

Manuscript Number:	GIGA-D-20-00195R2	
Full Title:	The chromosome-level draft genome of a diploid plum (<i>Prunus salicina</i>)	
Article Type:	Data Note	
Funding Information:	The Industry University Research Collaborative Innovation Major Projects of Guangzhou Science Technology Innovation Commission (201704020021)	Dr Yehua He
	Guangdong Key Laboratory of Innovation Method and Decision Management System (CN) (2016LM1128)	Dr Yehua He
Abstract:	<p>Background: Plums are one of the most economically important Rosaceae fruit crops, and contain dozens of species distributed across the world. Until now, only limited genomic information is available for the genetic studies and breeding programs of plums. <i>Prunus salicina</i>, an important diploid plum species, plays a predominant role in modern commercial plums production. Here we selected <i>P. salicina</i> for whole-genome sequencing and presented a chromosome-level genome assembly through the combination of PacBio sequencing, Illumina sequencing and Hi-C technology.</p> <p>Findings: The assembly had a total size of 284.2 Mb, with contig N50 of 1.8Mb and Scaffold N50 of 32.3 Mb. 96.56% of the assembled sequences were anchored onto eight pseudochromosomes and a total of 24,448 protein-coding genes were identified. Phylogenetic analysis showed that <i>P. salicina</i> had closer relationship with <i>P. mume</i> and <i>P. armeniaca</i>, with <i>P. salicina</i> diverging from their common ancestor approximately 9.05 million years ago (Mya). 146 gene families were expanded during <i>P. salicina</i> evolution, and some cell wall-related GO terms were significantly enriched. It was noteworthy that members in the DUF579 family, a new class involved in xylan biosynthesis, were significantly expanded in <i>P. salicina</i>, which provided new insight into the xylan metabolism in plums.</p> <p>Conclusions: We constructed the first high-quality chromosome-level plum genome using PacBio, Illumina and Hi-C technologies. This work provides a valuable resource for facilitating plum breeding programs and studying the genetic diversity mechanisms of plums and <i>Prunus</i> species.</p>	
Corresponding Author:	Yehua He Soth China Agricultural University Guangzhou, Guangdong Province CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Soth China Agricultural University	
Corresponding Author's Secondary Institution:		
First Author:	Chaoyang Liu	
First Author Secondary Information:		
Order of Authors:	Chaoyang Liu	
	Chao Feng	
	Weizhuo Peng	
	Jingjing Hao	
	Juntao Wang	
	Jianjun Pan	

	Yehua He
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Dear Editor and Reviewer,</p> <p>We would like to thank you for helpful suggestions on our manuscript entitled "The chromosome-level draft genome of a diploid plum (<i>Prunus salicina</i>)" (GIGA-D-20-00195R1). Following the comments and suggestions, we corrected the numbering of the <i>P. salicina</i> pseudochromosomes all over the manuscript. The <i>P. salicina</i> pseudo-chromosome names and gene IDs in Figure 1, 3 and 4 were modified in the revised manuscript. We also contacted the staff of GDR (Genome Database for Rosaceae) and GenBank to update our genomic information. We expected that it would meet the publication requirement of GigaScience.</p> <p>A point by point response to the reviewer' comments and questions and the main corrections in the paper were provided below.</p> <p>Reviewer #1: I thank the authors for further revising their manuscript and clarifying some outstanding issues in regards of English proofreading and MS layout. Thank you very much for the answers to my previous questions, even if I do not fully agree with soem of them.</p> <p>However, there is still a major revision necessary before the manuscript is ready for publication. I bet I overlooked it in the first version of the manuscript because of the other issues that were since corrected. My main concern relates to the chromosome nomenclature: the chromosome numbering is not in adequation with the <i>Prunus</i> genetic map. For exemple, Chromosome 1 in all <i>Prunus</i> species is always the largest one and following Figure 1, it appears that it is chromosome 2, here. The same remark applies to the other chromosomes, not only chromosome 2 (see figure 2B, chromosome 1 of <i>P. salicina</i> should in fact be chromosome 6 in the <i>Prunus</i> genetic map, chr3 should be chr4 and so on), and that's the reason why I was recommending using, even a few, <i>Prunus</i> genetic markers, to correct this discrepancy. This major issue is coming from the first release of the <i>P. mume</i> genome in 2012 and was reproduced in the <i>P. armeniaca</i> genome presented here. If colinearity has to be displayed (Figure 3) then it should be made clear that Chr2 here should be in fact Chr 1 in the genetic map. In fact, I would once again recommend the authors to re-order their chromosomes, according to the general acknowledged genetic map. Since the genetic maps were obtained by using molecular markers which are largely colinear and syntenic in between <i>Prunus</i> species (peach, <i>P. mume</i>, apricot and plum included) I would strongly recommend to right this issue, both within the <i>P. salicina</i> assembly and the following colinearity studies with the other genomes. Since genetic maps were released before genome assembly, the authors are expected to follow the internationally acknowledged nomenclature. Reproducing forever the mistake made initially for the <i>P. mume</i> genome would severely limit the interest of this de novo assembled genome and thus the impact of its release.</p> <p>In conclusion, I recommend the authors to correct the numbering of the <i>P. salicina</i> chromosome all over the MS (by using a few of plum markers and even better <i>Prunus</i> orthologous markers as published in https://doi.org/10.1371/journal.pone.0208032, for that they only need to do a ePCR with markers depicted in Table S2F) and the data available online (and therefore Figure 3, accordingly).</p> <p>Response : Thank you very much for your kindly suggestion. We further carefully read the literatures that you mentioned, and realized that the important role of the <i>Prunus</i> genetic maps was ignored in our first version of the manuscript. According to your kind advices, the numbering of the <i>P. salicina</i> pseudochromosome was corrected in the revised manuscript, using the markers that you recommended.</p> <p>The present pseudochromosome numbering is consisted with that in the published <i>P. salicina</i> genetic map (Ref 77 in the revised manuscript). The <i>P. salicina</i> pseudo-chromosomes names and gene IDs in Figure 1, 3 and 4 were modified in the revised manuscript. The detailed information about the <i>P. salicina</i> pseudochromosomes in Table S5 was corrected accordingly. Moreover, we also contacted the staff of GDR (Genome Database for Rosaceae) and GenBank to update our genomic information.</p>

Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p>	Yes

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

The chromosome-level draft genome of a diploid plum (*Prunus salicina*)

Chaoyang Liu^{1,3*}, Chao Feng^{2,*}, Weizhuo Peng^{1,3}, Jingjing Hao^{1,3}, Juntao Wang^{1,3}
Jianjun Pan⁴, Yehua He^{1,3}

¹ Key Laboratory of Biology and Germplasm Enhancement of Horticultural Crops in South China, Ministry of Agriculture, South China Agricultural University, Guangzhou 510642, China

² Key Laboratory of Plant Resources Conservation and Sustainable Utilization, South China Botanical Garden, Chinese Academy of Sciences, Guangzhou 510650, China

³ Maoming Branch, Guangdong Laboratory for Lingnan Modern Agriculture, Maoming 525000, China

⁴ Agricultural Technology Extension Center of Conghua District, Guangzhou. Guangzhou 510900, Guangdong Province, China

*Equal contribution

Corresponding author: Yehua He (email: heyehua@hotmail.com)

30 **Abstract**

31 **Background:** Plums are one of the most economically important Rosaceae fruit crops,
32 and contain dozens of species distributed across the world. Until now, only limited
33 genomic information is available for the genetic studies and breeding programs of
34 plums. *Prunus salicina*, an important diploid plum species, plays a predominant role
35 in modern commercial plums production. Here we selected *P. salicina* for
36 whole-genome sequencing and presented a chromosome-level genome assembly
37 through the combination of PacBio sequencing, Illumina sequencing and Hi-C
38 technology. **Findings:** The assembly had a total size of 284.2 Mb, with contig N50 of
39 1.8Mb and scaffold N50 of 32.3Mb. 96.56% of the assembled sequences were
40 anchored onto eight pseudochromosomes and a total of 24,448 protein-coding genes
41 were identified. Phylogenetic analysis showed that *P. salicina* had closer relationship
42 with *P. mume* and *P. armeniaca*, with *P. salicina* diverging from their common
43 ancestor approximately 9.05 million years ago (Mya). 146 gene families were
44 expanded during *P. salicina* evolution, and some cell wall-related GO terms were
45 significantly enriched. It was noteworthy that members in the DUF579 family, a new
46 class involved in xylan biosynthesis, were significantly expanded in *P. salicina*, which
47 provided new insight into the xylan metabolism in plums. **Conclusions:** We
48 constructed the first high-quality chromosome-level plum genome using PacBio,
49 Illumina and Hi-C technologies. This work provides a valuable resource for
50 facilitating plum breeding programs and studying the genetic diversity mechanisms of
51 plums and *Prunus* species.

52

53 **Background**

54 Plums are one of the most economically important Rosaceae fruit crops and are
55 produced throughout the world. About 12.6 million tons of plums (include sloes) are
56 produced per year (FAOSTAT 2018, <http://faostat.fao.org/>), and the fruits are widely
57 used for fresh consumption and processing like canning and beverages [1]. There are
58 19-40 species of plums distributed across Asia, Europe and America. Plums have
59 great diversity and are considered as a link between the major subgenera in the genus
60 *Prunus* [2].

61 *Prunus salicina*, commonly called the Japanese plum or Chinese plum, is an
62 important diploid ($2x=2n=16$) plum species that predominates in the modern
63 commercial production of plums (Fig. 1). *P. salicina* originates in China and its fruits
64 are mostly used for fresh consumption for their characteristic taste [3]. Cultivars of *P.*
65 *salicina* have wide variability in phenology, fruit size and shape, flavour, firmness,
66 aroma, texture, phenolic composition, antioxidant activity and both skin and pulp
67 color [4].

68 However, the genetic and genomic information for *P. salicina* as well as most plum
69 species were scarce [5]. The availability of a fully sequenced and annotated genome
70 will help to measure and characterize the genetic diversity and determine how this
71 diversity relates to the tremendous phenotypic diversity among plum cultivars. The
72 genomic information is essential to support many of the studies involved in
73 fundamental questions about plums biology and genetics. Moreover, genome-based
74 tools could be developed to improve breeding works of plums, which were usually
75 hindered by the high degree of heterozygosity, self-incompatible and long juvenile
76 stage [2, 5, 6].

77 The fruit firmness, one of the most important indices of plum quality, is closely
78 associated with cell wall compositions [2]. Xylan is a major component of secondary
79 cell walls [7], and the xylan metabolism is involved in various aspects of plant growth
80 and development like fruit ripening and softening [8]. According to previous studies,
81 the plum species presented more xylose (the main component of xylan) compared to

82 other *Prunus* species, and plums were regarded as one of the richest natural sources of
83 xyliot [9, 10]. The relatively high levels of xylan-related metabolites may be
84 associated with the distinct mechanisms of the xylan metabolism in plums, and the
85 available plum genomic information will be helpful to better understand the
86 mechanism at molecular level.

87 Genome resources are already available for a number of Rosaceae fruit crops [11],
88 including apple [12-15], peach [16], pear [17-20], strawberry [21, 22], almond [23,
89 24], black raspberry [25], sweet cherry [26, 27], apricot [28, 29], loquat [30] and
90 *Prunus mume* [31]. However, whole-genome sequencing and chromosome-level
91 assembly for plums have not been reported until now. In this study, *P. salicina* was
92 selected for the whole-genome sequencing as a genomic reference. A high-quality
93 chromosome-level *de novo* genome assembly of *P. salicina* was generated using an
94 integrated strategy that combines PacBio sequencing, Illumina sequencing and Hi-C
95 technology. The assembly has a total size of 284.2 Mb with contig N50 of 1.8Mb and
96 scaffold N50 of 32.3 Mb, and vast majority (96.56%) of the assembled sequence was
97 anchored onto eight pseudochromosomes. The availability of the high-quality
98 chromosome-scale genome sequences not only provides fundamental knowledge
99 regarding plum biology but also presents a valuable resource for genetic diversity
100 analysis and breeding programs of plums and other *Prunus* crops.

101

102 **Methods**

103 **Sample collection**

104 The *Prunus salicina* Lindl. cv. ‘Sanyueli’, a Japanese plum landrace originating from
105 Southern China, was selected for genome sequencing and assembly. ‘Sanyueli’ has a
106 cultivation history of more than 200 years and many distinctive characteristics,
107 including early-maturation, high-yield and low chilling requirements. The samples of
108 the ‘Sanyueli’ were kept at the Horticultural Germplasm Conservation Center of
109 South China Agricultural University (SCAU) for breeding and research in Guangzhou,
110 Guangdong Province, China (113°22'4" N, 23°9'5" E). Total genomic DNA was

111 extracted from fresh young leaves of 5-year-old *P. salicina* tree using the CTAB
112 method [32]. Samples from a total of six tissues, including leaf, flower, branch, young
113 fruit pericarp, young fruit pulp and matured fruit, were collected from the same *P.*
114 *salicina* tree. Total RNA was extracted from the six tissues using E.N.Z.A.[®] Plant
115 RNA kit (OMEGA).

116

117 **Library construction and sequencing**

118 A combination of PacBio single-molecule real-time (SMRT) sequencing, Illumina's
119 paired-end sequencing and Hi-C technology was applied. For PacBio sequencing,
120 SMRT libraries were constructed using the PacBio 20-kb protocol
121 (<https://www.pacb.com/>). The Illumina DNA paired-end libraries were constructed
122 with an insert size of 350 bp, and sequencing was performed on the Illumina HiSeq
123 4000 platform according to the manufacturer's instructions. Reads with adaptors, with
124 more than 10% unknown bases (N) and with more than 50% low-quality bases (≤ 5)
125 were filtered out to gain the clean data for further analysis.

126 The Hi-C library was prepared using the standard procedures. The young leaves of
127 the same *P. salicina* tree were used as starting materials. Nuclear DNA from young
128 leaves was cross-linked in situ, extracted, and digested with DpnII restriction
129 endonuclease. The 5' overhangs of the digested fragments were biotinylated, and the
130 resulting blunt ends were ligated. The cross-links were reversed after ligation,
131 proteins were removed to release the DNA molecules. The purified DNA was sheared
132 to a mean fragment size of 350 bp and ligated to adaptors, followed by purification
133 through biotin-streptavidin-mediated pull down. The quality of Hi-C sequencing was
134 evaluated with HiCUP [33].

135 The RNA-seq libraries for the six tissues of *P. salicina* were constructed according
136 to the manufacturer's protocols, and were sequenced by Illumina HiSeq 4000 in
137 paired-end 150bp mode.

138

139 **Genome size estimation and *de novo* assembly**

140 Sequencing data from the Illumina library were used to perform a k-mer analysis to

141 estimate the genome size of *P. salicina*. Quality-filtered reads were subjected to
142 17-mer frequency distribution analysis using SOAPdenovo (SOAPdenovo, RRID:
143 SCR_010752) [34].

144 The *de novo* assembly of the *P. salicina* genome was carried out using the
145 FALCON assembler (FALCON, RRID: SCR_016089) [35], followed by the polishing
146 with Quiver [36] and Pilon (Pilon, RRID: SCR_014731) [37]. The PacBio subreads
147 were subsequently processed by a self-correction of errors using FALCON [35]
148 according to the manufacturer's instructions with the following parameters:
149 length_cutoff = 7,000, length_cutoff_pr = 4,000, max_diff = 100, max_cov = 100. The
150 draft assembly was further polished using Quiver [36]. The pipeline of 'Purge
151 Haplotigs' was used to remove the redundant sequences caused by genomic
152 heterozygosity [38]. Finally, the Illumina reads were mapped back to the assembly
153 and the remaining errors were corrected by Pilon [37].

154 Clean Hi-C reads were aligned to the assembled genome with BWA aligner (BWA,
155 RRID: SCR_010910) with default parameters [39]. Only uniquely aligned read pairs
156 whose mapping quality more than 20 were remained for further analysis. Invalid read
157 pairs, including dangling-end and self-cycle, relegation, and dumped products, were
158 filtered by HiCUP [33]. The valid interaction pairs were used to cluster, order, and
159 orient the assembly contigs onto pseudochromosomes by LACHESIS (LACHESIS,
160 RRID:SCR_017644; parameters: CLUSTER_N = 8, CLUSTER_MIN_RE_SITES
161 = 1157, CLUSTER_MAX_LINK_DENSITY = 5, CLUSTER
162 _NONINFORMATIVE_RATIO = 0) [40]. The Juicebox [41] was applied to build the
163 interaction matrices and complete the visual correction.

164

165 **Genome quality evaluation**

166 To evaluate the coverage of the assembly, the paired-end Illumina short reads were
167 aligned to the assembly using BWA. RNA-seq reads from six tissues of *P. sacilina*
168 were mapped against our assembly using Hisat with default parameters [42]. The
169 SNPs were counted to evaluate the accuracy of the genome assembly. For CEGMA
170 (Core Eukaryotic Genes Mapping Approach; CEGMA, RRID: SCR_015055)

171 evaluation, a set of highly reliable conserved protein families that occur in a range of
172 model eukaryotes were build and then the 248 core eukaryotic genes were mapped to
173 the genome [43]. Genome completeness was also accessed using BUSCO
174 (Benchmarking Universal Single-Copy Orthologs; RRID: SCR_015008) analysis
175 which included a set of 1440 single-copy orthologous genes [44].

176

177 **Repeat annotations**

178 To annotate repeat elements in the *P. salicina* genome, a combined strategy based on
179 homology searching and *de novo* prediction was applied. For homology-based
180 prediction, interspersed repeats were identified using RepeatMasker
181 (<http://www.repeatmasker.org>) (RepeatMasker, RRID: SCR_012954) and
182 RepeatProteinMask (RepeatProteinMask, RRID: SCR 012954) [45] to search against
183 the Repbase database [46]. For *de novo* prediction, RepeatScout
184 (<http://www.repeatmasker.org/>) (RepeatScout, RRID:SCR 014653) [47],
185 RepeatModeler (<http://www.repeatmasker.org/RepeatModeler/>) RepeatModeler
186 (RRID:SCR_015027), and LTR_Finder (http://tlife.fudan.edu.cn/tlife/ltr_finder/)
187 (LTR_Finder, RRID:SCR_015247) [48] were used to identify *de novo* involved
188 repeats. Tandem repeats were also *de novo* predicted using Tandem Repeats Finder
189 (TRF) [49].

190 Telomere sequences were identified by BLASTN searches of both ends of the
191 pseudochromosomes using four tandem repeats of the telomere repeat motif
192 (TTTAGGG) with e-value cut-off of 0.003.

193

194 **Gene annotations**

195 A combination of three approaches, including homology-based prediction, *de novo*
196 prediction and transcriptome-based prediction, was used to predict the protein-coding
197 genes within *P. salicina* genome. For homology-based prediction, the homologous
198 protein sequences of *Prunus persica*, *Prunus avium*, *Prunus mume*, *Pyrus*
199 *bretschneideri*, *Malus domestica*, *Fragaria vesca* and *Arabidopsis thaliana* were
200 obtained from NCBI database and mapped onto the *P. salicina* genome using TblastN

201 (TBLASTN; RRID:SCR_011822) (E-value $\leq 1e-5$) [50], and then the matching
202 proteins were aligned to the homologous genome sequences for accurate spliced
203 alignments with GeneWise (GeneWise, RRID:SCR_015054) [51] to define gene
204 models. For *de novo* prediction, Augustus (Augustus, RRID: SCR_008417) [52],
205 GlimmerHMM (GlimmerHMM, RRID: SCR_002654) [53], SNAP (SNAP, RRID:
206 SCR_002127) [54], GeneID (GeneID, RRID: SCR_002473) [55] and Genescan
207 (Genescan, RRID: SCR_012902) [56] were used to predict the coding regions of
208 genes. For transcriptome-based predictions, RNA-seq data from six tissues were used
209 for genome annotation, processed by HISAT2 (HISAT2, RRID: SCR_015530) [42]
210 and Stringtie (StringTie, RRID: SCR_016323) [57]. RNA-seq data were also *de novo*
211 assembled with Trinity (Trinity, RRID: SCR_013048) [58]. The assembled sequences
212 were aligned against *P. salicina* genome with PASA (Program to Assemble Spliced
213 Alignment, PASA, RRID: SCR_014656) [59], and the effective alignments were
214 assembled to gene structures. Gene models predicted by all of the methods were
215 integrated by EVIDENCEModeler (EVIDENCEModeler, RRID: SCR_014659) [59]. To
216 update the gene models, PASA was further used to generate UTRs [59].

217

218 **Gene functions**

219 The functional annotation of protein-coding genes within *P. salicina* genome was
220 carried out by aligning protein sequences against SwissProt [60] and NR databases
221 using BLASTp (with a threshold of E-value $\leq 1e-5$). The protein motifs and domains
222 were annotated by searching against InterPro (InterPro, RRID: SCR_006695) [61] and
223 Pfam (Pfam, RRID: SCR_004726) database [62] with InterProScan (InterProScan,
224 RRID: SCR_005829) [63]. Gene Ontology (GO) terms for each gene were retrieved
225 according to the corresponding InterPro entry. KEGG pathway was mapped by the
226 constructed gene set to identify the best match for each gene [64].

227

228 **Non-coding RNA annotation**

229 The tRNAs were predicted using the program tRNAscan-SE (tRNAscan-SE, RRID:
230 SCR_010835) [65], and rRNA genes were annotated using BLASTN (BLASTN,

231 RRID: SCR_001598) tool with E-value of 1e-5 against rRNA sequences from several
232 relative plant species. miRNA and snRNA were identified by searching against the
233 Rfam (Rfam, RRID:SCR_007891) database [66] with default parameters using the
234 INFERNAL software (INFERNAL, RRID:SCR_011809)[67].

235

236 **Gene family construction**

237 OrthoFinder version 2.3.3 (OrthoFinder, RRID:SCR_017118) [68] was used to
238 classify the orthogroups of proteins from *P. salicina* and 16 other sequenced rosids
239 species, including *P. armeniaca*, *P. mume*, *P. persica*, *P. dulcis*, *P. avium*, *P. yedoensis*,
240 *M. domestica*, *P. bretschneideri*, *Pyrus communis*, *F. vesca*, *Potentilla micrantha*,
241 *Rosa chinensis*, *Rosa multiflora*, *Rubus occidentalis*, *Morus notabilis* and *A. thaliana*.

242

243 **Phylogenetic tree and divergence time estimation**

244 For phylogenetic tree construction, proteins of single-copy orthogroups (i.e., the
245 orthogroups which contain none or only one genes for each species) presented in at
246 least 70% of species were selected and aligned with MAFFT version 6.846b (MAFFT,
247 RRID: SCR_011811) [69]. After determination of the best substitution model for each
248 orthogroup with IQ-TREE version 1.7-beta12 (IQ-TREE, RRID: SCR_017254) [70],
249 the maximum likelihood phylogenetic tree across the 17 plant species was constructed
250 using IQ-TREE with the parameter (-p -bb 1000), setting *A. thaliana* as outgroup.

251 The divergence time of each node in the phylogenetic tree was estimated with
252 Bayesian Evolutionary Analysis Sampling Trees (BEAST, RRID: SCR_010228)
253 [71]. Two fossil constraints and a secondary calibration node were applied. The fossil
254 *Prunus wutuensis* (age: Early Eocene, minimum age of 55.0 Mya) and the fossil
255 *Rubus acutiformis* (age: Middle Eocene, minimum age of 41.3 Mya) were placed at
256 the stem *Prunus* and *Rubus*, respectively [72]. For the secondary calibration node, the
257 divergence of Rosoideae and Amygdaloideae at 100.7 Mya was dated according to
258 Xiang et al. [72]. The Markov chain Monte Carlo was reported 10,000,000 times with
259 1000 steps.

260

261 **Gene family expansion and contraction analysis**

262 For gene family expansion and contraction analysis, the ancestral gene content of
263 each cluster at each node was investigated with CAFÉ version 3.1 (CAFÉ, RRID:
264 SCR_005983) [73], basing on the phylogeny and gene numbers per orthogroup in
265 each species, the gene family expansions/contractions at each branch were determined
266 with p -value < 0.001 .

267

268 **Genome synteny analysis**

269 A Python version of MCScan (minspan=100; MCScan, RRID: SCR_017650;
270 <https://github.com/tanghaibao/jcvi/wiki/MCscan>) was employed to analyze the
271 synteny between the *P. salicina* genome and other genomes within *Prunus* following
272 the approaches of Haibao Tang [74].

273

274 **Positively selected gene analysis**

275 The ratios of nonsynonymous to synonymous substitutions (Ka/Ks) were calculated
276 using the Codeml program with the free-ratio model as implemented in the PAML
277 (PAML, RRID: SCR_014932) package [75]. The positive selection analysis was
278 performed using the Codeml program with the optimized branch-site model as
279 implemented in the PAML package. The positively selected genes were subjected to
280 GO functional annotation.

281

282 **Gene Ontology enrichment analysis**

283 The Gene Ontology (GO) enrichment analysis for the specific groups of genes (e.g.
284 tandem duplication and expanded genes) were performed using R package ‘topGO’
285 [76], setting all *P. salicina* genes as background. The lowest-level GO terms under
286 enrichment (p -value < 0.01) were focused, and p -value was calculated using a ‘classic’
287 algorithm with the ‘fisher’ test. The lowest-level GO terms was based on the directed
288 acyclic graph (DAG) of GO, with the parameter ‘nodeSize = 100’.

289

290 **The identification of the DUF579 family members**

291 For the identification of the DUF579 family members, the hidden Markov model
292 (HMM) profile corresponding to the DUF579 domain (PF04669) was downloaded
293 from Pfam database (<http://pfam.sanger.ac.uk/>), and subsequently exploited for the
294 genome of *P. salicina*, *P. persica*, *P. mume*, *P. armeniaca*, *P. dulcis* and *A. thaliana*
295 using HMMER 3.0. The default parameters were employed and the cutoff value was
296 set to 0.01.

297

298 **Results and Discussion**

299 **Genome sequencing and assembly**

300 We sequenced and assembled the genome of *P. salicina* using a combination of
301 short-read sequencing from Illumina Hiseq, SMRT sequencing from PacBio and Hi-C
302 technology. For the Illumina sequencing, a total of approximately 26.6 Gb ($85.4 \times$
303 coverage) short reads was obtained (Table S1). A total of ~53.0 Gb long-sequencing
304 reads were generated by PacBio Sequel platform. After removing adaptors within
305 sequences, about 52.9 Gb ($169.7 \times$ coverage) subreads were obtained (Table S1). The
306 subreads have a mean length of 13.2 kb (Table S2). About 59.1 Gb ($189.5 \times$ coverage)
307 sequencing data generated from Hi-C library was produced (Table S1). The quality of
308 Hi-C sequencing was evaluated with HiCUP [33], and the effect rate was
309 approximately 28.10% (Table S3).

310 In the genome assembly process, Illumina sequencing data were used for the
311 genome survey and polishing of preliminary contigs, PacBio long reads were used for
312 contig assembly and Hi-C reads were used for chromosome-level scaffolding. Based
313 on the total number of k-mers (19,341,904,177), the estimated *P. salicina* genome size
314 was calculated to be approximately 311.82 Mb (Figure S1). The heterozygous and
315 repeat sequencing ratios were 0.70% and 54.49%, respectively (Table S4). The *de*
316 *novo* genome assembly of *P. salicina* with a total length of 284.2 Mb (Table 1) was
317 yielded. As shown in Fig. 1, the Hi-C assisted genome assembly was anchored onto
318 the eight pseudochromosomes with lengths ranging from 23.70 to 54.53 Mb (Table
319 S5), which were designated according to the published genetic map of *P. salicina* [77].

320 Five regions of tandemly repeated telomeric repeat sequences were identified on three
321 pseudochromosomes (Table S5). The total length of pseudochromosomes accounted
322 for 96.56% of the genome sequences (Figure 1), with contig N50 of 1.78 Mb and
323 scaffold N50 of 32.32 Mb (Table1; Table S6).

324

325 **Evaluation of the genome assembly**

326 To assess the genome assembly quality, the Illumina clean data were aligned to the *P.*
327 *salicina* genome, with the mapping rate of 96.93%. A total of 98.81% assembled
328 genome was covered by the reads and the mapping coverage with at least 4×, 10×,
329 20× was 98.48 %, 98.06% and 97.13%, respectively (Table1; Table S7). The RNA-seq
330 reads were mapped against the genome assembly, and the percentage of aligned reads
331 ranged from 92.44% to 95.25% (Table1; Table S8). A total of 3,668 homozygous
332 SNPs were identified, accounting for only 0.0015% of the reference genome (Table
333 S9). The low rate of homozygous SNPs suggested that the assembly had a high base
334 accuracy. 234 Core Eukaryotic Genes (CEGs) out of the complete set of 248 CEGs
335 (94.35%) were covered by the assembly, and 229 (92.34%) of these were complete
336 (Table1; Table S10). BUSCO analysis based on single copy orthologs set showed that
337 95.7% of the expected genes were identified as complete, 1.3% were fragmented, and
338 only 3.0% were missing (Table1; Table S11). These results verified the high quality of
339 the presently generated *P. salicina* genome assembly

340

341 **Genome annotation**

342 The results of the repeat annotations found that 48.28% of the assembly was covered
343 with transposable elements (TE). Among them, long terminal repeat (LTR)
344 retrotransposons represented the greatest proportion, making up 42.10% of the
345 genome (Table1; Table S12). The TE percentage and density of duplicates resulted
346 from tandem duplications were shown in Figure 1. Tandem duplicates occurred for
347 9.8% of the genes (Table 1) and were preferentially enriched in ‘transferase activity
348 (GO: 0016758 and GO: 0016747)’ and ‘phloem development (GO: 0010088)’ (Figure

349 S2). The significant enrichment of the sieve element occlusion genes in ‘phloem
350 development’, which were involved in wound sealing of the phloem [78], might be
351 associated with specific requirements during the damage response in *P. salicina*.

352 For gene annotations, we predicted 24,448 non-redundant protein-coding genes in
353 *P. salicina*. There were 24,209 genes (~99.0%) that could be assigned to eight
354 pseudochromosomes (Table 1), and the gene density was shown in Figure 1. The
355 average number of exons per gene, and average CDS length were 4.97 and 1,157.42,
356 respectively (Table 2). Further gene functional annotation showed that 23,931 (97.9%)
357 protein-coding genes were successfully annotated (Table 1; Table S13). For the
358 identification of non-coding RNA (ncRNA) genes, a total of 627 miRNA, 960 tRNA,
359 273 rRNA and 2,023 snRNA in the *P. salicina* genome were predicted (Table S14).

360

361 **Evolution of the *P. salicina* genome**

362 The genome sequences of the representative sequenced rosid species were collected
363 and subjected to comparative genomic analysis with *P. salicina* to reveal the genome
364 evolution and divergence of *P. salicina*. A total of 15,751 orthogroups containing
365 23,265 genes were found in *P. salicina*. Moreover, 1,010 genes which were specific to
366 *P. salicina* were identified. A comparison of the predicted proteomes among the 17
367 species indicated that 9,616, 10,447, 11,098, 13,963 and 15,512 orthogroups were
368 shared between *P. salicina* and Rosids, Rosales, Rosaceae, Amygdaloideae and
369 *Prunus*, respectively.

370 The phylogenetic analysis confirmed the close relationship among *P. salicina*, *P.*
371 *mume* and *P. armeniaca*. The molecular clock of these plant genomes was also
372 calculated. The data indicated that *P. salicina* diverged from the ancestor of *P. mume*
373 and *P. armeniaca* approximately 9.05 Mya, from the ancestor of *P. persica* and
374 *P. dulcis* 11.12 Mya (Figure 2).

375 We also explored the genome syntenic blocks between *P. salicina* and the other
376 representative *Prunus* species. As shown in Fig. 3, our genome assembly of *P.*
377 *salicina* exhibited a high level of genome synteny with all the other *Prunus* genomes,
378 especially the genomes of *P. avium* and *P. dulcis*. Significantly fewer inversions were

379 found in *P. salicina* vs *P. avium* and *P. salicina* vs *P. dulcis* than that in *P. salicina* vs *P.*
380 *mume* and *P. salicina* vs *P. armeniaca*.

381

382 **Expansion and contraction of gene families in *P. salicina***

383 The gene family analysis showed that during the evolution of *P. salicina*, 146 gene
384 families were expanded and 500 gene families were contracted. The functional
385 enrichment on Gene Ontology of those expanded gene families identified 60
386 significantly enriched GO terms (p -value < 0.05) (Table S15; Figure S3).

387 It was noteworthy that genes from the expanded families were enriched in a series
388 of cell wall related processes, such as ‘cell wall polysaccharide metabolic process
389 (GO: 0010383)’, ‘hemicellulose metabolic process (GO: 0010410)’ and ‘regulation of
390 cellular biosynthetic process (GO: 0031326)’. Specially, genes in ‘xylan biosynthetic
391 process (GO: 0045492)’, which corresponded to the DUF579 family [79], were
392 significantly expanded. Further investigation showed that the major copy differences
393 were found in Clade II, which consisted of orthologs of IRX15/IRX15L [79], with
394 seven members in *P. salicina* and only two to four members in other *Prunus* species
395 (Figure 4). It was reported that IRX15 and IRX15L defined a new class of genes
396 involved in xylan biosynthesis [80, 81]. The species-specific expansion of this new
397 subclade might contribute to the relatively high content of xylan-related metabolites
398 (like xylose and xyliot) in plum [9, 10], which provided new insight into the xylan
399 metabolism in plum.

400 Moreover, the FRS (FAR1-related sequence) gene family, which played multiple
401 roles in a wide range of cellular processes [82], was also significantly expanded in the
402 phylogeny (GO: 000945), and the family expansion may be related to the genetic and
403 phenotypic diversity in *P. salicina*.

404

405 **Positively selected genes in *P. salicina***

406 The K_a/K_s ratios for all the 2,314 single-copy orthologs shared with the sequenced
407 *Prunus* species were calculated. A total of 213 candidate genes in *P. salicina*
408 underwent positive selection ($P < 0.05$). Most of them were enriched in the GO terms

409 involved in ‘monooxygenase activity (GO: 0004497)’ and ‘enzyme inhibitor activity
410 (GO: 0004857)’ (Figure S4). It was noteworthy that the category ‘monooxygenase
411 activity’ was also found in the enriched GO terms for the expanded gene families in *P.*
412 *salicina*, which might provide valuable candidate genes for further functional
413 investigations.

414

415 **Conclusions**

416 To our knowledge, this is the first report of the chromosome-level genome assembly
417 of plums using Illumina and PacBio sequencing platforms with Hi-C technology. The
418 assembly had a total size of 284.2 Mb, the contig and scaffold N50 reached 1.8 Mb
419 and 32.3 Mb, respectively. A total of 24,448 protein-coding genes were predicted, and
420 23,931 genes (97.9%) have been annotated. Phylogenetic analysis indicated that *P.*
421 *salicina* was closely related to *P. mume* and *P. armeniaca*. Expanded gene families in
422 *P. salicina* were significantly enriched in several cell-wall related processes.
423 Remarkably, the *P. salicina*-specific expansion of the xylan biosynthesis-related
424 DUF579 family provided new insight into the xylan metabolism in plums. Given the
425 economic and evolutionary importance of *P. salicina*, the genomic data in this study
426 offer a valuable resource for facilitating plum breeding programs and studying the
427 genetic basis for agronomic and adaptive divergence of plum and *Prunus* species.

428

429 **Availability of supporting data and materials**

430 This Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank
431 under the accession WERZ00000000. The version described in this paper is version
432 WERZ01000000. The raw sequencing data are available through the NCBI Sequence
433 Read Archive (SRA) via accession numbers from SRR10233497 to SRR10233505,
434 via the Project PRJNA574159. The transcriptome data are available through the NCBI
435 SRA (from SRR10235674 to SRR10235679). The genome data have also been
436 submitted to *Genome Database for Rosaceae* (Accession number: tfGDR1044). All
437 annotation tables containing results of an analysis of the draft genome are available at

438 Figshare [83]. Supporting data is also available via the *GigaScience* database GigaDB
439 [84].

440

441 **Additional files**

442 **Table S1** Statistics of *P. salicina* genome sequencing data.

443 **Table S2** Statistics of characteristics of PacBio long-reads.

444 **Table S3** Statistics of Hi-C sequencing data.

445 **Table S4** Estimation of the genome size using k-mer analysis.

446 **Table S5** Summary of assembled 8 pseudochromosomes of *P. Salicina*.

447 **Table S6** Summary of the genome assembly of *P. Salicina*.

448 **Table S7** Statistics of mapping ratio in genome.

449 **Table S8** Summary of the transcriptome and their mapping rate on the genome
450 assembly.

451 **Table S9** Number and density of SNPs in *P. salicina* genome.

452 **Table S10** Assessment of CEGMA.

453 **Table S11** Summary of BUSCO analysis results according to prediction.

454 **Table S12** Detailed classification of repeat sequences.

455 **Table S13** Statistics of functional annotation.

456 **Table S14** Summary of non-coding RNA.

457 **Table S15** List of the Gene ontology terms significantly enriched in the expanded
458 gene families of *P. salicina*

459 **Figure S1** 17-mer frequency distribution in *P. salicina* genome.

460 **Figure S2** Gene ontology enrichment of the tandemly duplicated genes in *P. salicina*.

461 **Figure S3** Gene ontology enrichment of *P. salicina*-expanded genes.

462 **Figure S4** Gene ontology enrichment of the positively selected genes in *P. salicina*.

463

464 **Authors' Contributions**

465 Y.H.H. conceived the study. C.Y.L., C.F. and J.T.W. performed bioinformatics
466 analysis. W.Z.P., J.J.H. and J.J.P. collected the samples and extracted the DNA. C.Y. L.
467 and C. F. wrote the manuscript. All authors read and approved the final manuscript.

468

469 **Abbreviations**

470 BLAST: Basic Local Alignment Search Tool; BEAST: Bayesian Evolutionary
471 Analysis Sampling Trees; bp: base pair; BUSCO: Benchmarking Universal
472 Single-Copy Orthologs; CEGMA: Core Eukaryotic Genes Mapping Approach; CTAB:
473 cetyltrimethylammonium bromide; EVM: EVIDENCEModeler; Gb: gigabase pair; GO:
474 Gene Ontology; Hi-C: high-throughput chromosome conformation capture; kb:
475 kilobase pair; KEGG: Kyoto Encyclopedia of Genes and Genomes; Mb: megabase
476 pair; miRNA: microRNA; Mya: million years ago; NCBI: National Center for
477 Biotechnology Information; PacBio: Pacific Biosciences; PAML: phylogenetic
478 analysis by maximum likelihood; PASA: Program to Assemble Spliced Alignments;
479 RNA-seq: RNA sequencing; rRNA: ribosomal RNA; SMRT: single-molecule
480 real-time; SnRNA, small nuclear RNA; SNP: single-nucleotide polymorphism; TRF:
481 Tandem Repeats Finder; tRNA: transfer RNA.

482

483 **Funding**

484 This work was financially supported by The Industry University Research
485 Collaborative Innovation Major Projects of Guangzhou Science Technology
486 Innovation Commission (201704020021) and Modern Agricultural Industry
487 Technology System of Guangdong Province (2016LM1128).

488

489 **Competing interests**

490 The authors declare no competing interests.

491

492 **Figure Legends**

493 **Figure 1** The genome and photograph of *P. salicina*. Landscape of the *P. salicina*
494 genome, comprising 8 pseudochromosomes that cover ~96.56% of assembly (A);
495 Concentric circles, from outermost to innermost, showing TE percentage (red; B);
496 gene density (green; C); density of duplicates resulted from tandem duplications (blue;
497 D); (E) photograph of *P. salicina*.

498

499 **Figure 2** Evolution of *P. salicina* genome and orthogroups. (A) The phylogeny,
500 divergence time and orthogroup expansions/contractions for 17 rosids species. The
501 tree was constructed by maximum likelihood method using 341 single copy
502 orthogroups. All nodes have 100% bootstrap support. Divergence time was estimated
503 on a basis of three calibration points (blue circles). Blue bar indicates 95% HPD
504 (highest posterior density) for each node. The numbers in red and green indicate the
505 numbers of orthogroups that have expanded and contracted along particular branches,
506 respectively. (B) The comparison of genes among 17 rosids. The grey bars indicate
507 the genes belonging to 9,616 rosids-shared orthogroups in each of 17 rosids. The grey
508 + green bars indicate the genes belonging to 10,447 rosales-shared orthogroups in
509 each of 16 rosales. The grey + green + pink bars indicate the genes belonging to
510 11,098 Rosaceae-shared orthogroups in each of 15 Rosaceae. The grey + green + pink
511 + yellow bars indicate the genes belonging to 13,963 rosaceae-shared orthogroups in
512 each of ten Amygdaloideae. The grey + green + pink + yellow + blue bars indicate the
513 genes belonging to 15,512 *Prunus*-shared orthogroups in each of seven *Prunus*
514 species. The red and stripe bars indicate the genes in species-specific orthogroups and
515 unassigned genes, respectively. The white bars indicate the remaining genes for each
516 genome.

517

518 **Figure 3** Chromosome-level collinearity patterns between *P. salicina*, *P. mume* and *P.*
519 *armeniaca* (A) and between *P. salicina*, *P. avium* and *P. dulcis* (B). The numbers
520 indicate the pseudochromosome order generated from the original genome sequence.

521 The pseudochromosome 2 and 6 in *P. armeniaca* and *P. mume* are reversed. Each gray
522 line represents one block. The inverted regions are highlighted with brown color.

523 .

524 **Figure 4** The significant expansion of the DUF579 family members in *P. salicina*. (A)
525 Phylogenetic tree of the DUF579 proteins from *P. salicina* (red circle), *P. persica*
526 (hollow inverted triangle), *P. mume* (solid triangle), *P. armeniaca* (hollow diamond), *P.*
527 *dulcis* (solid diamond) and *A. thaliana* (solid square). (B) The summary of the
528 numbers of clade members in DUF579 family.

529

530 **Table 1** Summary of genome assembly and annotation for *P. salicina*

531

	Number or percentage
Assembly feature	
Total length of scaffolds (bp)	284,209,110
Number of scaffolds	75
N50 of scaffolds (bp)	32,324,625
Total length of contigs (bp)	284,189,410
Number of contigs	272
N50 of contigs (bp)	1,777,944
Mapping rate by reads from short-insert libraries	96.93%
Assembled CEGs	94.35%
Completely assembled CEGs	92.34%
Complete BUSCOs	95.7%
Complete and single-copy BUSCOs	86.5%
Complete and duplicated BUSCOs	9.2%
Fragmented BUSCOs	1.3%
Missing BUSCOs	3.0%
RNA-Seq evaluation	92.44%-95.25%
Genome annotation	
Percentage of transposable elements (TE)	48.28%
Percentage of long terminal repeat (LTR) retrotransposon	42.1%
No. of predicted protein-coding genes	24,448
No. of genes assigned to pseudochromosomes	24,209 (99.0%)
No. of genes annotated to public database	23,930 (97.9%)
No. of genes annotated to GO database	13,484 (55.2%)
No. of genes duplicated by tandem duplications	2,384(9.8%)

532

533 **Table 2** Statistics of predicted protein-coding genes.

	Gene set	Number	Average transcript length (bp)	Average CDS length (bp)	Average exons per gene	Average exons length (bp)	Average intron length (bp)
<i>De novo</i> prediction	Augustus	23,592	2,627.71	1167.83	4.80	243.43	384.45
	GlimmerHMM	39,985	5,450.51	747.07	3.14	238.12	2200.59
	SNAP	24,882	2,876.50	728.45	4.22	172.73	667.66
	Geneid	33,780	3,829.40	899.99	4.44	202.74	851.78
	Genscan	21,882	8,251.09	1355.87	6.34	213.98	1292.13
Homolog prediction	<i>Pyrus bretschneideri</i>	20,265	3,119.83	1356.17	4.74	286.35	472.06
	<i>Malus domestica</i>	20,010	2,920.17	1361.30	4.65	292.56	426.72
	<i>Prunus mume</i>	23,064	3,038.66	1346.19	4.78	281.67	447.84
	<i>Prunus persica</i>	28,915	2,296.51	1099.56	4.06	270.55	390.64
	<i>Arabidopsis thaliana</i>	28,284	2,071.73	973.28	3.67	265.51	412.07
	<i>Fragaria vesca</i>	22,927	2,994.24	1380.61	4.59	300.66	449.24
RNA-seq	<i>Prunus avium</i>	22,715	3,077.20	1351.28	4.74	284.86	461.03
	PASA	196,264	3,913.86	1008.68	5.16	195.60	698.88
	Transcripts	42,450	11,076.28	2360.92	6.85	344.83	1490.64
	EVM	27,981	2,736.70	1061.73	4.57	232.52	469.68
	PASA-update*	27,594	2,784.15	1092.82	4.64	235.59	464.83
	Final set*	24,448	2,988.45	1157.42	4.97	233.09	461.72

534 * UTR regions were contained

535

536 **References**

- 537 1. Roussos PA, Efstathios N, Intidhar B, Denaxa N-K and Tsafouros A. Plum (*Prunus*
538 *domestica* L. and *P. salicina* Lindl.). In: Monique Simmonds VRP, editor. Nutritional
539 Composition of Fruit Cultivars. Elsevier; 2016. p. 639 - 666.
- 540 2. Topp BL, Russell DM, Neumüller M, Dalbó MA and Liu W. Plum. In: Maria Luisa Badenes
541 DHB, editor. Fruit Breeding. Springer; 2012. p. 571-621.
- 542 3. Hartmann W and Neumüller M. Plum breeding. In: Shri Mohan Jain PMP, editor. Breeding
543 Plantation Tree Crops: Temperate Species. Springer; 2009. p. 161-231.
- 544 4. Okie W and Hancock J. Plums. In: Hancock JF, editor. Temperate Fruit Crop Breeding.
545 Springer Science & Business Media; 2008. p. 337-358.
- 546 5. Esmenjaud D and Dirlewanger E. Plum. In: Kole C, editor. Genome Mapping and Molecular
547 Breeding in Plants. Springer; 2007. p. 119-135.
- 548 6. Guerra M and Rodrigo J. Japanese plum pollination: A review. *SCI Hortic-Amsterdam*
549 2015;**197**:674-686.
- 550 7. Rennie EA and Scheller HV. Xylan biosynthesis. *Curr Opin Biotech* 2014;**26**:100-107.
- 551 8. Brummell DA and Schröder R. Xylan metabolism in primary cell walls. *NZ J Forestry Sci.*
552 2009;**39**:125-143.
- 553 9. Renard CMGC and Ginies C. Comparison of the cell wall composition for flesh and skin
554 from five different plums. *Food Chem* 2009;**114**(3):1042-1049.
- 555 10. Arcaño YD, García ODV, Mandelli D, Carvalho WA and Pontes LAM. Xylitol: A review on
556 the progress and challenges of its production by chemical route. *Catal Today* 2020;**344**:2-14.
- 557 11. Aranzana MJ, Decroocq V, Dirlewanger E, Eduardo I, Gao ZS, Gasic K, et al. *Prunus*
558 genetics and applications after *de novo* genome sequencing: achievements and prospects.
559 *Hortic Res* 2019;**6** (1):1-25.
- 560 12. Velasco R, Zharkikh A, Affourtit J, Dhingra A, Cestaro A, Kalyanaraman A, et al. The
561 genome of the domesticated apple (*Malus* × *domestica* Borkh.). *Nat Genet* 2010;**42**
562 (10):833-839.
- 563 13. Chen X, Li S, Zhang D, Han M, Jin X, Zhao C, et al. Sequencing of a wild apple (*Malus*
564 *baccata*) genome unravels the differences between cultivated and wild apple species
565 regarding disease resistance and cold tolerance. *G3: Genes, Genomes, Genet* 2019;**9**
566 (7):2051-2060.
- 567 14. Zhang L, Hu J, Han X, Li J, Gao Y, Richards CM, et al. A high-quality apple genome
568 assembly reveals the association of a retrotransposon and red fruit colour. *Nat Commun*
569 2019;**10** (1):1-13.
- 570 16. Verde I, Abbott AG, Scalabrin S, Jung S, Shu S, Marroni F, et al. The high-quality draft
571 genome of peach (*Prunus persica*) identifies unique patterns of genetic diversity,
572 domestication and genome evolution. *Nat Genet* 2013;**45**(5):487-494.
- 573 17. Linsmith G, Rombauts S, Montanari S, Deng CH, Celton J-M, Guérif P, et al.

- 574 Pseudo-chromosome-length genome assembly of a double haploid “Bartlett” pear (*Pyrus*
575 *communis* L.) GigaScience 2019; 8 (12):giz138.
- 576 18. Wu J, Wang Z, Shi Z, Zhang S, Ming R, Zhu S, et al. The genome of the pear (*Pyrus*
577 *bretschneideri* Rehd.). Genome Res 2013;**23**(2):396-408.
- 578 19. Chagné D, Crowhurst RN, Pindo M, Thrimawithana A, Deng C, Ireland H, et al. The draft
579 genome sequence of European pear (*Pyrus communis* L. ‘Bartlett’). PloS One 2014;**9**
580 (4):e92644.
- 581 20. Dong X, Wang Z, Tian L, Zhang Y, Qi D, Huo H, et al. *De novo* assembly of a wild pear
582 (*Pyrus betuleafolia*) genome. Plant Biotechnol J 2020;**18**(2):581-595
- 583 21. Shulaev V, Sargent DJ, Crowhurst RN, Mockler TC, Folkerts O, Delcher AL, et al. The
584 genome of woodland strawberry (*Fragaria vesca*). Nat Genet 2011;**43**(2):109-116.
- 585 22. Edger PP, Poorten TJ, VanBuren R, Hardigan MA, Colle M, McKain MR, et al. Origin and
586 evolution of the octoploid strawberry genome. Nat Genet 2019;**51**(3):541-547.
- 587 23. Alioto T, Alexiou KG, Bardil A, Barteri F, Castanera R, Cruz F, et al. Transposons played a
588 major role in the diversification between the closely related almond and peach genomes:
589 Results from the almond genome sequence. Plant J 2020;**101**(2):455-472.
- 590 24. Sánchez-Pérez R, Pavan S, Mazzeo R, Moldovan C, Cigliano RA, Del Cueto J, et al.
591 Mutation of a bHLH transcription factor allowed almond domestication. Science 2019; **364**
592 (6445):1095-1098.
- 593 25. VanBuren R, Bryant D, Bushakra JM, Vining KJ, Edger PP, Rowley ER, et al. The genome of
594 black raspberry (*Rubus occidentalis*). Plant J 2016;**87**(6):535-547.
- 595 26. Shirasawa K, Isuzugawa K, Ikenaga M, Saito Y, Yamamoto T, Hirakawa H, et al. The
596 genome sequence of sweet cherry (*Prunus avium*) for use in genomics-assisted breeding.
597 DNA Res 2017;**24**(5):499-508.
- 598 27. Wang J, Liu W, Zhu D, Hong P, Zhang S, Xiao S, et al. Chromosome-scale genome assembly
599 of sweet cherry (*Prunus avium* L.) cv. Tieton obtained using long-read and Hi-C sequencing.
600 Hort Res 2020;**7** (1):1-11.
- 601 28. Jiang F, Zhang J, Wang S, Yang L, Luo Y, Gao S, et al. The apricot (*Prunus armeniaca* L.)
602 genome elucidates Rosaceae evolution and beta-carotenoid synthesis. Hort Res 2019;**6**
603 (1):1-12.
- 604 29. Campoy JA, Sun H, Goel M, Jiao W-B, Folz-Donahue K, Kukat C, et al. Chromosome-level
605 and haplotype-resolved genome assembly enabled by high-throughput single-cell sequencing
606 of gamete genomes. BioRxiv. 2020.
- 607 30. Jiang S, An H, Xu F and Zhang X. Chromosome-level genome assembly and annotation of
608 the loquat (*Eriobotrya japonica*) genome. GigaScience 2020;**9**(3):giaa015.
- 609 31. Zhang Q, Chen W, Sun L, Zhao F, Huang B, Yang W, et al. The genome of *Prunus mume*.
610 Nat Commun 2012;**3**:1318.
- 611 32. Lodhi MA, Ye G-N, Weeden NF and Reisch BI. A simple and efficient method for DNA
612 extraction from grapevine cultivars and *Vitis* species. Plant Mol Biol Rep. 1994;**12** (1):6-13.

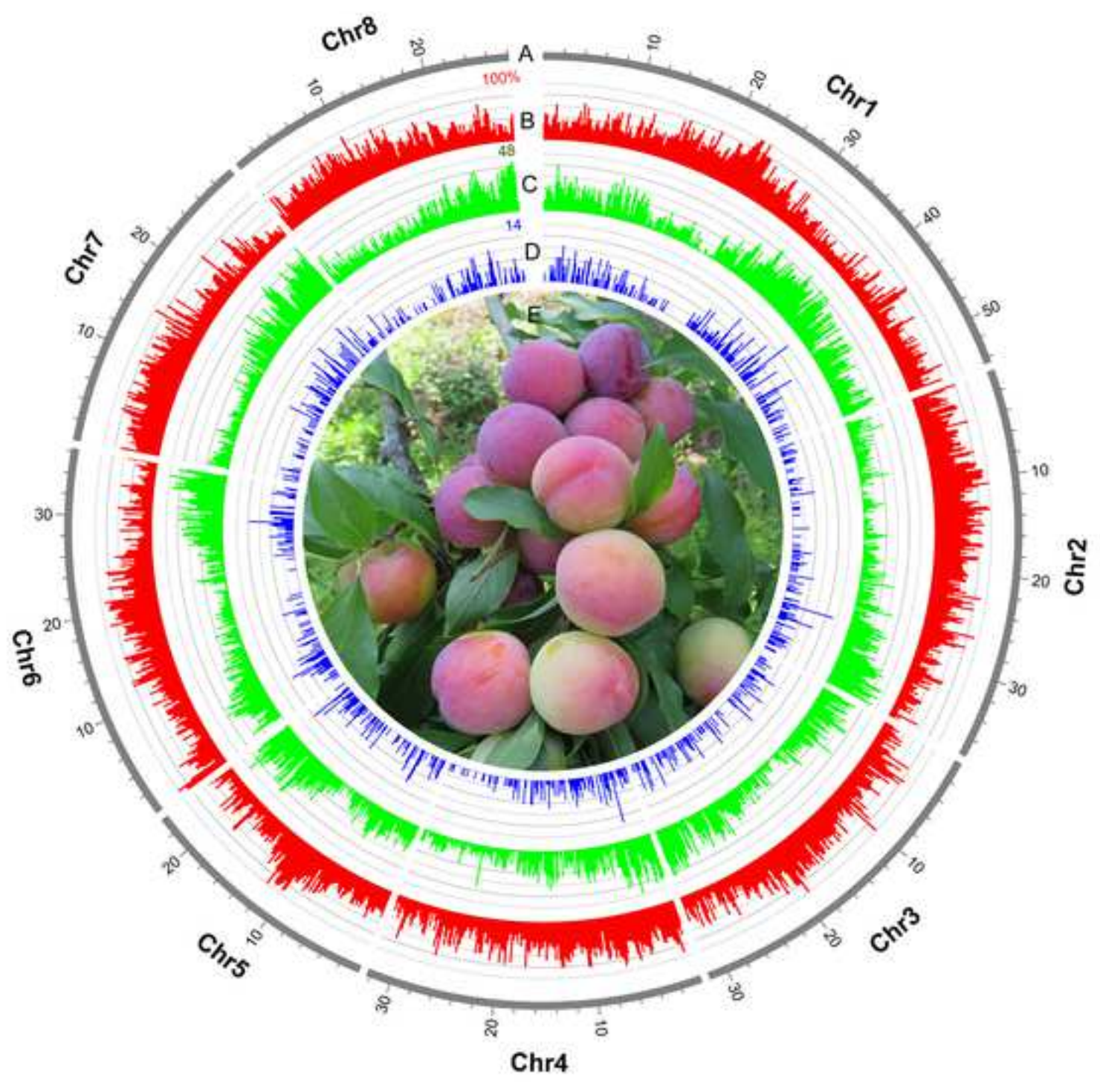
- 613 33. Wingett S, Ewels P, Furlan-Magaril M, Nagano T, Schoenfelder S, Fraser P, et al. HiCUP:
614 pipeline for mapping and processing Hi-C data. *F1000Res* 2015;**4**:1310.
- 615 34. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved
616 memory-efficient short-read *de novo* assembler. *Gigascience* 2012;**1** (1):2047-217X-1-18.
- 617 35. Chin C-S, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. Phased
618 diploid genome assembly with single-molecule real-time sequencing. *Nat Methods* 2016;**13**
619 (12):1050-1054.
- 620 36. Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Nonhybrid,
621 finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods*
622 2013;**10**(6):563-569.
- 623 37. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated
624 tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS*
625 *One* 2014;**9** (11):e112963.
- 626 38. Roach MJ, Schmidt SA and Borneman ARJB. Purge Haplotigs: allelic contig reassignment
627 for third-gen diploid genome assemblies. *BMC Bioinformatics* 2018;**19** (1):460.
- 628 39. Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
629 *Bioinformatics* 2009;**25** (14):1754-1760.
- 630 40. Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO and Shendure J. Chromosome-scale
631 scaffolding of *de novo* genome assemblies based on chromatin interactions. *Nat Biotechnol*
632 2013;**31** (12):1119-1125.
- 633 41. Robinson JT, Turner D, Durand NC, Thorvaldsdottir H, Mesirov JP and Aiden EL. Juicebox.
634 js provides a cloud-based visualization system for Hi-C data. *Cell Syst* 2018;**6**(2):256-258.
635 e1.
- 636 42. Kim D, Langmead B and Salzberg SL. HISAT: a fast spliced aligner with low memory
637 requirements. *Nat Methods* 2015;**12** (4):357-360.
- 638 43. Parra G, Bradnam K and Korf I. CEGMA: a pipeline to accurately annotate core genes in
639 eukaryotic genomes. *Bioinformatics* 2007;**23** (9):1061-1067.
- 640 44. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO:
641 assessing genome assembly and annotation completeness with single-copy orthologs.
642 *Bioinformatics* 2015;**31**(19):3210-3212.
- 643 45. Tarailo-Graovac M and Chen N. Using RepeatMasker to identify repetitive elements in
644 genomic sequences. *Curr Protoc Bioinf* 2009;**25** (1):4-10.
- 645 46. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O and Walichiewicz J. Repbase
646 Update, a database of eukaryotic repetitive elements. *Cytogenet Genome Res* 2005;**110**
647 (1-4):462-467.
- 648 47. Price AL, Jones NC and Pevzner PA. *De novo* identification of repeat families in large
649 genomes. *Bioinformatics* 2005;**21** (suppl_1):i351-i358.
- 650 48. Xu Z and Wang H. LTR_FINDER: an efficient tool for the prediction of full-length LTR
651 retrotransposons. *Nucleic Acids Res* 2007;**35** (suppl_2):W265-W268.

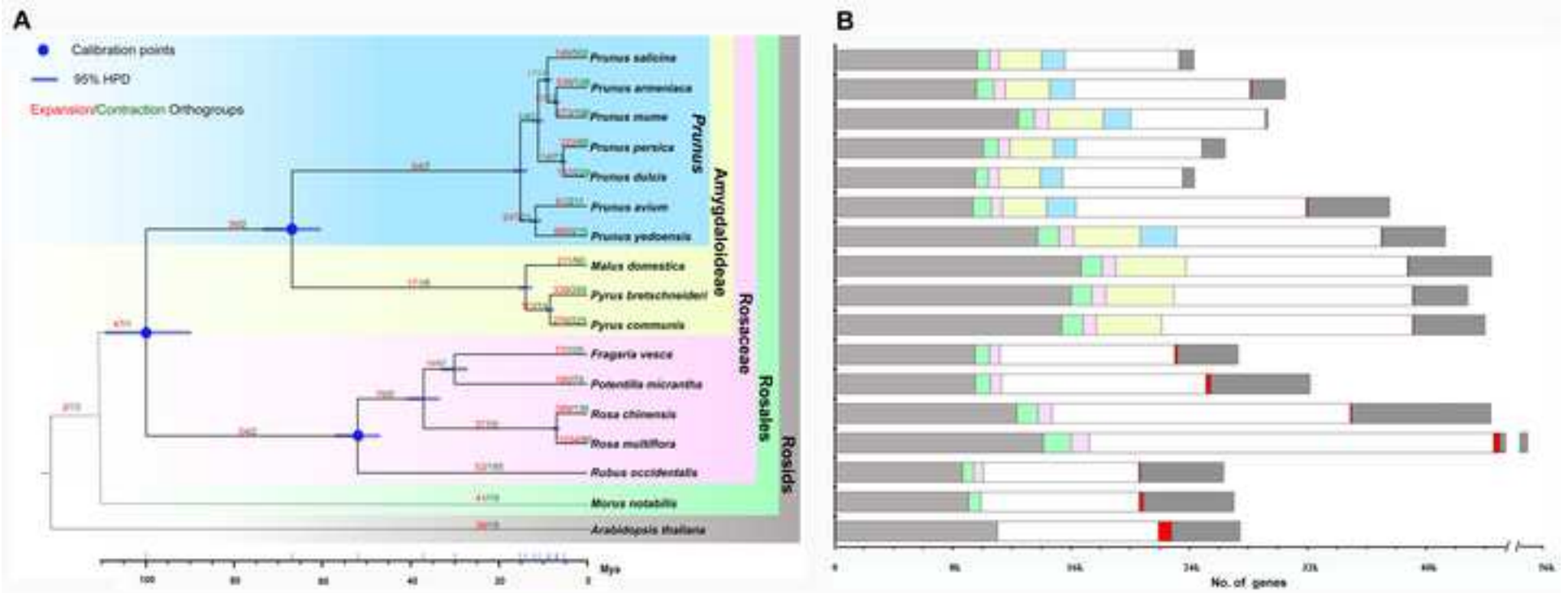
- 652 49. Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*
653 1999;**27** (2):573-580.
- 654 50. Gertz EM, Yu Y-K, Agarwala R, Schäffer AA and Altschul SF. Composition-based statistics
655 and translated nucleotide searches: improving the TBLASTN module of BLAST. *BMC Biol*
656 2006;**4** (1):1-14.
- 657 51. Birney E, Clamp M and Durbin R. GeneWise and genomewise. *Genome Res* 2004;**14**
658 (5):988-995.
- 659 52. Stanke M, Steinkamp R, Waack S and Morgenstern B. AUGUSTUS: a web server for gene
660 finding in eukaryotes. *Nucleic Acids Res* 2004;**32** (suppl_2):W309-W312.
- 661 53. Majoros WH, Pertea M and Salzberg SL. TigrScan and GlimmerHMM: two open source ab
662 initio eukaryotic gene-finders. *Bioinformatics* 2004;**20** (16):2878-2879.
- 663 54. Korf I. Gene finding in novel genomes. *BMC Bioinf* 2004;**5** (1):59.
- 664 55. Blanco E, Parra G and Guigó R. Using geneid to identify genes. *Curr Protoc Bioinf* 2007;**18**
665 (1):4.3. 1-4.3. 28.
- 666 56. Burge C and Karlin S. Prediction of complete gene structures in human genomic DNA. *J Mol*
667 *Biol* 1997;**268** (1):78-94.
- 668 57. Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT and Salzberg SL. StringTie
669 enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*
670 2015;**33**(3):290-295.
- 671 58. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. *De novo*
672 transcript sequence reconstruction from RNA-seq using the Trinity platform for reference
673 generation and analysis. *Nat Protoc* 2013;**8** (8):1494-1512.
- 674 59. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic gene
675 structure annotation using EVIDENCEModeler and the Program to Assemble Spliced
676 Alignments. *Genome Biol* 2008;**9** (1):R7.
- 677 60. Bairoch A and Apweiler R. The SWISS-PROT protein sequence database and its supplement
678 TrEMBL in 2000. *Nucleic Acids Res* 2000;**28** (1):45-48.
- 679 61. Mulder N and Apweiler R. InterPro and InterProScan: tools for protein sequence classification
680 and comparison. *Methods Mol Biol* 2007; **396**:59-70.
- 681 62. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, et al. Pfam: the protein
682 families database. *Nucleic Acids Res* 2013;**42** (D1):D222-D230.
- 683 63. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5:
684 genome-scale protein function classification. *Bioinformatics* 2014;**30** (9):1236-1240.
- 685 64. Kanehisa M and Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids*
686 *Res* 2000;**28** (1):27-30.
- 687 65. Lowe TM and Eddy SR. tRNAscan-SE: a program for improved detection of transfer RNA
688 genes in genomic sequence. *Nucleic Acids Res* 1997;**25** (5):955-964.
- 689 66. Griffiths-Jones S, Bateman A, Marshall M, Khanna A and Eddy SR. Rfam: an RNA family
690 database. *Nucleic Acids Res* 2003;**31**(1):439-441.

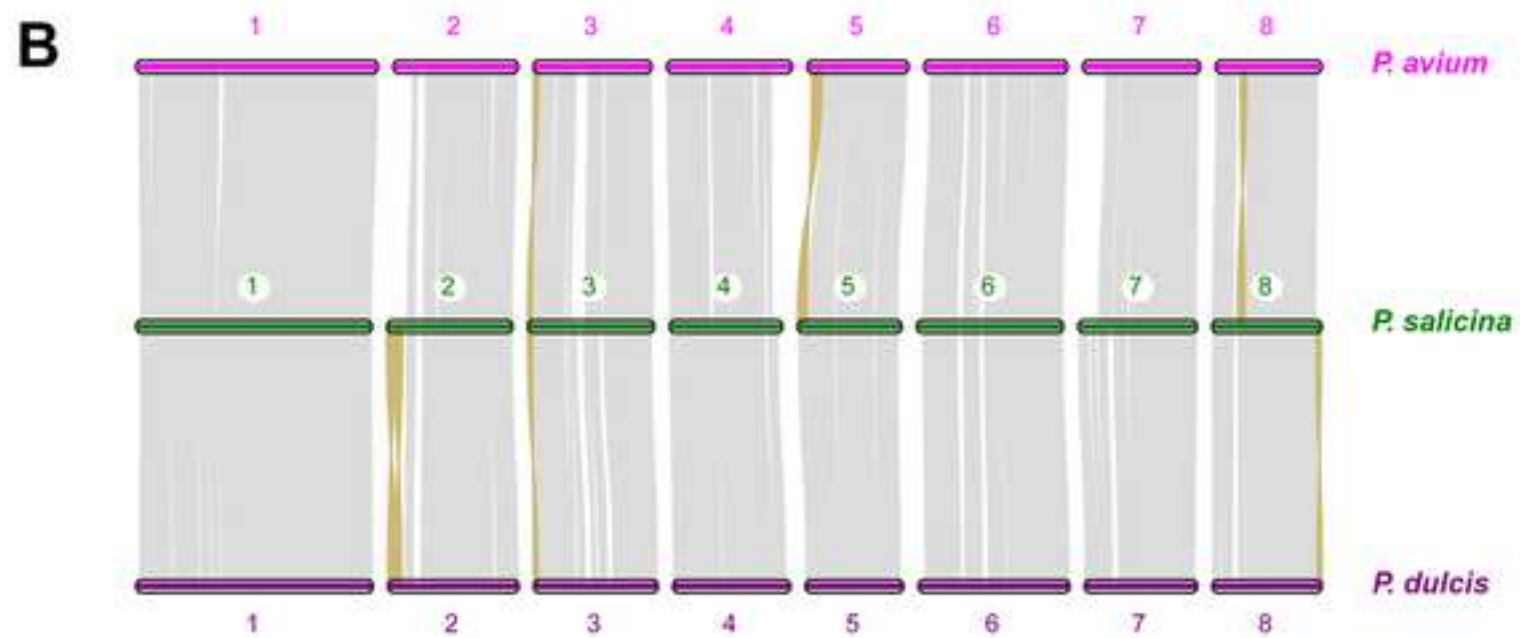
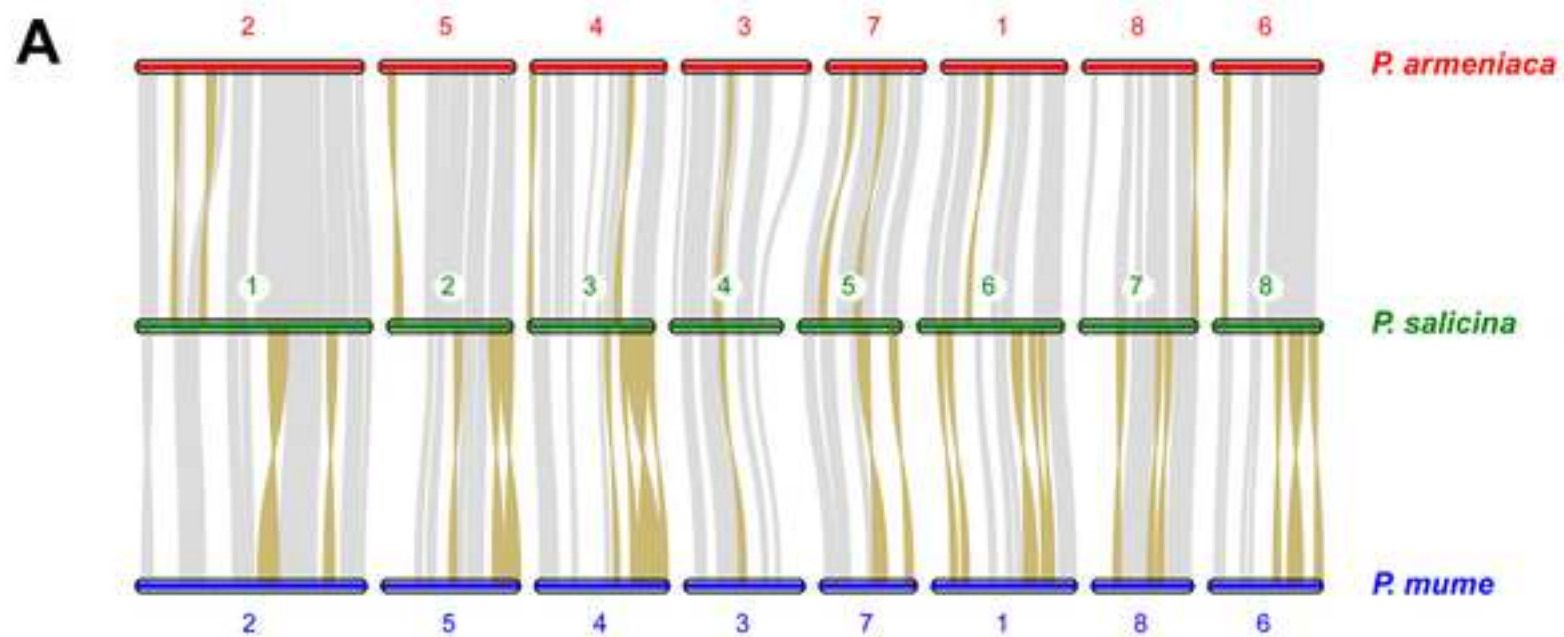
- 691 67. Nawrocki EP and Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches.
692 Bioinformatics 2013;**29** (22):2933-2935.
- 693 68. Emms DM and Kelly S. OrthoFinder: solving fundamental biases in whole genome
694 comparisons dramatically improves orthogroup inference accuracy. Genome Biol 2015;**16**
695 (1):157.
- 696 69. Katoh K and Standley DM. MAFFT multiple sequence alignment software version 7:
697 improvements in performance and usability. Mol Biol Evol 2013;**30** (4):772-780.
- 698 70. Nguyen L-T, Schmidt HA, Von Haeseler A and Minh BQ. IQ-TREE: a fast and effective
699 stochastic algorithm for estimating maximum-likelihood phylogenies. Mol Biol Evol 2015;**32**
700 (1):268-274.
- 701 71. Drummond AJ and Rambaut A. BEAST: Bayesian evolutionary analysis by sampling trees.
702 BMC Evol Biol 2007;**7** (1):1-8.
- 703 72. Xiang Y, Huang C-H, Hu Y, Wen J, Li S, Yi T, et al. Evolution of Rosaceae fruit types based
704 on nuclear phylogeny in the context of geological times and genome duplication. Mol Biol
705 Evol 2017;**34** (2):262-281.
- 706 73. De Bie T, Cristianini N, Demuth JP and Hahn MW. CAFE: a computational tool for the study
707 of gene family evolution. Bioinformatics 2006;**22** (10):1269-1271.
- 708 74. Tang H. Multiple collinearity scan—mcsan. 2009.
- 709 75. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. Mol Biol Evol
710 2007;**24**(8):1586-1591.
- 711 76. Alexa A and Rahnenführer J. Gene set enrichment analysis with topGO. Bioconductor
712 Improv. 2009;27.
- 713 77. Carrasco B, González M, Gebauer M, García-González R, Maldonado J and Silva HJPo.
714 Construction of a highly saturated linkage map in Japanese plum (*Prunus salicina* L.) using
715 GBS for SNP marker calling. PloS one 2018;**13** (12): e0208032.
- 716 78. Ernst AM, Jekat SB, Zielonka S, Müller B, Neumann U, Rüping B, et al. Sieve element
717 occlusion (SEO) genes encode structural phloem proteins involved in wound sealing of the
718 phloem. P Natl Acad Sci USA 2012;**109** (28): E1980-E1989.
- 719 79. Temple H, Mortimer JC, Tryfona T, Yu X, Lopez - Hernandez F, Sorieul M, et al. Two
720 members of the DUF 579 family are responsible for arabinogalactan methylation in
721 Arabidopsis. Plant Direct 2019;**3** (2):e00117.
- 722 80. Jensen JK, Kim H, Cocuron JC, Orlor R, Ralph J and Wilkerson CG. The DUF579 domain
723 containing proteins IRX15 and IRX15-L affect xylan synthesis in Arabidopsis. Plant J
724 2011;**66** (3):387-400.
- 725 82. Brown D, Wightman R, Zhang Z, Gomez LD, Atanassov I, Bukowski JP, et al. Arabidopsis
726 genes IRREGULAR XYLEM (IRX15) and IRX15L encode DUF579 - containing proteins
727 that are essential for normal xylan deposition in the secondary cell wall. Plant J 2011;**66**
728 (3):401-413.
- 729 82. Ma L and Li G. FAR1-related sequence (FRS) and FRS-related factor (FRF) family proteins

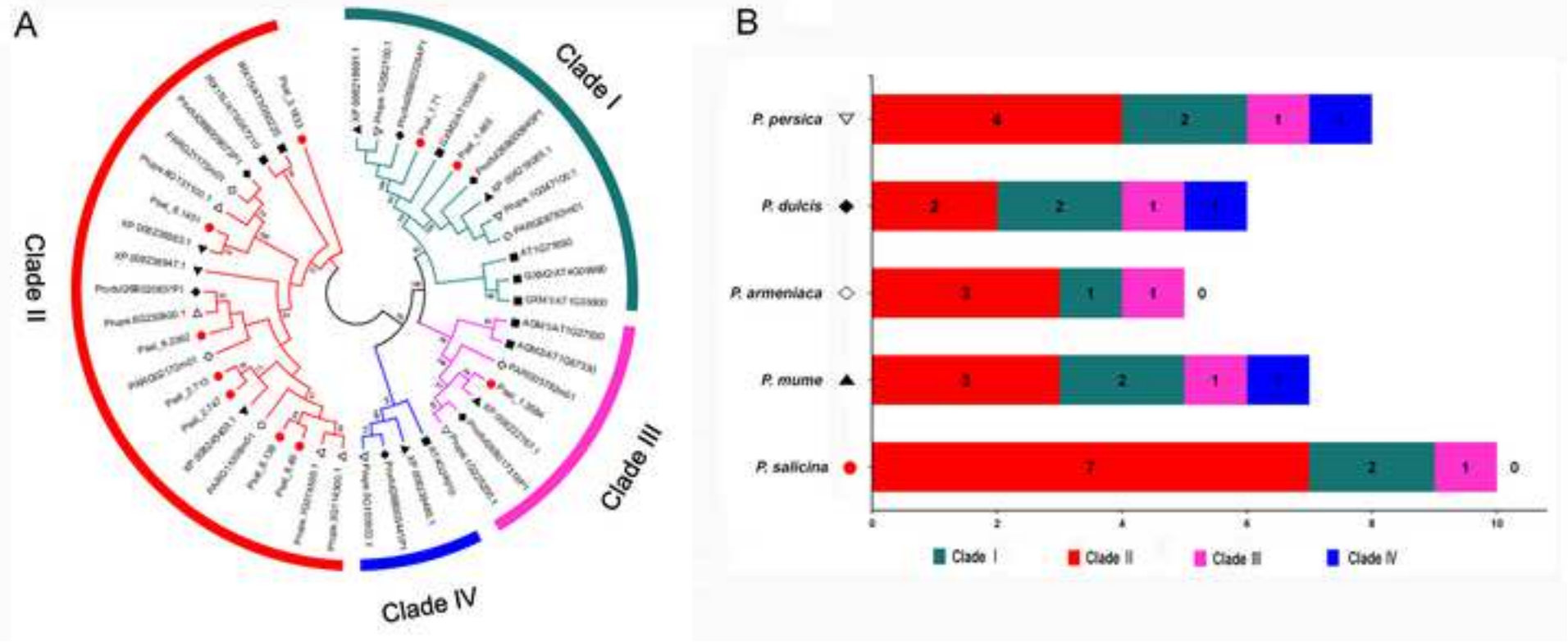
730 in *Arabidopsis* growth and development. Front Plant Sci 2018; 9:692.
731 83. Liu C, Feng C, Peng W, Hao J, Pan J, He Y. Annotation results of *Prunus salicina* genome
732 Figshare 2020. <https://doi.org/10.6084/m9.figshare.9973469>.
733 84 Liu C, Feng C, Peng W, Hao J, Wang J, Pan J, He Y Supporting data for "The
734 chromosome-level draft genome of a diploid plum (*Prunus salicina*)"
735 GigaScience Database. 2020. <http://dx.doi.org/10.5524/100811>
736

Figure 1













Click here to access/download
Supplementary Material
Supplementary Tables-9-28.xlsx

