 SNPs (85.3%)

309 were shared in all four samples while the other 35,886 SNPs were only presented in 1-

310 3 samples (Figure 6A). Interestingly, we found nine hotspots of heterozygous SNPs

311 on chromosomes 1, 2, 3, 4, 7, 10, 11 and 16 (Figure 5A, Figure S6). To further

312 investigate the highly heterozygous region, we examined the top 5% heterozygosity

313  windows and identified a total of nine large continuous fragments (chromosome 1:

314  1.1-1.3 M, chromosome 2: 4.2-7.2 M, chromosome 3: 9.4-9.9 M, chromosome 4:

315  21.8-22.9 M, chromosome 7: 15.3-26.2 M, chromosome 10: 0.7-6.5 M, 17.6-18.3 M,

316  chromosome 11: 7.1-7.8 M, chromosome 16: 13.0-13.5 M). The GO enrichment

317  analysis on the genes in these regions showed that the most significantly enriched

318  terms were response to water deprivation, protein phosphorylation, cell division,

319  response to oxidative stress and response to salt stress, which were closely associated

320  with key physiological activities in plants (Table S9, Figure 6C, Figure S7).

321  **Discussion**

322  The complete reference genome is essential for crop genetics and breeding. The

323  previous versions of the grapevine reference genome have thousands of gaps with

324  errors in repetitive regions and missing centromeres and telomeres, which limited the

325  access of variants within these regions. Sometimes, such unreachable regions are

326  underlying QTL of important agronomic traits, such as the berry color and sex

327  determination on chromosome 2 (Fournier-Level et al., 2009; Zhou et al. 2017; Zhou

328  et al. 2019; Zou et al. 2021) and disease resistance on chromosome 14 (Riaz et al.,

329  2008; Morales-Cruz et al. 2022). The complete reference genome has great potential

330  to reveal the missing heritability of important polygenic agronomic traits, therefore, it

331  could increase the genetic gain in grapevine breeding. More and more investigations

332  suggested the important functions of gene clusters, a total of 377 gene clusters were

333  detected in the PN_T2T. The grapevine genome is also widely used in studies of plant

334  evolution and comparative genomics because of its important phylogenetic position

335  on the evolution of eudicots (Jaillon et al., 2007). The T2T version could be widely

336  used in plant evolutionary genomics, especially at the repetitive sequences,

337  centromeres and telomeres. The T2T gap-free reference genome had incorporated

338  gene annotations of previous versions with more accurate TE annotation (up to ~ 67%

339  of the genome), which will be an important resource for grapevine functional

340  genomics and breeding.

**The architecture and context of plant centromeres**

The centromeric region ranges from Kbs to Gbs in length, including > 90% tandem repeats (McKinley and Cheeseman, 2016). The centromere is among the last pieces of great unknowns in genomics since it was inaccessible by previous sequencing technologies. The assemblies often collapse due to the highly repetitive nature of the centromeric region. We assembled and annotated centromeres for all 19 chromosomes of the grapevine genome (Figure 1). Most of the chromosomes have a single centromere while others could have multiple centromeric regions so called holocentromere (Steiner and Henikoff, 2014; Hofstatter et al., 2022). On chromosomes 16 and 18, we found tandem repeats in many regions while on other chromosomes only a single peak was detected (Figure 2B), suggesting that the structure of the centromeric region might be more complicated and requires further investigations.

In the PN40024 grapevine reference genome, there are three major repetitive patterns in the 19 chromosomes, suggesting different chromosomal evolutionary histories (Figure 3D-H). On chromosomes 3, 14 and 18, we found 135 bp, 56bp and 66 bp tandem repeats, respectively (Figure S4), while on other chromosomes, the major unit of tandem repeats is 107 bp (Figure 3D-H, Figure 4B, Figure 4D). The evolutionary histories of the grapevine centromere of each chromosome are still an open question to be addressed with all *Vitis* genomes. Previous comparative genomic analyses suggested that centromere is conservative among closely related species with a constant number of chromosomes (Liao et al., 2018). The transformation of centromeric structures happens during the chromosomal division and fusion when the number of chromosomes changes in evolution. The muscadine grape has 20 chromosomes with chromosomes 7 and 13 collinear with *Vitis* chromosome 7, which is associated with a chromosome fusion event (Cochetel et al., 2021). Only one centromeric region is left on chromosome 7 in our grapevine reference genome (Figure 2B, Figure S4), suggesting one centromere was lost during the evolvement of the genus *Vitis*.

370  The centromeric architecture shaped the content within the genome, population

371  genetic diversity within species and genetic differentiation among species. Population

372  genetic analyses revealed that the genetic variants in the centromeric region are highly

373  linked with much lower genetic diversity compared to chromosome arms (Kawabe et

374  al., 2008). The centromeres capture tens to thousands of genes that are highly linked

375  with the centromeric tandem repeats. These genes along with the centromeric region

376  are functional as supergenes. In total, we found 343 captured genes (Table S7) in the

377  centromeric region in the grapevine reference genome. Interestingly, the genes are

378  mainly involved in the ethylene (ETH), gibberellin (GA), and abscisic acid (ABA)

379  signaling pathways (Saibo et al., 2007; Shi et al., 2013).

380  **The hotspots of heterozygous variants in selfing plants**

381  The current grapevine reference genome was generated by a Pinot Noir sample

382  (PN40024) selfed for nine generations with high homozygosity at 99.8% of the

383  genome (Figure S2A-D). However, we still interested in the remained heterozygous

384  sites. Thus, we collected Illumina resequencing reads for four clones of PN40024

385  maintained in different international labs. Interestingly, the heterozygous SNPs and

386  SVs were enriched in specific regions when mapped to the PN_T2T. In total, we

387  found 208,330 heterozygous SNPs shared by the four samples, and 35,886 SNPs

388  specific to 1-3 samples. The former is more likely the original variants of PN40024

389  after nine generations of selfing while the latter could be somatic mutations generated

390  during the distribution and tissue culture in different labs. Interestingly, we found the

391  hotspots of common variants were enriched in central biological processes including

392  the oxidation-reduction process and protein phosphorylation. The hotspots on

393  chromosome 2 also cover the sex-determination QTL region (Figure 6), which

394  complicated the mining of the sex-determination genes (Zhou et al. 2019; Zou et al.,

395  2021), because the candidate genes were not presented in the old version of the

396  reference genome. It has been reported that during the clonal reproduction of fruit

397  trees, such heterozygous deleterious variants accumulate in the genome (Zhou et al.

398  2019; Xiao et al., 2023). The clonal processes hide recessive deleterious variants

399   including small SNPs and indels and large structural variants in a heterozygous state

400   (Zhou et al., 2017; Zhou et al., 2019). A strong inbreeding depression has been

401   commonly observed in clonal crops, including potato, cassava, citrus and grapevine

402   (Zhou et al., 2017; Ramu et al., 2017; Zhang et al., 2021; Wang et al., 2022b;) since

403   the strongly deleterious variants in these genomic regions has been exposed to lethal

404   or strong recessive selection during selfing cycles. In grapevine breeding, the

405   inbreeding and outcrossing depression were commonly detected because the hidden

406   heterozygous recessive deleterious variants increased during clonal propagation has

407   been exposed during sex reproduction.

408   **Methods**

409   **Sample collection and genome sequencing**

410   PN40024 is a line that belongs to one of the near homozygous lines originally derived

411   from Pinot Noir by successive selfing steps, estimated the close to 97% homozygosity

412   tested by SSR markers (Jaillon et al., 2007). We got this inbred material from INRAE

413   under Material Transfer Agreement (MTA) and transplanted it in the greenhouse

414   belonged to AGIS (Agricultural Genomics Institute at Shenzhen, Chinese Academy of

415   Agricultural Sciences, Shenzhen, China) for subsequent experiments. Young leaves

416   and ovules from PN40024 were flash-frozen in liquid nitrogen. Genomic DNA and

417   RNA were isolated using the DNeasy Plant Mini kit (Qiagen) following the

418   manufacturer's instructions. For PacBio HIFI sequencing, two single-molecule real-

419   time cells were sequenced on a PacBio Sequel II platform, and a total of 21 Gb of

420   HiFi read was generated using CCS (https://github.com/PacificBiosciences/ccs) with

421   the default parameter for the sequenced accessions. From each RNA-seq sample,

422   isolate poly (A) mRNA 10 μg of total RNA was used to prepare Illumina RNA-seq

423   libraries. These libraries were then sequenced using the Illumina HiSeqTM 2000

424   system in accordance with the manufacturer's instructions.

425   **T2T genome assembly**

426    Initially, PN40024 was assembled genome by incorporating PacBio single-molecule

427    real-time long-read sequences. Reads generated by the PacBio Sequel II platform

428    were self-corrected, trimmed and assembled by hifiasm, using default parameters

429    (https://github.com/chhylp123/hifiasm, Cheng et al., 2021). The initial output of

430    hifiasm (v.0.13) yielded the p_ctg draft assembly. Genome heterozygosity was

431    estimated using a k-mer-based approach by GenomeScope2.0 (Ranallo-Benavidez et

432    al., 2020), estimated close to 99.8% homozygosity (Figure S2A-D). Then, homology-

433    based scaffolds were generated with MUMmer (v.4.0.0) (Marcais et al., 2018)

434    "scaffold", using the 12X.v2 reference genome (Figure S3). By applying MUMmer

435    tools, we order and orient the contig-level assemblies into 19 chromosomes, and join

436    the adjacent contigs to generate a scaffold with 100 N. Finally, we adjusted the

437    assembly manually through aligning the genome sequencing data from previous

438    version of PN40024, which was mapped to the genome assembly by Minimap2

439    (v.2.21) and visualized in IGV (v.2.12.3) software to observe whether the gap regions

440    were supported by reads (Figure S1). The filling and close of the gaps with the

441    selected and assigned contigs were performed by mapping the 50 bp-sequences

442    around the gap to continuous long reads (CLR) of PN40024.v4 and obtaining the

443    gapless telomere-to-telomere PN40024 assembly for all 19 grape chromosomes.

444    Assembly was inspected based on BUSCO (Simão et al., 2015) completeness and the

445    duplication score.

446    **The annotation of genes and TEs**

447    We have used an self-developed method for genome annotation. The putative genes

448    were first searched for by using transcripts and uniprot as evidence. A preliminary

449    gene model was then built for the putative genes and the further search was performed

450    using AUGUSTUS (V3.4.0) (Stanke et al., 2006). All the found putative genes

451    fragments were then filtered, including genes involving duplicated regions, genes with

452    CDS lengths shorter than 90 and genes not supported by any evidence. The missing

453    genes were attempted to be complemented and the complete genes were subjected to

454    the alternative splicing analyses. Finally, all the results were examined by hidden

455 Markov models downloaded from the Pfam database to obtain the final gene models.

456 Interproscan (v.5.56-89.0, Jones et al., 2014) was used to function annotation for our

457 assembly, Pfam (v.34.0, Mistry et al., 2021) and Coils (v.2.2.1, Fitzkee et al., 2005)

458 was used for the identification of structural domains (https://github.com/unavailable-

459 2374/Genome-Wide-Annotation-Pipeline).

460 The primary repeat analysis was outlined in Figure 2A and began with the

461 construction of a Pan-Vitis database of repeat families by RepeatModeler (open-2.0.3,

462 Flynn et al., 2020) and a series of scripts, which was then applied with RepeatMasker

463 (open-4.1.2). For building this Pan-Vitis repeat database we download 17 Vitis

464 genomes from NCBI, then use RepeatModeler2 to identify TE family. After that, we

465 got 17 consensus fasta files of TE family, by removing the single copy and failed

466 annotations we aggregated these files. We used NCBI-BLAST+2.9.0 (Altschul et al.,

467 1990) to remove some redundancy sequences (-i 80%, -l 80%). After all, we got the

468 final file of repeat identity, then we used deepTE (Yan et al., 2020) with the Plant

469 model to classify those unclassified repeat elements. Finally, the repetitive sequence

470 of the complete reference genome was annotated by RepeatMasker.

471 **Genome comparison between different versions of the grapevine reference**

472 **genome**

473 To compare previous versions of grapevine genome with PN_T2T, we align the

474 genomes using minimap2 and index the alignment BAM file using

475 samtools(minimap2 -ax asm5 -t 4 --eqx A.fa B.fa | samtools sort -O BAM - >

476 A_B.bam, samtools index A_B.bam). Next, to detect structural variations between

477 genomes, we need to find synteny and structural rearrangements between the genomes.

478 For this, we use SyRI: (syri -c A_B.bam -r A.fa -q B.fa -F B --prefix A_B). Finally,

479 Plotsr were used to generate the graph: plotsr --sr A_Bsyri.out --sr B_Csyri.out --sr

480 C_Dsyri.out --genomes genomes.txt -o output_plot.pdf,

481 https://github.com/schneebergerlab/plotsr). MUMmer (v.4.0.0) was used to compare

482 the 12X.v2 genome with the reference genome PN_T2T using whole-genome

483 alignments (Marçais et al., 2018). First, we aligned the two genome sequences using

484    nucmer (nucmer --mum) and then filtered one-to-one alignments with a minimum

485    alignment length of 10,000 bp (delta-filter -i 95 -l 10000).

486    Samtools (v.1.7) were used to extract the sequence of chromosome 18: 25.0-26.0 Mb

487    in 12X.v2 and aligned the sequence in PN_T2T. The gap information was detected

488    with a python script (getgaps.py) and finally used LINKVIEW2

489    (https://github.com/YangJianshun/LINKVIEW2) to visualize the alignment results.

490    **The identification of telomeres and centromeres**

491    The telomere repeat units were explored by using the TIDK (v.0.2.0)

492    (https://github.com/tolkit/telomeric-identifier) with options: tidk explore -f genome.fa

493    --minimum 5 --maximum 12 -o tidk_explore -t 2 --log --dir telemere_find --extension

494    TSV. Then the whole genome was searched using the parameter: tidk search -f

495    genome.fa -s TTTAGGG -o tidk_search --dir telemere_find. Finally, we completed

496    the rapid statistics of telomere based on the tidk plot and used R script to visualize the

497    telomere peak.

498    For centromere annotation, the TRF (v.4.09) (Benson, 1999) was used to finish

499    tandem repeats annotation with the parameter: trf genome.fa 2 7 7 80 10 50 500 -f -d -

500    m, and then we merged the results of annotation by using TRF2GFF

501    (https://github.com/Adamtaranto/TRF2GFF). To complete data statistics and

502    visualization, we performed information extracted by using awk command in the

503    linux system and analyzed the results in IGV (v.2.12.3) (Thorvaldsdóttir et al., 2013).

504    We used four software to show more details about the centromeric region: Iqtree (v.

505    2.1.4-beta) (Minh et al., 2020) was used to achieve the phylogenetic tree (options: -m

506    GTR+I+G -bb 1000 -bnni -alrt 1000); itol (v.6) (Letunic and Bork, 2021) was used to

507    visualize the phylogenetic tree; GeneDoc (v.2.7.0)

508    ( https://github.com/karlnicholas/genedoc) was to achieve multiple sequence

509    alignment; R script was used to plot the data statistics and typeset details respectively.

510    To detect the functions of the genes captured in the centromeric regions, we

511    downloaded the protein sequence library of Swissprot (2022/08/30,

512    https://ftp.ncbi.nlm.nih.gov/blast/db/FASTA/ ）for a local blast. After this, we

513   extracted all the protein sequences of PN_T2T blasted by diamond (v.2.0.15)

514   (parameter: - k 1 - e 0.00001, https://github.com/python-diamond/Diamond). We

515   further uploaded the SwissProt ID to DAVID (https://david.ncifcrf.gov/tools.jsp) and

516   completed Gene Ontology enrichment and annotation. Finally, Data visualization was

517   completed by our R scripts.

**The identification of gene clusters**

519   To define the clustered genes in the reference genome, protein sequences were

520   extracted using gffread and then filtered by e-value less than 1e-5 and similarity

521   greater than 30% using blastp for all-vs-all alignments. The filtered alignment results

522   were combined with functional annotations to filter out alignment results that did not

523   share the same structural domains. Finally, we determined the presence of gene

524   clusters by identifying three consecutive identical PF numbers, using such PF

525   numbers as seeds, and going up and down 30 genes to find genes with the same PF

526   number. In total, 377 gene clusters were found (Table S8).

**The heterozygosity in selfed PN40024 clones**

528   Four resequencing samples were downloaded from NCBI database (SRR6156373,

529   SRR8835144, SRR8835157, SRR8835168) and mapped to newly assembled PN_T2T

530   genome for SNP calling. Quality-controlled reads were mapped to the genome using

531   bwa (v.0.7.15) with the default parameters. SAMtools (v.1.4) and GATK (v.4.1.8)

532   were used for sorting and indexing the bam file with no duplicates. The gvcf files

533   were combined in GATK and were used to join calling SNPs across all samples. To

534   obtain high-quality SNPs, we performed strict filtering of the SNP calls based on the

535   following criteria: (1) the SNPs with more than two alleles were removed in all

536   samples in vcftools with parameters --min-alleles 2 --max-alleles 2; (2) we removed

537   the SNPs with quality scores (GQ) less than 30 (--minGQ 30) and the missing rate is 0

538   (--max-missing 1); (3) SNPs had minor allele frequencies (MAFs) $\geq 0.01$ to remove

539   the invariable sites.

**Data availability**

541    All PacBio sequence data have been deposited to the NCBI Sequence Read Archive

542    under the project number: PRJNA882193 and the National Genomics Data Center

543    (NGDC) Genome Sequence Archive (GSA) (https://ngdc.cncb.ac.cn/gsa/), with

544    BioProject number PRJCA012093. The assembly and annotation have been deposited

545    to zenodo: https://zenodo.org/record/7751391#.ZBgVmcJBy3A.

546    **Code availability**

547    All the scripts and pipelines used in this study have been achieved in GitHub:

548    https://github.com/zhouyflab

549

550    **Reference**

551    **Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J.** (1990). Basic

552    local alignment search tool. Journal of molecular biology **215**:403-410.

553    10.1016/s0022-2836(05)80360-2.

554    **Belser, C., Baurens, F.C., Noel, B., Martin, G., Cruaud, C., Istace, B., Yahiaoui,**

555    **N., Labadie, K., Hřibová, E., Doležel, J., et al.** (2021). Telomere-to-telomere

556    gapless chromosomes of banana using nanopore sequencing. Communications

557    biology **4**:1047. 10.1038/s42003-021-02559-3.

558    **Benson, G.** (1999). Tandem repeats finder: a program to analyze DNA sequences.

559    Nucleic acids research **27**:573-580. 10.1093/nar/27.2.573.

560    **Canaguier, A., Grimplet, J., Di Gaspero, G., Scalabrin, S., Duchêne, E., Choisne,**

561    **N., Mohellibi, N., Guichard, C., Rombauts, S., Le Clainche, I., et al.** (2017). A

562    new version of the grapevine reference genome assembly (12X.v2) and of its

563    annotation (VCost.v3). Genomics data **14**:56-62. 10.1016/j.gdata.2017.09.002.

564    **Castro, C., Carvalho, A., Gaivão, I., and Lima-Brito, J.** (2021). Evaluation of

565    copper-induced DNA damage in Vitis vinifera L. using Comet-FISH. Environmental

566    science and pollution research international **28**:6600-6610. 10.1007/s11356-020-

567    10995-7.

568    **Cheng, H., Concepcion, G.T., Feng, X., Zhang, H., and Li, H.** (2021). Haplotype-

569    resolved de novo assembly using phased assembly graphs with hifiasm. Nature

methods **18**:170-175. 10.1038/s41592-020-01056-5.

Chin, C.S., Peluso, P., Sedlazeck, F.J., Nattestad, M., Concepcion, G.T., Clum, A., Dunn, C., O'Malley, R., Figueroa-Balderas, R., Morales-Cruz, A., et al. (2016). Phased diploid genome assembly with single-molecule real-time sequencing. Nature methods **13**:1050-1054. 10.1038/nmeth.4035.

Cochetel, N., Minio, A., Massonnet, M., Vondras, A.M., Figueroa-Balderas, R., and Cantu, D. (2021). Diploid chromosome-scale assembly of the Muscadinia rotundifolia genome supports chromosome fusion and disease resistance gene expansion during Vitis and Muscadinia divergence. G3 (Bethesda, Md.) **11**10.1093/g3journal/jkab033.

Coulon, S., and Vaurs, M. (2020). Telomeric Transcription and Telomere Rearrangements in Quiescent Cells. Journal of molecular biology **432**:4220-4231. 10.1016/j.jmb.2020.01.034.

Deng, Y., Liu, S., Zhang, Y., Tan, J., Li, X., Chu, X., Xu, B., Tian, Y., Sun, Y., Li, B., et al. (2022). A telomere-to-telomere gap-free reference genome of watermelon and its mutation library provides important resources for gene discovery and breeding. Molecular plant **15**:1268-1284. 10.1016/j.molp.2022.06.010.

Engin, A.B., and Engin, A. (2021). The Connection Between Cell Fate and Telomere. Advances in experimental medicine and biology **1275**:71-100. 10.1007/978-3-030-49844-3_3.

Fajkus, J., Sýkorová, E., and Leitch, A.R. (2005). Telomeres in evolution and evolution of telomeres. Chromosome research: an international journal on the molecular, supramolecular and evolutionary aspects of chromosome biology **13**:469-479. 10.1007/s10577-005-0997-2.

Fernandes, J.B., Wlodzimierz, P., and Henderson, I.R. (2019). Meiotic recombination within plant centromeres. Current opinion in plant biology **48**:26-35. 10.1016/j.pbi.2019.02.008.

Fitzkee, N.C., Fleming, P.J., and Rose, G.D. (2005). The Protein Coil Library: a structural database of nonhelix, nonstrand fragments derived from the PDB. Proteins

599    **58**:852-854. 10.1002/prot.20394.

600    **Flynn, J.M., Hubley, R., Goubert, C., Rosen, J., Clark, A.G., Feschotte, C., and**

601    **Smit, A.F.** (2020). RepeatModeler2 for the automated genomic discovery of

602    transposable element families. Proceedings of the National Academy of Sciences of

603    the United States of America **117**:9451-9457. 10.1073/pnas.1921046117.

604    **Fournier-Level, A., Le Cunff, L., Gomez, C., Doligez, A., Ageorges, A., Roux, C.,**

605    **Bertrand, Y., Souquet, J.M., Cheynier, V., and This, P.** (2009). Quantitative genetic

606    bases of anthocyanin variation in grape (Vitis vinifera L. ssp. sativa) berry: a

607    quantitative trait locus to quantitative trait nucleotide integrated study. Genetics

608    **183**:1127-1139. 10.1534/genetics.109.103929.

609    **Giani, A.M., Gallo, G.R., Gianfranceschi, L., and Formenti, G.** (2020). Long walk

610    to genomics: History and current approaches to genome sequencing and assembly.

611    Computational and structural biotechnology journal **18**:9-19.

612    10.1016/j.csbj.2019.11.002.

613    **Grassi, F., and De Lorenzis, G.** (2021). Back to the Origins: Background and

614    Perspectives of Grapevine Domestication. International journal of molecular sciences

615    **22**10.3390/ijms22094518.

616    **Guo, X., Su, H., Shi, Q., Fu, S., Wang, J., Zhang, X., Hu, Z., and Han, F.** (2016).

617    De Novo Centromere Formation and Centromeric Sequence Expansion in Wheat and

618    its Wide Hybrids. PLoS genetics **12**:e1005997. 10.1371/journal.pgen.1005997.

619    **Hofstatter, P.G., Thangavel, G., Lux, T., Neumann, P., Vondrak, T., Novak, P.,**

620    **Zhang, M., Costa, L., Castellani, M., Scott, A., et al.** (2022). Repeat-based

621    holocentromeres influence genome architecture and karyotype evolution. Cell

622    **185**:3153-3168.e3118. 10.1016/j.cell.2022.06.045.

623    **Holt, C., and Yandell, M.** (2011). MAKER2: an annotation pipeline and genome-

624    database management tool for second-generation genome projects. BMC

625    bioinformatics **12**:491. 10.1186/1471-2105-12-491.

626    **Jaillon, O., Aury, J.M., Noel, B., Policriti, A., Clepet, C., Casagrande, A.,**

627    **Choisne, N., Aubourg, S., Vitulo, N., Jubin, C., et al.** (2007). The grapevine

628    genome sequence suggests ancestral hexaploidization in major angiosperm phyla.

629    Nature **449**:463-467. 10.1038/nature06148.

630    **Jones, P., Binns, D., Chang, H.Y., Fraser, M., Li, W., McAnulla, C., McWilliam,**

631    **H., Maslen, J., Mitchell, A., Nuka, G., et al.** (2014). InterProScan 5: genome-scale

632    protein function classification. Bioinformatics (Oxford, England) **30**:1236-1240.

633    10.1093/bioinformatics/btu031.

634    **Kawabe, A., Forrest, A., Wright, S.I., and Charlesworth, D.** (2008). High DNA

635    sequence diversity in pericentromeric genes of the plant Arabidopsis lyrata. Genetics

636    **179**:985-995. 10.1534/genetics.107.085282.

637    **Kille, B., Balaji, A., Sedlazeck, F.J., Nute, M., and Treangen, T.J.** (2022). Multiple

638    genome alignment in the telomere-to-telomere assembly era. Genome biology **23**:182.

639    10.1186/s13059-022-02735-6.

640    **Kobayashi, T.** (2011). How does genome instability affect lifespan?: roles of rDNA

641    and telomeres. Genes to cells : devoted to molecular & cellular mechanisms **16**:617-

642    624. 10.1111/j.1365-2443.2011.01519.x.

643    **Korf, I.** (2004). Gene finding in novel genomes. BMC bioinformatics **5**:59.

644    10.1186/1471-2105-5-59.

645    **Lander, E.S., Linton, L.M., Birren, B., Nusbaum, C., Zody, M.C., Baldwin, J.,**

646    **Devon, K., Dewar, K., Doyle, M., FitzHugh, W., et al.** (2001). Initial sequencing

647    and analysis of the human genome. Nature **409**:860-921. 10.1038/35057062.

648    **Letunic, I., and Bork, P.** (2021). Interactive Tree Of Life (iTOL) v5: an online tool

649    for phylogenetic tree display and annotation. Nucleic acids research **49**:W293-w296.

650    10.1093/nar/gkab301.

651    **Li, W., and Godzik, A.** (2006). Cd-hit: a fast program for clustering and comparing

652    large sets of protein or nucleotide sequences. Bioinformatics (Oxford, England)

653    **22**:1658-1659. 10.1093/bioinformatics/btl158.

654    **Liao, Y., Zhang, X., Li, B., Liu, T., Chen, J., Bai, Z., Wang, M., Shi, J., Walling,**

655    **J.G., Wing, R.A., et al.** (2018). Comparison of Oryza sativa and Oryza brachyantha

656    Genomes Reveals Selection-Driven Gene Escape from the Centromeric Regions. The

657    Plant cell **30**:1729-1744. 10.1105/tpc.18.00163.

658    **Logsdon, G.A., Vollger, M.R., and Eichler, E.E.** (2020). Long-read human genome

659    sequencing and its applications. Nature reviews. Genetics **21**:597-614.

660    10.1038/s41576-020-0236-x.

661    **Logsdon, G.A., Vollger, M.R., Hsieh, P., Mao, Y., Liskovykh, M.A., Koren, S.,**

662    **Nurk, S., Mercuri, L., Dishuck, P.C., Rhie, A., et al.** (2021). The structure, function

663    and evolution of a complete human chromosome 8. Nature **593**:101-107.

664    10.1038/s41586-021-03420-7.

665    **Maestri, S., Gambino, G., Lopatriello, G., Minio, A., Perrone, I., Cosentino, E.,**

666    **Giovannone, B., Marcolungo, L., Alfano, M., Rombauts, S., et al.** (2022).

667    'Nebbiolo' genome assembly allows surveying the occurrence and functional

668    implications of genomic structural variations in grapevines (Vitis vinifera L.). BMC

669    genomics **23**:159. 10.1186/s12864-022-08389-9.

670    **Marçais, G., Delcher, A.L., Phillippy, A.M., Coston, R., Salzberg, S.L., and**

671    **Zimin, A.** (2018). MUMmer4: A fast and versatile genome alignment system. PLoS

672    computational biology **14**:e1005944. 10.1371/journal.pcbi.1005944.

673    **Mascher, M., Wicker, T., Jenkins, J., Plott, C., Lux, T., Koh, C.S., Ens, J.,**

674    **Gundlach, H., Boston, L.B., Tulpová, Z., et al.** (2021). Long-read sequence

675    assembly: a technical evaluation in barley. The Plant cell **33**:1888-1906.

676    10.1093/plcell/koab077.

677    **Massonnet, M., Cochetel, N., Minio, A., Vondras, A.M., Lin, J., Muyle, A.,**

678    **Garcia, J.F., Zhou, Y., Delledonne, M., Riaz, S., et al.** (2020). The genetic basis of

679    sex determination in grapes. Nature communications **11**:2902. 10.1038/s41467-020-

680    16700-z.

681    **Magris, G., Di Gaspero, G., Marroni, F., Zenoni, S., Tornielli, G.B., Celii, M., De**

682    **Paoli, E., Pezzotti, M., Conte, F., Paci, P., et al.** (2019). Genetic, epigenetic and

683    genomic effects on variation of gene expression among grape varieties. The Plant

684    journal : for cell and molecular biology 99:895-909. 10.1111/tpj.14370.

685    **McKinley, K.L., and Cheeseman, I.M.** (2016). The molecular basis for centromere

686    identity and function. Nature reviews. Molecular cell biology **17**:16-29.

687    10.1038/nrm.2015.5.

688    **Melters, D.P., Bradnam, K.R., Young, H.A., Telis, N., May, M.R., Ruby, J.G.,**

689    **Sebra, R., Peluso, P., Eid, J., Rank, D., et al.** (2013). Comparative analysis of

690    tandem repeats from hundreds of species reveals unique insights into centromere

691    evolution. Genome biology **14**:R10. 10.1186/gb-2013-14-1-r10.

692    **Miga, K.H., and Sullivan, B.A.** (2021). Expanding studies of chromosome structure

693    and function in the era of T2T genomics. Human molecular genetics **30**:R198-r205.

694    10.1093/hmg/ddab214.

695    **Miga, K.H., Koren, S., Rhie, A., Vollger, M.R., Gershman, A., Bzikadze, A.,**

696    **Brooks, S., Howe, E., Porubsky, D., Logsdon, G.A., et al.** (2020). Telomere-to-

697    telomere assembly of a complete human X chromosome. Nature **585**:79-84.

698    10.1038/s41586-020-2547-7.

699    **Minh, B.Q., Schmidt, H.A., Chernomor, O., Schrempf, D., Woodhams, M.D., von**

700    **Haeseler, A., and Lanfear, R.** (2020). IQ-TREE 2: New Models and Efficient

701    Methods for Phylogenetic Inference in the Genomic Era. Molecular biology and

702    evolution **37**:1530-1534. 10.1093/molbev/msaa015.

703    **Minio, A., Massonnet, M., Figueroa-Balderas, R., Castro, A., and Cantu, D.**

704    (2019a). Diploid Genome Assembly of the Wine Grape Carménère. G3 (Bethesda,

705    Md.) **9**:1331-1337. 10.1534/g3.119.400030.

706    **Minio, A., Cochetel, N., Vondras, A.M., Massonnet, M., and Cantu, D.** (2022).

707    Assembly of complete diploid-phased chromosomes from draft genome sequences.

708    G3 (Bethesda, Md.) **12**10.1093/g3journal/jkac143.

709    **Minio, A., Massonnet, M., Figueroa-Balderas, R., Vondras, A.M., Blanco-Ulate,**

710    **B., and Cantu, D.** (2019b). Iso-Seq Allows Genome-Independent Transcriptome

711    Profiling of Grape Berry Development. G3 (Bethesda, Md.) **9**:755-767.

712    10.1534/g3.118.201008.

713    **Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A.,**

714    **Sonnhammer, E.L.L., Tosatto, S.C.E., Paladin, L., Raj, S., Richardson, L.J., et al.**

715   (2021). Pfam: The protein families database in 2021. Nucleic acids research **49**:D412-

716   d419. 10.1093/nar/gkaa913.

717   **Morales-Cruz, A., Aguirre-Liguori, J. Massonnet, M.,Minio. A., Zaccheo, M.,**

718   **Cochetel, Noe., Walker, A., Riaz, S., Zhou, Y.F., Cantu, D., Gaut, B.S.** (2022).

719   Multigenic resistance to Xylella fastidiosa in wild grapes (Vitis sps.) and its

720   implications within a changing climate. bioRxiv, doi:

721   https://doi.org/10.1101/2022.10.08.511428

722   **Naish, M., Alonge, M., Wlodzimierz, P., Tock, A.J., Abramson, B.W., Schmücker,**

723   **A., Mandáková, T., Jamge, B., Lambing, C., Kuo, P., et al.** (2021). The genetic and

724   epigenetic landscape of the Arabidopsis centromeres. Science (New York, N.Y.)

725   **374**:eabi7489. 10.1126/science.abi7489.

726   **Navarro-Payá, D., Santiago, A., Orduña, L., Zhang, C., Amato, A., D'Inca, E.,**

727   **Fattorini, C., Pezzotti, M., Tornielli, G.B., Zenoni, S., et al.** (2021). The Grape

728   Gene Reference Catalogue as a Standard Resource for Gene Selection and Genetic

729   Improvement. Frontiers in plant science **12**:803977. 10.3389/fpls.2021.803977.

730   **Nurk, S., Koren, S., Rhie, A., Rautiainen, M., Bzikadze, A.V., Mikheenko, A.,**

731   **Vollger, M.R., Altemose, N., Uralsky, L., Gershman, A., et al.** (2022). The complete

732   sequence of a human genome. Science (New York, N.Y.) **376**:44-53.

733   10.1126/science.abj6987.

734   **Pertea, M., Kim, D., Pertea, G.M., Leek, J.T., and Salzberg, S.L.** (2016).

735   Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie

736   and Ballgown. Nature protocols **11**:1650-1667. 10.1038/nprot.2016.095.

737   **Pertea, M., Pertea, G.M., Antonescu, C.M., Chang, T.C., Mendell, J.T., and**

738   **Salzberg, S.L.** (2015). StringTie enables improved reconstruction of a transcriptome

739   from RNA-seq reads. Nature biotechnology **33**:290-295. 10.1038/nbt.3122.

740   **Podlevsky, J.D., and Chen, J.J.** (2016). Evolutionary perspectives of telomerase

741   RNA structure and function. RNA biology **13**:720-732.

742   10.1080/15476286.2016.1205768.

743   **Ramu, P., Esuma, W., Kawuki, R., Rabbi, I.Y., Egesi, C., Bredeson, J.V., Bart,**

744 **R.S., Verma, J., Buckler, E.S., and Lu, F.** (2017). Cassava haplotype map highlights

745 fixation of deleterious mutations during clonal propagation. Nature genetics **49**:959-

746 963. 10.1038/ng.3845.

747 **Ranallo-Benavidez, T.R., Jaron, K.S., and Schatz, M.C.** (2020). GenomeScope 2.0

748 and Smudgeplot for reference-free profiling of polyploid genomes. Nature

749 communications **11**:1432. 10.1038/s41467-020-14998-3.

750 **Riaz, S., Tenscher, A.C., Rubin, J., Graziani, R., Pao, S.S., and Walker, M.A.**

751 (2008). Fine-scale genetic mapping of two Pierce's disease resistance loci and a major

752 segregation distortion region on chromosome 14 of grape. TAG. Theoretical and

753 applied genetics. Theoretische und angewandte Genetik **117**:671-681.

754 10.1007/s00122-008-0802-7.

755 **Rice, E.S., and Green, R.E.** (2019). New Approaches for Genome Assembly and

756 Scaffolding. Annual review of animal biosciences **7**:17-40. 10.1146/annurev-animal-

757 020518-115344.

758 **Roach, M.J., Johnson, D.L., Bohlmann, J., van Vuuren, H.J.J., Jones, S.J.M.,**

759 **Pretorius, I.S., Schmidt, S.A., and Borneman, A.R.** (2018). Population sequencing

760 reveals clonal diversity and ancestral inbreeding in the grapevine cultivar

761 Chardonnay. PLoS genetics **14**:e1007807. 10.1371/journal.pgen.1007807.

762 **Rudd, M.K., Wray, G.A., and Willard, H.F.** (2006). The evolutionary dynamics of

763 alpha-satellite. Genome research **16**:88-96. 10.1101/gr.3810906.

764 **Sasaki, M., and Kobayashi, T.** (2021). Gel Electrophoresis Analysis of rDNA

765 Instability in Saccharomyces cerevisiae. Methods in molecular biology (Clifton, N.J.)

766 **2153**:403-425. 10.1007/978-1-0716-0644-5_28.

767 **Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., and Zdobnov,**

768 **E.M.** (2015). BUSCO: assessing genome assembly and annotation completeness with

769 single-copy orthologs. Bioinformatics (Oxford, England) **31**:3210-3212.

770 10.1093/bioinformatics/btv351.

771 **Song, J., Logeswaran, D., Castillo-González, C., Li, Y., Bose, S., Aklilu, B.B., Ma,**

772 **Z., Polkhovskiy, A., Chen, J.J., and Shippen, D.E.** (2019). The conserved structure

773    of plant telomerase RNA provides the missing link for an evolutionary pathway from

774    ciliates to humans. Proceedings of the National Academy of Sciences of the United

775    States of America **116**:24542-24550. 10.1073/pnas.1915312116.

776    **Song, J.M., Xie, W.Z., Wang, S., Guo, Y.X., Koo, D.H., Kudrna, D., Gong, C.,**

777    **Huang, Y., Feng, J.W., Zhang, W., et al.** (2021). Two gap-free reference genomes

778    and a global view of the centromere architecture in rice. Molecular plant **14**:1757-

779    1767. 10.1016/j.molp.2021.06.018.

780    **Stanke, M., Keller, O., Gunduz, I., Hayes, A., Waack, S., and Morgenstern, B.**

781    (2006). AUGUSTUS: ab initio prediction of alternative transcripts. Nucleic acids

782    research **34**:W435-439. 10.1093/nar/gkl200.

783    **Steiner, F.A., and Henikoff, S.** (2014). Holocentromeres are dispersed point

784    centromeres localized at transcription factor hotspots. eLife **3**:e02025.

785    10.7554/eLife.02025.

786    **Talbert, P.B., and Henikoff, S.** (2020). What makes a centromere? Experimental cell

787    research **389**:111895. 10.1016/j.yexcr.2020.111895.

788    **Thorvaldsdóttir, H., Robinson, J.T., and Mesirov, J.P.** (2013). Integrative

789    Genomics Viewer (IGV): high-performance genomics data visualization and

790    exploration. Briefings in bioinformatics **14**:178-192. 10.1093/bib/bbs017.

791    **Turner, K.J., Vasu, V., and Griffin, D.K.** (2019). Telomere Biology and Human

792    Phenotype. Cells **8**10.3390/cells8010073.

793    **Venter, J.C., Adams, M.D., Myers, E.W., Li, P.W., Mural, R.J., Sutton, G.G.,**

794    **Smith, H.O., Yandell, M., Evans, C.A., Holt, R.A., et al.** (2001). The sequence of

795    the human genome. Science (New York, N.Y.) **291**:1304-1351.

796    10.1126/science.1058040.

797    **Vondras, A.M., Lerno, L., Massonnet, M., Minio, A., Rowhani, A., Liang, D.,**

798    **Garcia, J., Quiroz, D., Figueroa-Balderas, R., Golino, D.A., et al.** (2021).

799    Rootstock influences the effect of grapevine leafroll-associated viruses on berry

800    development and metabolism via abscisic acid signalling. Molecular plant pathology

801    **22**:984-1005. 10.1111/mpp.13077.

802 **Wang, B., Yang, X., Jia, Y., Xu, Y., Jia, P., Dang, N., Wang, S., Xu, T., Zhao, X.,**
803 **Gao, S., et al.** (2022a). High-quality Arabidopsis thaliana Genome Assembly with
804 Nanopore and HiFi Long Reads. Genomics, proteomics & bioinformatics **20**:4-13.
805 10.1016/j.gpb.2021.08.003.

806 **Wang, N., Song, X., Ye, J., Zhang, S., Cao, Z., Zhu, C., Hu, J., Zhou, Y., Huang,**
807 **Y., Cao, S., et al.** (2022b). Structural variation and parallel evolution of apomixis in
808 citrus during domestication and diversification. National Science Review
809 10.1093/nsr/nwac114.

810 **Xu, Y., Wu, Y., Wang, L., Qian, C., Wang, Q., and Wan, W.** (2020). Identification
811 of curcumin as a novel natural inhibitor of rDNA transcription. Cell cycle
812 (Georgetown, Tex.) **19**:3362-3374. 10.1080/15384101.2020.1843817.

813 **Yan, H., Bombarely, A., and Li, S.** (2020). DeepTE: a computational method for de
814 novo classification of transposons with convolutional neural network. Bioinformatics
815 (Oxford, England) **36**:4269-4275. 10.1093/bioinformatics/btaa519.

816 **Yuan, X., Dai, M., and Xu, D.** (2020). Telomere-related Markers for Cancer. Current
817 topics in medicinal chemistry **20**:410-432. 10.2174/1568026620666200106145340.

818 **Yue, J., Chen, Q., Wang, Y., Zhang, L., Ye, C., Wang, X., Cao, S., Lin, Y., Huang,**
819 **W., Xian, H., et al.** (2022). Telomere-to-telomere and gap-free reference genome
820 assembly of the kiwifruit Actinidia chinensis. Horticulture Research uhac264.
821 10.1093/hr/uhac264.

822 **Zhang, C., Yang, Z., Tang, D., Zhu, Y., Wang, P., Li, D., Zhu, G., Xiong, X.,**
823 **Shang, Y., Li, C., et al.** (2021). Genome design of hybrid potato. Cell **184**:3873-
824 3883.e3812. 10.1016/j.cell.2021.06.006.

825 **Zhang, Y., Fu, J., Wang, K., Han, X., Yan, T., Su, Y., Li, Y., Lin, Z., Qin, P., Fu,**
826 **C., et al.** (2022). The telomere-to-telomere gap-free genome of four rice parents
827 reveals SV and PAV patterns in hybrid rice breeding. Plant biotechnology journal
828 **20**:1642-1644. 10.1111/pbi.13880.

829 **Zhou, Y., Massonnet, M., Sanjak, J.S., Cantu, D., and Gaut, B.S.** (2017).
830 Evolutionary genomics of grape (Vitis vinifera ssp. vinifera) domestication.

831 Proceedings of the National Academy of Sciences of the United States of America

832 **114**:11715-11720. 10.1073/pnas.1709257114.

833 **Zhou, Y., Minio, A., Massonnet, M., Solares, E., Lv, Y., Beridze, T., Cantu, D.,**

834 **and Gaut, B.S.** (2019). The population genetics of structural variants in grapevine

835 domestication. Nature plants **5**:965-979. 10.1038/s41477-019-0507-8.

836 **Gil, P., Dewey, E., Friml, J., Zhao, Y., Snowden, K.C., Putterill, J., Palme, K.,**

837 **Estelle, M., and Chory, J.** (2001). BIG: a calossin-like protein required for polar

838 auxin transport in Arabidopsis. Genes & development **15**:1985-1997.

839 10.1101/gad.905201.

840 **Saibo, N.J., Vriezen, W.H., De Grauwe, L., Azmi, A., Prinsen, E., and Van der**

841 **Straeten, D.** (2007). A comparative analysis of the Arabidopsis mutant amp1-1 and a

842 novel weak amp1 allele reveals new functions of the AMP1 protein. Planta **225**:831-

843 842. 10.1007/s00425-006-0395-9.

844 **Shi, H., Ye, T., Wang, Y., and Chan, Z.** (2013). Arabidopsis ALTERED MERISTEM

845 PROGRAM 1 negatively modulates plant responses to abscisic acid and dehydration

846 stress. Plant physiology and biochemistry : PPB **67**:209-216.

847 10.1016/j.plaphy.2013.03.016.

848 **Zubimendi, J.P., Martinatto, A., Valacco, M.P., Moreno, S., Andreo, C.S.,**

849 **Drincovich, M.F., and Tronconi, M.A.** (2018). The complex allosteric and redox

850 regulation of the fumarate hydratase and malate dehydratase reactions of Arabidopsis

851 thaliana Fumarase 1 and 2 gives clues for understanding the massive accumulation of

852 fumarate. The FEBS journal **285**:2205-2224. 10.1111/febs.14483.

853 **Hölscher, C., Lutterbey, M.C., Lansing, H., Meyer, T., Fischer, K., and von**

854 **Schaewen, A.** (2016). Defects in Peroxisomal 6-Phosphogluconate Dehydrogenase

855 Isoform PGD2 Prevent Gametophytic Interaction in Arabidopsis thaliana. Plant

856 physiology **171**:192-205. 10.1104/pp.15.01301.

857 **Zou, C., Massonnet, M., Minio, A., Patel, S., Llaca, V., Karn, A., Gouker, F.,**

858 **Cadle-Davidson, L., Reisch, B., Fennell, A., et al.** (2021). Multiple independent

859 recombinations led to hermaphroditism in grapevine. Proceedings of the National

860    Academy of Sciences of the United States of America 11810.1073/pnas.2023548118.

861    **Hua Xiao, Zongjie Liu, Nan Wang, Shuo Cao , Guizhou Huang , Wenwen Liu ,**

862    **Yanling Peng , Qiming Long , Summaira Riaz , Andrew M. Walker , Brandon S.**

863    **Gaut, and Yongfeng Zhou.** (2023). The adaptive and maladaptive introgression in

864    grapevine domestication. PNAS, in press

865

## 866    Acknowledgments

## 878    Author contributions

879    Y.Z. conceived and designed the project with H.X., Z.C. and C.R.. Z.C. provided the

880    PN40024 sample. X.S. W.L., X.X and Z.M. performed the tissue culture of the

881    sample in the greenhouse. X.S., X.W., H.X., N.W., F.Z., H.X. and Y.W. performed

882    the bioinformatic analyses. A.V., K.A., D.H., J. G., J.T., D.W. Z.L., X.L. and W.L.

883    performed the gene annotation. Y.P., S.H. Z.L., W.L., X.W., Y.F., Y.W and C.L.

884    assisted in bioinformatics analyses. X.S., S.C., X.W., H.X. and Y.Z. wrote the

885    manuscript with comments and inputs from all authors.

## 886    Conflict of interests
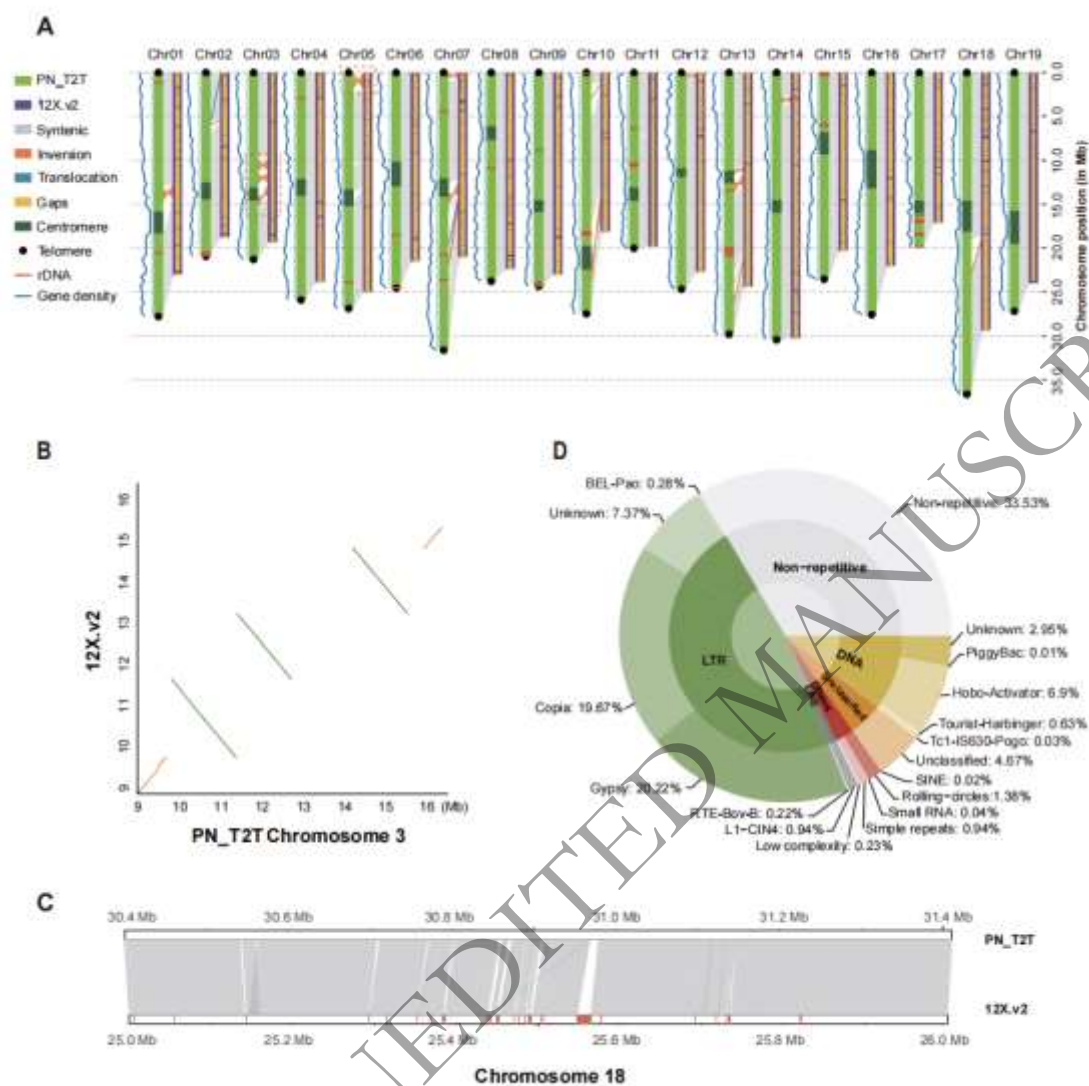
887    The authors declare no conflict of interest.

888    **Table 1. Comparison of genomic features of 12X.v2 and PN_T2T assembly**

|  | 12X.v2 | PN_T2T |
|---|---|---|
| Total sequence length (bp) | 426,176,009 | 494,873,210 |
| Number of chromosomes | 19 | 19 |
| Contig N50 (bp) | 102,700 | 26,899,771 |
| Max length | 30,274,277 | 36,684,271 |
| Number of gaps | 9429 | 0 |
| Centromere | - | 19/19 |
| Telomere | - | 36/38 |
| Bases masked (bp) | 303,719,475 | 328,929,883 |
| Retroelements (bp) | 217,819,122 | 241,027,616 |
| LTR (bp) | 212,117,752 | 235,245,099 |
| The number of genes | 28,516 | 37,534 |
| The number of TE | 942,096 | 935,783 |
| BUSCO | 93.70% | 98.50% |

889

890

891

892

**Figure 1. The T2T gap-free assembly of the grapevine reference genome.** (A) An overview of the genome assemblies (12X.v2 right, PN_T2T left). The red dashed boxes on Chromosome 3 and Chromosome 5 indicated differences in large inversions between the two versions of genomic assemblies. (B) A zoomed-in portion of the red dashed box region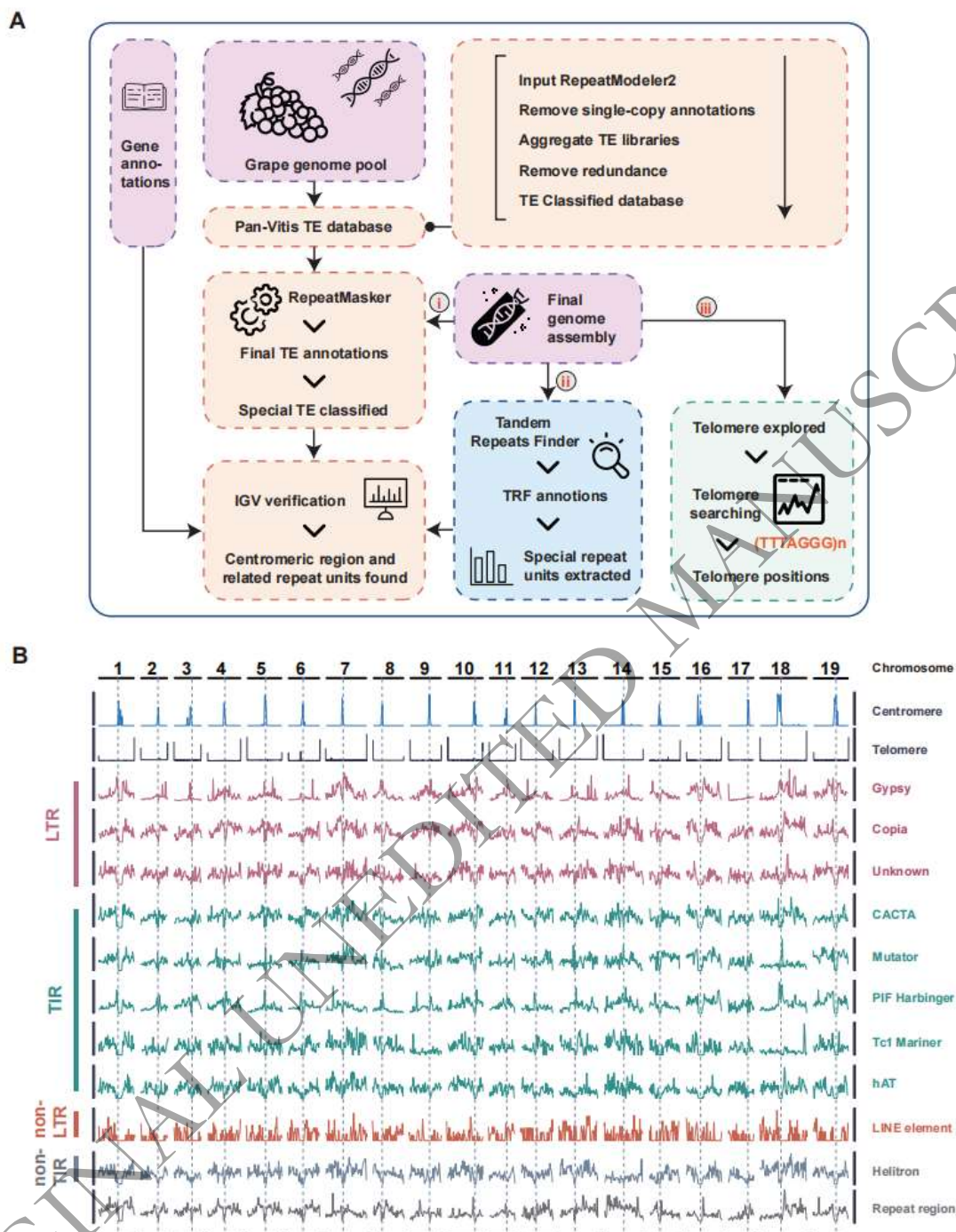 on Chromosome 3 in A. (C) Plot showing 1 Mb syntenic region between the 12X.v2 and PN_T2T assemblies on Chromosome 18. Grey bands connected corresponding collinear regions, and red boxes at the bottom showed the gaps in 12X.v2. (D) Types and percentages of different TE families detected in the PN_T2T genome.

902

**Figure 2. The repeats annotation in PN_T2T reference genome.** (A) Dataflow of centromere and telomere predictions. (B) Chromosomal distribution of telomeres, centromeres and different types of TEs. Dotted vertical lines indicated the center locations of predicted centromeres.

910

911

912  **Figure 3. The schematic illustration of centromeric repeat units in the PN_T2T**

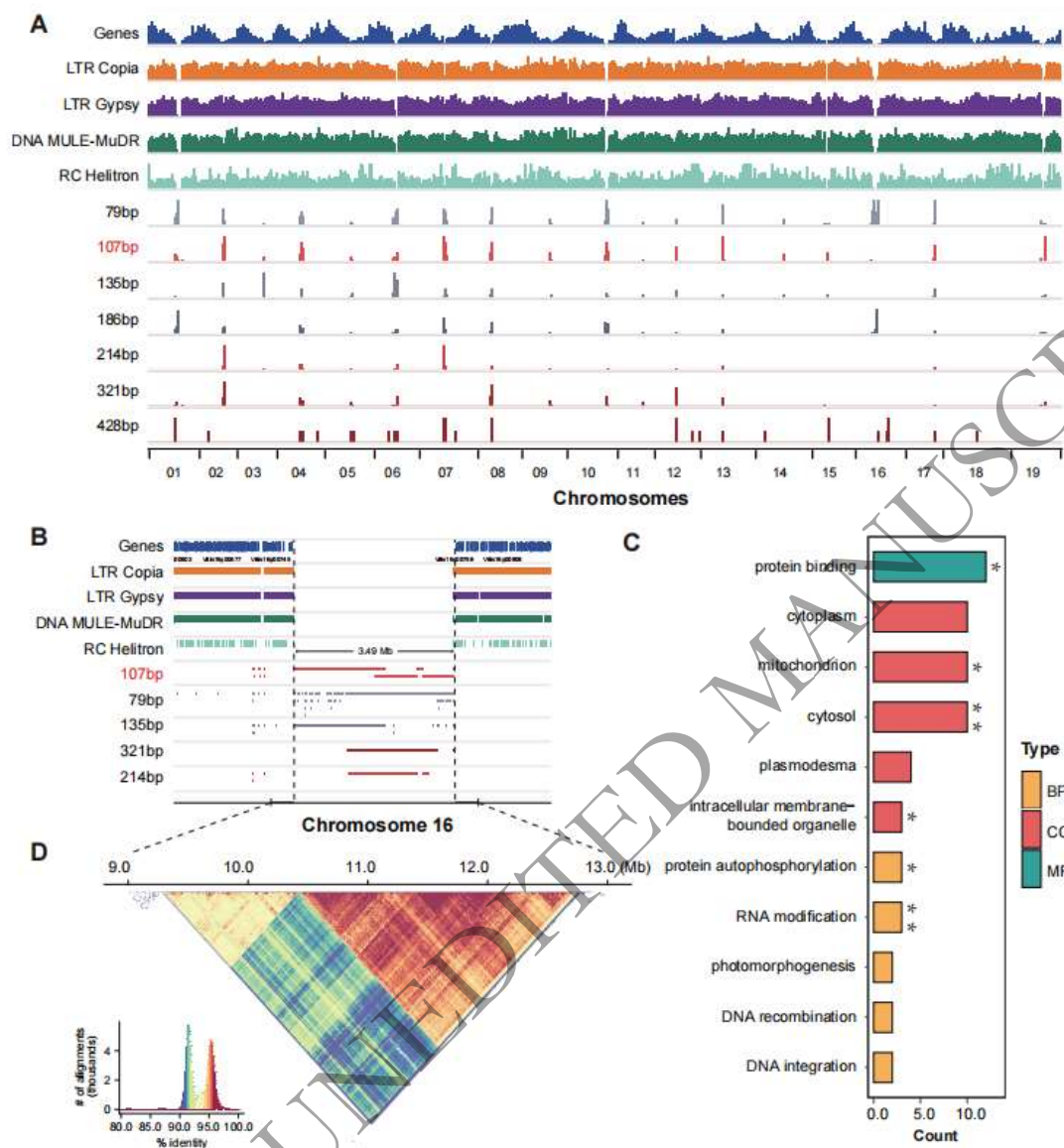913  **genome.** (A) The distribution of different repeat units' lengths in the whole genome.

914  The number of different repeat unit copies was indicated in the upper part of the

915  graphs while the chromosomal percentage of different repeat units was shown in the

916  lower part. (B) The total length of 107 bp repeat unit copies in each chromosome. (C)

917  The alignment of the 107 bp repeat units among 19 chromosomes. (D-H) The total

918  length of different repeat units in chromosomes 16, 17, 3, 14 and 18, respectively.

919

920

**Figure 4. Characteristics and distribution of repeat unit copies in centromeres.**
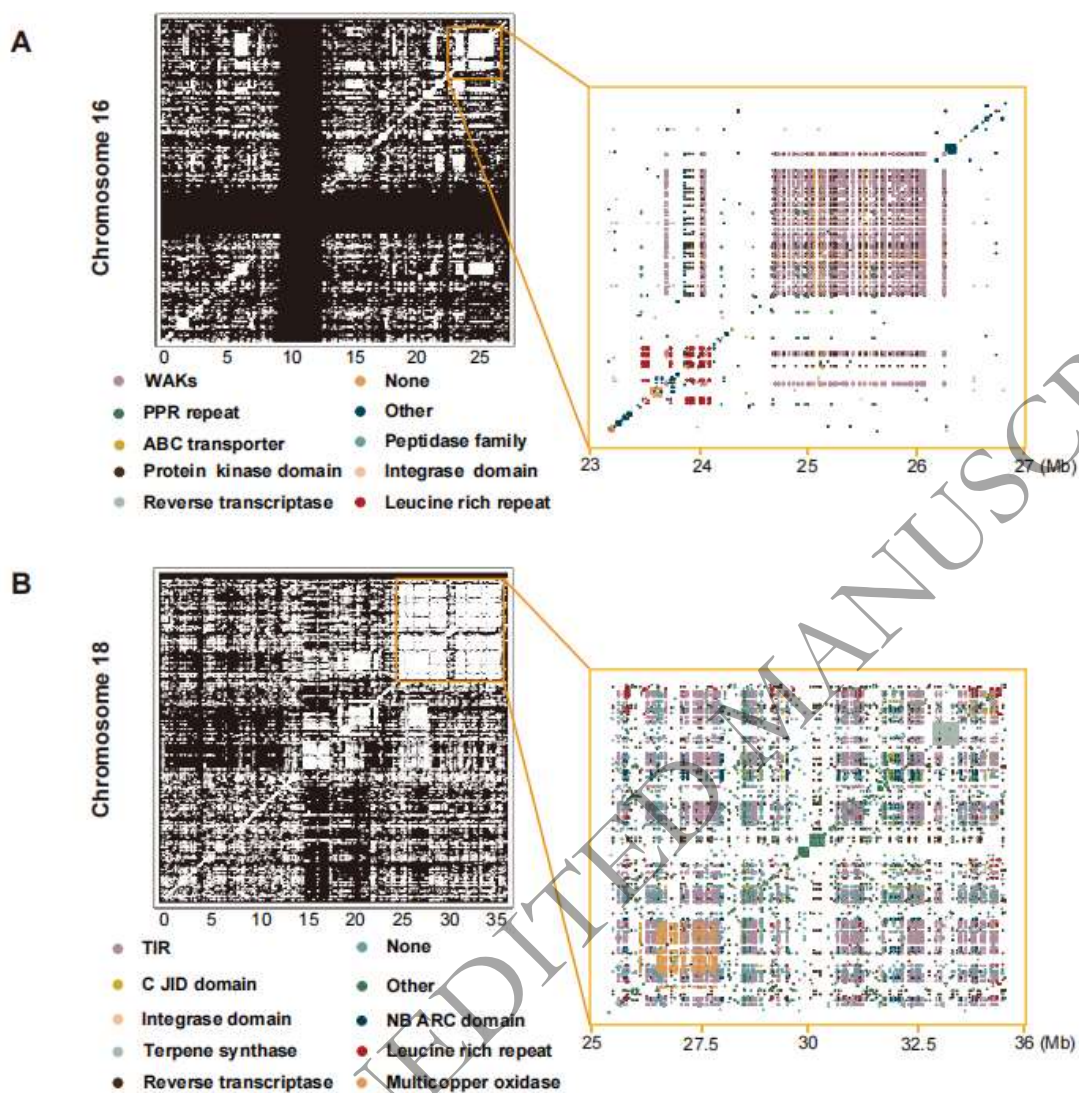(A) The distribution of genes, TEs and different repeat units in the whole genome. (B)
Visualization of the predicted centromeric region on Chromosome 16 in IGV. (C) GO
functional annotation of genes captured in centromeres. MF: molecular function, CC:
cellular component, BP: biological process. Enrichment significant p-value: *, P<
0.05. **, P<0.001. (D) The triangle shows sequence similarity within each haplotype
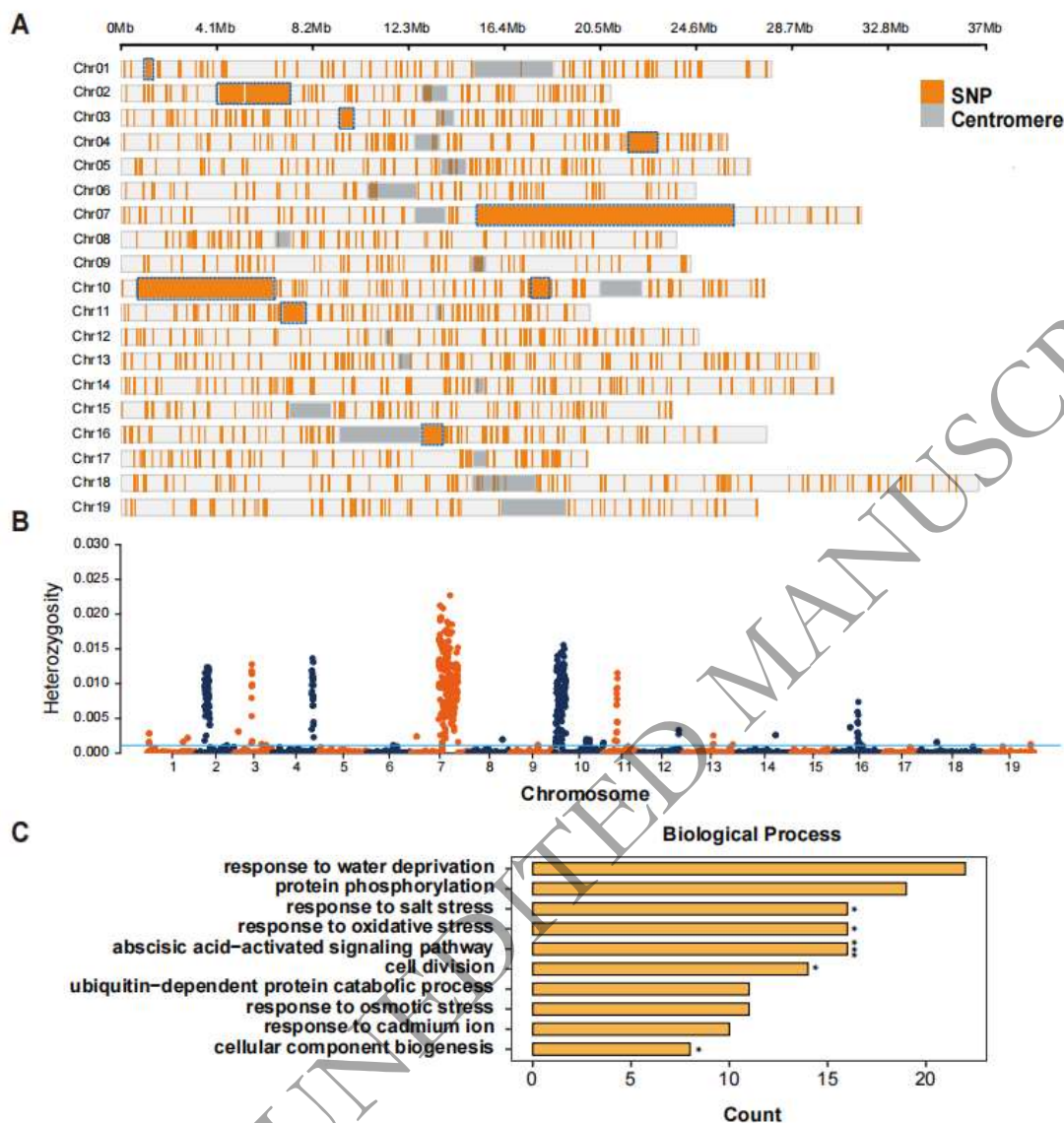and colored by identity.

928

929

**Figure 5. Schematic of identified gene clusters. (A)** The gene clusters in Chromosome 16 and Chromosome 16: 22-27 Mb. **(B)** The gene clusters in Chromosome 18 and Chromosome 18:25-36 Mb. The graphs on the right were the enlargement of regions in white boxes on the left. Different color indicated the different gene clusters. Both split and compound.

935

**Figure 6. The characterization of heterozygous regions in PN40024.** (A) The heterozygous sites were shared in all four PN40024 samples. The Grey bar indicated the centromere region while the orange lines indicated the heterozygous sites that existed in all samples. Blue boxes picked out the large heterozygous fragments. (B) The heterozygosity in PN40024 genome calculated with no overlapping 100 kb windows across four samples. (C) The GO enrichment analysis of genes contained heterozygous sites shown in A. Enrichment significant p-value: *, $P < 0.05$. **, $P < 0.001$. ***, $P < 0.001$.