

1 **The structure of the tetraploid sour cherry 'Schattenmorelle' (*Prunus cerasus* L.)**
2 **genome reveals insights into its segmental allopolyploid nature**

3

4 Thomas W. Wöhner¹, Ofere F. Emeriewen¹, Alexander H.J. Wittenberg², Koen
5 Nijbroek², Rui Peng Wang², Evert-Jan Blom², Jens Keilwagen⁴, Thomas Berner⁴,
6 Katharina J. Hoff³, Lars Gabriel³, Hannah Thierfeldt³, Omar Almolla⁵, Lorenzo Barchi⁵,
7 Mirko Schuster¹, Janne Lempe¹, Andreas Peil¹, Henryk Flachowsky¹

8

9 ¹Julius Kühn Institute (JKI) – Federal Research Centre for Cultivated Plants, Institute
10 for Breeding Research on Fruit Crops, Pillnitzer Platz 3a, D-01326, Dresden, Germany

11 ²KeyGene N.V., P.O. Box 216, 6700 AE Wageningen, The Netherlands

12 ³Institute of Mathematics and Computer Science, University of Greifswald, Walther-
13 Rathenau-Str. 47, 17489 Greifswald, Germany

14 ⁴Julius Kühn Institute (JKI) – Federal Research Centre for Cultivated Plants, Institute
15 for Biosafety in Plant Biotechnology, Erwin-Baur-Str. 27, D-06484 Quedlinburg,
16 Germany

17 ⁵DISAFA – Plant genetics, University of Turin, Grugliasco (TO), 10095 Italy

18

19

20 *corresponding author: thomas.woehner@julius-kuehn.de

21

22 **Keywords: genome assembly, *P. cerasus*, sour cherry, tetraploid**

23

24 **Abstract**

25 Sour cherry (*Prunus cerasus* L.) is an economically important allotetraploid cherry
26 species believed to have evolved in the Caspian Sea and Black Sea regions. How,
27 when and where exactly the evolution of this species took place is unclear. It resulted
28 from a hybridization of the tetraploid ground cherry (*Prunus fruticosa* Pall.) and an
29 unreduced (2n) pollen of the diploid ancestor sweet cherry (*P. avium* L.). Some
30 indications implement that the genome of sour cherry is segmental allopolyploid, but
31 how it is structured and to what extent is unknown. To get an insight, the genome of
32 the sour cherry cultivar 'Schattenmorelle' was sequenced at ~400x using Illumina

33 NovaSeq™ short-read and Oxford Nanopore long-read technologies (ONT R9.4.1
34 PromethION). Additionally, the transcriptome of ‘Schattenmorelle’ was sequenced
35 using PacBio Sequel II SMRT cell sequencing at ~300x. The final assembly resulted
36 in a ~629 Mbp long pseudomolecule reference genome, which could be separated
37 into two subgenomes each split into eight chromosomes. Subgenome *Pce_S_a* which
38 originates from *P. avium* has a length of 269 Mbp, whereas subgenome *Pce_S_f* which
39 originates from *P. fruticosa* has a length of 299.5 Mbp. The length of unassembled
40 contigs was 60 Mbp. The genome of the sour cherry shows a size-reduction
41 compared to the genomes of its ancestral species. It also shows traces of
42 homoeologous sequence exchanges throughout the genome. Comparative positional
43 sequence and protein analyses provided evidence that the genome of sour cherry is
44 segmental allotetraploid and that it has evolved in a very recent event in the past.

45

46 **Background**

47 Cherries include several species of the genus *Prunus*, which belong to the sub-family
48 *Spiraeoideae* in the plant family Rosaceae (Potter et al. 2007). Although the fruits of
49 several cherry species are used for consumption, only a few of them are grown and
50 marketed on an economically significant scale (Quero-García et al. 2019). The two
51 economically most important cherry species worldwide are the sweet cherry (*Prunus*
52 *avium* L.) and the sour cherry (*Prunus cerasus* L.). Both species are thought to have
53 originated in the Caspian Sea and Black Sea region (Quero-García et al. 2019).
54 Sour cherries commercial cultivation is mainly localized in Eastern and Central Eu-
55 rope, North America, and Central and Western Asia on an area of 217,960 ha. They
56 are mainly cultivated for processing jams, juices, and whole fruits in preserved or dried
57 forms. Furthermore, they are used in dairy products and baked goods. Their global
58 production in 2021 was 1.51M tons of fruit and the gross production value in 2020
59 was 1.2 billion \$US (<https://www.fao.org/faostat/en/#data>). The variability of morpho-
60 logical and fruit characteristics is very high in sour cherry, including fruit and juice
61 colour, size, firmness, and fruit compounds (Schuster et al. 2017). Furthermore, tree
62 growth varies from rather slender and upright to small and bushy types. Variation
63 within ecotypes, differing e.g. in cold tolerance or growth habit, have been selected
64 across Europe over time (Dirlewanger et al. 2007; Hancock 2008). However, just a
65 small number of cultivars actually dominates the cultivation of sour cherry.
66 'Schattenmorelle' is the dominant cultivar (cv) in Middle Europe (Figure 1), whereas
67 sour cherry production in the United States is still based on 'Montmorency' (Quero-
68 Garcia et al. 2019).

69 'Schattenmorelle' was first described in France and today it is known in many coun-
70 tries with different names. In Poland, for example, it is called 'Łutovka' and in France
71 'Griotte du Nord' or 'Griotte Noir Tardive'.

72 The sour cherry is an allotetraploid with $2n=4x=32$ chromosomes. It originated as a
73 hybrid of an unreduced $2n$ pollen grain of *P. avium* ($2n=2x=16$) and a $1n$ egg cell of
74 the tetraploid ground cherry *P. fruticosa* ($2n=4x=32$) (Kobel 1927; Olden and Nybom
75 1968). Evidence of hybridization events between sweet and ground cherries has al-
76 ready been found several times in areas where both species occur simultaneously
77 (Macková et al. 2018, Hrotkó et al. 2020). The resulting hybrids were usually triploid

78 and were assigned to the secondary species *P. ×mohacsyana* Kárpáti. Natural occurrences of tetraploid sour cherries can be found in Eastern Turkey and the Caucasus region. There, they grow in forests and are used as wild forms for fruit production. The real area of origin is not known so far. Although *P. cerasus* can also be found in the wild in Europe, it is rather unlikely that those sour cherries are spontaneous hybrids. Since sour cherries are cultivated almost in many areas of the Northern hemisphere, they are often rather allochthonous individuals. The origin of the sour cherry thus seems to be based on a few hybridisation events. The results obtained by Olden and Nybom (1968) in experiments on the resynthesis of the species *P. cerasus* confirmed this hypothesis. The progeny from crosses between sweet and ground cherry showed the characteristic phenotype of the sour cherry. Studies based on chloroplast DNA markers strongly suggest that hybridisation between *P. avium* and *P. fruticosa* led to the emergence of *P. cerasus* at least twice (Dirlewanger et al. 2007). Furthermore, the hypothesis could also be confirmed by genomic *in situ* hybridisation (Schuster and Schreiber 2000) and transcriptome sequencing (Bird et al. 2022).

93 The sour cherry genome is presumed to have a genome size of 599 Mbp (Dirlewanger et al. 2007). It consists of the two subgenomes, each with eight chromosomes in the haploid set of chromosomes. One subgenome comes from the sweet cherry (*Pce_a*), while the other is from the ground cherry (*Pce_f*). However, the genome does not appear to be completely allopolyploid, since it has long been suspected that parts of the sour cherry genome might be segmental (Beaver and Lezzoni 1993; Olden and Nybom 1968; Raptopoulus 1941; Schuster and Wolfram 2005). The impact of hybridization and polyploidization between sweet and ground cherry has not been investigated so far, nor when this event originates. Cai et al. (2018) assumes that a mix of multi- and bivalent pairing led to imbalance segregation of chromosomes during meiosis in sour cherry, which in turn suggests that the sour cherry genome has not yet stabilized (Mason et al. 2020), or is still in the process of doing so. Recent advances in genome sequencing shed light into the complex structure and shape of polyploid genomes and their evolution (Zhang et al. 2021, Edger et al. 2019, Bertioli et al. 2019, Wu et al. 2021, Wang et al. 2019).

108 Here we report a high quality pseudo-chromosome-level genome assembly of the tetraploid sour cherry 'Schattenmorelle' (hereinafter referred to as *Pce_S*) generated

110 with a combination of Illumina NovaSeq short-read and Oxford Nanopore long-read
111 sequencing technology. The sequences were scaffolded to chromosomes by Hi-C. In
112 parallel, a full-length transcriptome of 'Schattenmorelle' was generated with the Pac-
113 Bio Sequel II SMRT cell long-read technology. Comparative sequence and amino acid
114 analyses between data sets of *Prunus avium* cv 'Tieton' (hereinafter referred to as *Pa_T*)
115 and *Prunus fruticosa* ecotype Hármashtárhegy (hereinafter referred to as *Pf_{EH}*) as
116 representatives of the two ancestral species (Wang et al. 2019, Wöhner et al. 2021)
117 and the two subgenomes of 'Schattenmorelle' *Pce_S_a* and *Pce_S_f* shed light into the
118 evolution of sour cherry. Finally, possible HE within the sub-genomic structure of sour
119 cherry were spotted, explaining the sour cherry genome segmental allopolyploidy.

120 **Results**

121

122 *De novo assembly and scaffolding*

123

124 A total of 68 Gb of paired-end Illumina sequencing data was obtained, corresponding
125 to ~114x coverage of the estimated genome size of 599 Mbp. Using two PromethION
126 flow cells, a total of 178 Gb was produced (~300x coverage). The longest ONT reads
127 that together resulted in a 20x coverage were selected for assembly, having a
128 minimum read length of 64,214 bp. Table S1 summarizes the properties of the 20-
129 WGS-PCE.1.0 assembly after polishing. The *Prunus avium* and *Prunus fruticosa*
130 contigs were then separated successfully by read mapping and contig selection that
131 fit the hypothesis of 1 or more clear coverage peaks from the 20-WGS-PCE.1.0
132 assembly. The resulting two datasets, representing the subgenomes *Pce_s_a* and
133 *Pce_s_f*, were purged and used for scaffolding using HI-C. After manual curation of
134 both datasets, the final consensus genome assembly was scaffolded from 935 and
135 865 contigs of the *Pce_s_a* and *Pce_s_f* subgenomes, respectively. Eight clusters ideally
136 representing the eight chromosomes were obtained for each subgenome (Figure S1).
137 The final genome sequence is 628.5 Mbp long and consists of eight chromosomes
138 for each subgenome (Figure 2). A total of 269 Mbp were assigned to subgenome
139 *Pce_s_a* (N50 of 31.5 Mbp) and 299.5 Mbp (N50 of 39.4 Mbp) to *Pce_s_f*. Eighty-six and
140 134 unassembled contigs were unassigned to chromosomes for *Pce_s_a* (22.7 Mbp)
141 and *Pce_s_f* (37.3 Mbp), respectively.

142 The longest scaffold from *Pce_s_a* is 52.8 Mbp and 53.5 Mbp from subgenome *Pce_s_f*
143 (Table S2). Except for chromosome five, all scaffolds obtained from subgenome
144 *Pce_s_f* are longer compared to the corresponding chromosome of subgenome
145 *Pce_s_a*. The chloroplast sequence obtained was 158,178 bp and the mitochondrial
146 sequence was 343,516 bp long (Figure S2).

147

148 *Transcriptome sequencing, isoseq-analysis, structural and functional annotation*

149

150 The total repeat content of the entire *Pce_s* genome sequence was 49.7%. The total
151 repeat content of subgenome *Pce_s_a* was 48.3% and of subgenome *Pce_s_f* 50.9%,

152 respectively (Table 1). The class I elements Gypsy comprised the largest fractions of
153 repetitive elements in the *Pce_S* genome sequence. A quantitative reduction between
154 elements of this family was also detected in the *Pce_S_f* subgenome with a difference
155 of 10.7%. Several elements could only be detected in one genotype of the two
156 ancestral species. The TAD1 class I element only occurred in *Pf_{eH}*, while class II, order
157 TIR - IS3EU, P, and Sola-3 were specifically detected in the genome of *Pa_T*. No
158 element was found, which was only present in one of the two subgenomes of *Pce_S*.
159 Several elements occurred in both subgenomes (class I, LINE – R1-LOA, RTE-X, SINE
160 – tRNA-DEU- L2, class II, TIR – TcMar- Mariner, and DADA elements) but were not
161 detected in *Pf_{eH}* and *Pa_T*. The class I elements of the order LTR (ERV1, Pao) and
162 Academ/-2 were only detected in one of the two genomes representing the ancestral
163 species and in *Pce_S_a* and *Pce_S_f*. Iso-Seq results are summarized in Table S3. In
164 total 248,218 high quality isoforms have been identified. Both, the high and the low
165 quality isoforms have been used for genome annotation where each gene might be
166 represented by multiple isoforms. A total of 107,508 transcripts (*Pce_S_a*: 53,497;
167 *Pce_S_f*: 54,011) were predicted from the 60,123 gene models (*Pce_S_a*: 29,069; *Pce_S_f*:
168 31,054) obtained by structural annotation procedures (Table S4). Interproscan
169 analysis detected 1,381,841 functional annotations (*Pce_S_a*: 649,310; *Pce_S_f*:
170 687,531) using 16 databases. Two third (71,870) of the transcripts were assigned with
171 GO terms and 9,114 were found to be involved in annotated pathways.

172

173 *Completeness and quality of the genome and transcriptome*

174

175 BUSCO completeness of the *Pce_S* genome was 99.0% (S: 16.7%, D: 82.3%, F: 0.4%,
176 M: 0.6%, n: 1,614) respectively and comparable with *P. persica* 'Lovell' (99.3%) and
177 *P. avium* 'Tieton' (98.3%, Figure S3). Completeness of subgenome *Pce_S_a* was higher
178 (C: 89.4%, S: 84.8%, D: 4.6%, F: 1.5%, M: 9.1%, n: 1614) compared to subgenome
179 *Pce_S_f* (C: 87.1%, S: 80.9%, D: 6.2%, F: 1.2%, M: 11.7%, n: 1,614). The calculated
180 LAI index was 6.3 and low in comparison to other genomes (*Pp_L*: 17.6, *Pa_T*: 10.3, *Pf_{eH}*:
181 13.1). The LAI index for subgenome *Pce_S_a* was 7.1. The LAI index for subgenome
182 *Pce_S_f* was 5.6 (Figure S4). The nucleotide heterozygosity rates were 94.9% for aaaa,
183 2.39 for aaab, 2.4 for aabb; 0.001 for aabc and 0.308 for abcd (Figure 3). The

184 comparison of genetic position and physical position of up to 1,856 markers of the
185 five genetic sour cherry maps (Table S5a) showed a good co-linearity to the genome
186 sequence (Figure S5). Busco evaluation on completeness of the annotated proteins
187 resulted in 99.2% [C: 99.2% (S: 8.4%, D: 90.8%), F: 0.4%, M: 0.4%, n: 1,614]. The
188 chloroplast sequence obtained contained 427 genes, 21 rRNAs and 136 tRNAs,
189 whereas the mitochondrial sequence contained 188 genes, 3 rRNAs and 152 tRNAs
190 (Figure S2). An ab-initio and homology-based gene prediction with 14 reference
191 species was performed (IAA). Based on the homology prediction, thirty-four percent
192 of the proteins showed the highest IAA towards *Prunus fruticosa* and 17.9% towards
193 *P. avium*. Only 5.2% of the proteins showed no IAA to any of the used reference
194 datasets used, which was due to ab-initio prediction. The data is summarized in
195 Figure S6.

196

197 A comparison of transcripts of *Pces* and the annotation datasets of *PfeH* and *PaT*
198 enabled a quantitative comparison of shared transcripts within the datasets (Table 2).
199 A total number of 26,532 shared transcripts were found between the two subgenomes
200 *Pces_a* and *Pces_f* and the genomes of *PfeH* and *PaT*. Thirty-eight percent of the *P.*
201 *cerasus* proteins had a greater IAA to *PfeH*, whereas 54% showed a greater IAA to *PaT*.
202 Eight percent showed an identical IAA to both ancestral species. A larger number of
203 transcripts of both sour cherry subgenomes (Table 2) were assigned to the annotation
204 data set of *PfeH*. A total of 13,425 transcripts from the *Pces_a* subgenome and 13,107
205 from the *Pces_f* subgenome were found in the genome sequences of *PfeH* and *PaT*.
206 Seventy-five percent of the pool from the *Pces_a* subgenome showed a higher IAA to
207 *PaT* and 17% to *PfeH*, while 59% from the pool originating from the *Pces_f* subgenome
208 showed a higher IAA to *PfeH* and 32% to *PaT*.

209

210 *Identification of syntenic regions and inversions*

211

212 The sequences of the two subgenomes *Pces_a* and *Pces_f* and the genotypes *PaT*
213 and *PfeH* of the two ancestral species *P. avium* and *P. fruticosa* were screened for
214 duplicated regions using DAGchainer as previously published for peach (Verde et al.
215 2012). The seven major triplicated regions were found nearly one to one in *P. avium*

216 but not in *P. fruticosa*, which lacked regions 4 and 7 corresponding to Verde et al.
217 (2012). *P. avium* and *P. fruticosa* seem to derive from the same paleohexaploid event
218 like peach, but with a loss of the fourth and seventh paleoset of paralogs in *P.*
219 *fruticosa*. The graphical analysis is summarized in figure S7).

220

221 Thirteen inversions were detected through positional co-linearity comparison between
222 the two subgenomes using the molecular markers from the 9+6k SNP array (Figure
223 S8). Five inversions were found between subgenome *Pce_S_a* and the genome
224 sequence of *Pa_T*. Eleven inversions were found between subgenome *Pce_S_f* and *Pf_{eH}*
225 (Table S6). By comparing the position of amino acid sequences of orthologous
226 proteins (synteny), we found 21 inversions when comparing *Pce_S_f* with *Pf_{eH}*. Only
227 seven were found between *Pce_S_a* and *Pa_T* and 16 between both subgenomes *Pce_S_a*
228 and *Pce_S_f* (Figure S9).

229

230 *Detection of de novo homoeologous exchanges*

231 For the detection of de novo homoeologous exchanges we used three approaches by
232 comparing inter- and intraspecific %-covered bases (genomic and transcriptomic)
233 and %-IAA between proteins of *Pce_S* to *Pa_T* and *Pf_{eH}* (Figure 4, Figure S10, S11). *Pce_S*
234 short-reads were mapped against *Pa_T* and *Pf_{eH}* and only species specific reads (*Pa_T*
235 and *Pf_{eH}*) were filtered into read-subsets. The obtained read-subsets were re-mapped
236 against *Pce_S_a* and *Pce_S_f* and base coverage was calculated. A total of 1024 regions
237 (100k window) were intraspecific %-covered bases from mapped reads (*Pce_S_a* to
238 *Pa_T*, *Pce_S_f* to *Pf_{eH}*) was < than interspecific %-covered bases from mapped reads
239 (*Pce_S_a* to *Pf_{eH}*, *Pce_S_f* to *Pa_T*) were discovered. In a second approach, translocations
240 between the two subgenomes were localized by a short-read mapping analyses.
241 Short-reads (RNAseq) from *P. cerasus*, *P. avium* and *P. fruticosa* obtained from Bird
242 et al. (2022) were mapped on *Pce_S*. A total of 148 regions were intraspecific difference
243 of %-covered bases from obtained RNAseq reads (*Pa* and *Pce_S_a*, *Pf* and *Pce_S_f*) >
244 than interspecific difference of %-covered bases from obtained RNAseq reads (*Pf* and
245 *Pce_S_a*, *Pa* and *Pce_S_f*) indicated homoeologous exchanges between the two
246 subgenomes. Finally, 367 regions were the proportion of transcripts with intraspecific
247 amino acid identity (*Pa_T* and *Pce_S_a*, *Pf_{eH}* and *Pce_S_f*) < than the proportion of

248 transcripts with interspecific amino acid identity (*Pf_{eh}* and *Pce_s_a*, *Pa_T* and *Pce_s_f*)
249 were identified (Figure 4). Several regions were confirmed by calculating the 70%
250 quantile of the IAA value within a window of 1 Mbp windows (Note S2). This confirms
251 that there are transcripts in the *Pce_s* genome whose IAA to the homoeologous
252 representative genome (*Pf_{eh}*) is greater than to the homologous (*Pa_T*) representative.
253 A total of 16 in *Pce_s_a* and 12 in *Pce_s_f* regions spanning 49 250k-windows were
254 finally identified that match all three criteria indicating *denovo* homoelogous
255 exchanges within the subgenomes (Figure 4, Figure S10). No evidence for an
256 introgression of other *Prunus* species was found (Note S2, Note S3).

257

258 *LTR dating and divergence of time estimation.*

259

260 Left and right LTR identity of a subset of 2,385 (*Pce_s_a*), 3,028 (*Pce_s_f*), 3,130 (*Pa_T*)
261 and 3,992 (*Pf_{eh}*) LTRs was analysed. The homologous genomes shared 200 (*Pce_s_a*
262 versus *Pa_T*) and 100 (*Pce_s_f* versus *Pf_{eh}*) LTRs whereas 12 LTRs were shared by
263 *Pce_s_a* versus *Pf_{eh}* and *Pce_s_f* versus *Pa_T*. Only five common LTRs were found
264 between *Pce_s_a* and *Pce_s_f* and 13 between *Pa_T* and *Pf_{eh}*. A summary of the LTRs
265 insertion time is shown in figure 5A.

266

267 The youngest shared LTRs between *Pce_s_a* and *Pa_T* were calculated with 103,896.1
268 generations and between *Pce_s_f* and *Pf_{eh}* with 97,402.6 generations. When comparing
269 the homeologous chromosomes, the youngest shared LTRs between *Pa_T* and *Pf_{eh}*
270 was calculated with 116,883.1 generations, between *Pce_s_a* and *Pce_s_f* with
271 194,805.2 generations. LTRs of *Pa_T* were also found in subgenome *Pce_s_f* and
272 calculated with 207,792.2 generations. LTRs of *Pf_{eh}* were detected in subgenome of
273 *Pce_s_a* and calculated with 149,350.6 generations. This indicates an exchange of
274 LTRs between the two subgenomes.

275

276 A total of 834 single copy orthogroups among nine genomes were found and used for
277 single protein alignments. Single alignments were concatenated and a final alignment
278 with nine amino acid sequences representing each species with 419,586 amino acid
279 positions was used for phylogenetic tree construction. Using the RelTime method, the

280 estimated divergence time between the genera *Malus* and *Prunus* was 50.4 Mya. The
281 species groups *P. persica* and *P. mume* diverged from the *P.*
282 *yedonensis*/*P. avium*/*P. fruticosa* group 11.6 Mya. Based on this model, the divergence
283 of the two subgenomes of *P. cerasus* compared to the genome sequences of *Pa_T* and
284 *Pf_{EH}* was estimated with 2.93 Mya and 5.5 Mya respectively (Figure 5B).

285 **Discussion**

286 The genome of the economically most important sour cherry 'Schattenmorelle' in
287 Europe was sequenced using a combination of Oxford Nanopore R9.4.1 PromethION
288 long read technology and Illumina NovaSeq™ short read technology. After
289 assignment of the long-reads to the two subgenomes and Hi-C analysis, the final
290 assembly was 629 Mbp.

291 This sequence was used to study structural changes present in the allotetraploid sour
292 cherry genome after its emergence. Therefore, the sour cherry genome sequence was
293 compared to the published genome sequences of *Prunus avium* 'Tieton' (*Pa_T*, Wang
294 et al. 2019) and *Prunus fruticosa* ecotype Hármashatárhegy (*Pf_{EH}*, Wöhner et al. 2021a)
295 representing genotypes of the two ancestral species. The size of the subgenome *Pce_S*
296 originating from *P. avium* was 269 Mbp. A similar genome size (271 Mbp) is described
297 for the *Prunus avium* 'Big Star' (Pinosio et al. 2020) and 'Sato Nishiki' (Shirarsawa et
298 al. 2017). Larger differences were found to 'Regina' with 279 Mbp and 'Tieton' with
299 344 Mbp (Wang et al. 2020, Le Dantec et al. 2019). Differences were also found
300 between the size of subgenome *Pce_S_f* (299 Mbp) and the genome of the ground
301 cherry genotype *P. fruticosa* ecotype Hármashatárhegy (366 Mbp, Wöhner et al.
302 2021a). These differences indicate a reduction of the subgenome *Pce_S_a* by 0.49%-
303 21%, whereas for subgenome *Pce_S_f* a reduction of 18.29% was found. The reduction
304 in genome size for allotetraploid species in comparison to their ancestral genomes is
305 common, as already reported for *Nicotiana tabacum* (1.9%-14.3%) and *Gossypium*
306 species (Leitch et al., 2008; Hawkins et al., 2009; Renny-Byfield et al., 2011), with an
307 overall downsizing rate for angiosperms calculated as 0%-30% (Zenil-Ferguson et al.
308 2016). Genome downsizing in response to a genome hybridisation event can be
309 explained with evolutionary advantages, which give these species with smaller

310 genomes a selection advantage in the long-term (Knight et al. 2005, Zenil-Ferguson
311 et al. 2016).

312 Although downsizing of the *P. cerasus* subgenomes is most probable, enlargement
313 and expansion of the genomes of ancestral species during evolution would be another
314 possibility. However, an increase in genome size during the evolution of a species has
315 only rarely been documented (Leitch et al. 2008, Jackobs et al. 2004, Kim et al. 2014).
316 BUSCO analysis provides additional evidence for the reduction of genome size.
317 Although the number of genes does not correlate with genome size in eukaryotes
318 (Pierce 2012), differences between the ancestral genomes and sour cherry could be
319 observed when looking at BUSCO-completeness. Considering both subgenomes and
320 the genomes of *Pa_T* and *Pf_{eH}*, a completeness of > 96.4% was obtained. However, the
321 completeness of the single subgenomes was only 89.4% for *Pce_S_a* and 87.1% for
322 *Pce_S_f*. Structural differences between the *P. cerasus* subgenomes and the genomes
323 of *Pa_T* and *Pf_{eH}* were also found by comparing the number of repetitive elements. While
324 the content of repetitive elements differs by only 0.86% between *Pf_{eH}* and *Pce_S_f*, it is
325 17.8% between *Pa_T* and *Pce_S_a*. Whether this is a consequence of hybridisation
326 remains speculative and would deserve further studies. An increase of class I
327 elements Gypsy from six percent in the *Pa_T* genome to 7.3% in subgenome *Pce_S_a*
328 indicates an expansion of this class following the formation of the sour cherry genome
329 or a possible reduction of non-repetitive sequences in the corresponding subgenome
330 resulting in a smaller genome size.

331

332 A comparison of syntetic regions showed a high degree of collinearity between *Pce_S*,
333 *Pa_T* and *Pf_{eH}* genomes (Figure 2, Figure S7-9), with single inversions between the
334 respective chromosome pairs. Using the genome of *P. persica*, seven triplicated
335 regions were detected in *Pce_S_a*, confirming that these genomes descend from a
336 palaeohexploid ancestor. However, the triplicated regions 4 and 7 in *Pce_S_f* were only
337 detected in highly fragmented form or have been lost. Hao et al. (2022), who described
338 a rapid loss of homoeologs immediately after polyploidy events, described similar
339 finding.

340

341 Based on Ranallo-Benavidez (2020) the results from the k-mer analysis confirm that
342 the genome of sour cherry can be considered as highly heterozygous and segmental
343 allotetraploid. Furthermore, genomes of segmental allopolyploids may possess a mix
344 of auto- and allopolyploid segments through duplication-deletion events as a result of
345 homoeologous exchanges leading to either reciprocal translocations or
346 homoeologous non-reciprocal translocations (Mason and Wendel, 2020). Whereas
347 autotetraploids have an aaab > aabb rate, allotetraploids are considered to have aaab
348 < aabb. The near identical rate between aaab and aabb in *Pce_S* provides strong
349 evidence that the sour cherry is a segmental mix of auto- and allotetraploidy.

350

351 Due to this assumption of segmental allotetraploidy, homoeologous recombination
352 between homoeologous chromosomes is very likely. This is confirmed by the
353 coverage and amino acid identity analyses. Homoeologous exchange events between
354 the chromosomes of subgenome *Pce_S_a* and *Pce_S_f* were detected (Figure 4, Figure
355 S10). These exchanges are not balanced but probably a product of a
356 duplication/deletion event as described by Mason et al. (2020), generating the
357 proposed mosaic of genomic regions representing one or the other subgenome.
358 Additional mapping of transcriptomes from six other *Prunus* species indicates that no
359 major introgression from one of these species occurred. Using 14 reference species,
360 60,123 gene models were annotated. Almost the same number was assigned to the
361 two *P. cerasus* subgenomes (Table 2). No evidence was found for large introgressions
362 from any of the reference species (Note S2, S3). By comparing the amino acid identity
363 of the proteins of *Pa_T* and *Pf_{eH}* with the respective sour cherry subgenome, the
364 identified translocations via read mapping could be confirmed. The majority of the
365 transcripts (51%) could be assigned to the genotypes *Pa_T* and *Pf_{eH}* of the two ancestral
366 species *P. avium* and *P. fruticosa* (Figure S6). 5.2% of the transcripts could not be
367 assigned to any of the reference species. Only < 1% of the transcripts could not be
368 assigned to one of the ancestral species. They showed equivalent matches to both
369 species and are probably a product of ab-initio prediction. A total of 49,698 proteins
370 in subgenome *Pce_S_a* and 48,576 proteins in *Pce_S_f* shared only 13,435 and 13,107
371 proteins with *Pa_T* and *Pf_{eH}*, respectively. A total of 75% of the proteins of subgenome
372 *Pce_S_a* matched better to *Pa_T* compared to *Pf_{eH}*, whereas only 59% of *Pce_S_f* mapped

373 better to Pf_{eH} than to Pa_T (Table 2). Subgenome Pce_{S_a} seems to be closer to Pa_T than
374 Pce_{S_f} to Pf_{eH} . This was confirmed by an evolutionary approach that calculated the
375 separation of the subgenome Pce_{S_f} from *P. fruticosa* 5.5 Mya, while subgenome
376 Pce_{S_a} separated from *P. avium* 2.93 Mya (Figure 5B).

377 To validate these results, the insertion events of long terminal repeats between *P.*
378 *avium*, *P. fruticosa* and the subgenomes were calculated. Assuming a *Prunus*-specific
379 rate of 7.7×10^{-9} mutations per generation (Xie et al. 2016), LTRs of the same type with
380 the same insertion time were identified in the same positional order in the different
381 (sub)genomes. The most recent co-occurring LTRs between the genomes of Pa_T and
382 Pf_{eH} could be dated at 116,883.1 generations. Exact data on the duration of the
383 generation time of *Prunus* species in natural habitats do not exist. Although the
384 juvenile phase of many *Prunus* species is usually completed after 5 years (Besford et
385 al. 1996), it can be assumed that the times for a generation are considerably higher.
386 Many fruit species are hardly or not at all able to rejuvenate by seeds under natural
387 conditions (Coart et al. 2003), or they rejuvenate mainly by root suckers (Li et al. 2022).
388 Other studies on *Prunus* therefore assume a duration of 10 years per generation
389 (Wang et al. 2021), although even this seems rather too little.

390 Assuming that a generation change is to be expected after 10 to 60 years (Besford et
391 al. 1996), this would correspond to a time period of ~1 to 6 Mya. The youngest co-
392 occurring LTR could be estimated at 1.9 Mya. This suggests that *P. fruticosa* and *P.*
393 *avium* probably shared a gene pool between ~1 Mya and 2 Mya. It should be noted
394 that this estimate can vary greatly depending on the number of years per generation
395 used in the calculation (Figure 5A). Based on the results of the protein dating a
396 generation time of 30 years is more likely for *P. avium*. For *P. fruticosa*, which occurs
397 less frequently in natural habitats and reproduces mainly via root suckers, the
398 generation time seems to be somewhat longer at 55 to 60 years. Some LTRs present
399 in Pa_T , but absent in Pf_{eH} were found in subgenome Pce_{S_f} only and vice versa. Other
400 class I elements (LTR - ERV1, Pao) and Academ/-2 specifically detected in one of the
401 two genotypes Pa_T and Pf_{eH} representing the two ancestral species of sour cherry and
402 in Pce_{S_a} and Pce_{S_f} , which indicates a transfer of these elements between the two
403 subgenomes following the formation of the allotetraploid *P. cerasus* genome. This is
404 a further indication for a segmental exchange between the two sour cherry

405 subgenomes. Mason et al. (2020) speculated that uni-directional homeologous
406 exchange was observed in recent or synthetic allopolyploids. However, our results
407 confirm this hypothesis by the evidence that sour cherry is a recent allopolyploid with
408 autoploid segments derived from uni-directional homoeologous exchanges.

409

410 **Conclusion**

411 After sequencing of the genome of the sour cherry 'Schattenmorelle', the following
412 can be concluded. The genome of sour cherry is segmental allotetraploid. It consists
413 of two subgenomes, one derived from the sweet cherry *P. avium* and one from the
414 ground cherry *P. fruticosa*. DNA sequences have been repeatedly exchanged
415 between the two subgenomes. At the same time, a reduction in genome size has
416 taken place. Other *Prunus* species have not contributed to the evolution of this
417 species. No evidence was found for introgressions in the sour cherry genome derived
418 from *Prunus* species other than *P. avium* and *P. fruticosa*. Sour cherry is a very young
419 *Prunus* species. The origin of this species is estimated about 1 mya earliest.

420

421 **Material and Methods**

422

423 *Plant Material, DNA and RNA extraction, sequencing and iso-seq analysis*

424

425 Snap frozen *Prunus cerasus* L. 'Schattenmorelle' (accession KIZC99-2, Figure 1,
426 supplements 1.1) young leaf material was sent to KeyGene N.V. (Wageningen, The
427 Netherlands). High molecular weight extracted DNA (Wöhner et al. 2021a) was used
428 to generate 1D ligation (SQK-LSK109) libraries which were subsequently sequenced
429 on two Oxford Nanopore Technologies (ONT) R9.4.1 PromethION flow cells. The
430 same material was used to generate an Illumina PCR free paired-end library (insert
431 size of ~550 bp) which was sequenced on a HiSeq 4000™ platform using 150bp and
432 125bp paired-end sequencing.

433

434 Snap frozen tissues from buds, flowers, leaves and fruits were collected and in the
435 field and total RNA was extracted with Maxwell® RSC Plant RNA Kit (Promega). Two
436 pools were generated and used for PacBio Iso-Seq library preparation (Procedure &

437 Checklist – Iso-Seq™ Express Template Preparation for Sequel® and Sequel II
438 Systems, PN 101-763-800 Version 02). Each library pool was sequenced on a single
439 8M ZMW PacBio Sequel II SMRT cell (supplements 1.1). Obtained full-length reads
440 with 5'-end primer, the 3'-end primer and the poly-A tail were filtered and these
441 sequences were trimmed off. Transcripts containing (artificial) concatemers were
442 completely discarded. Isoforms (consensus sequence) generated by full-length read
443 clustering (based on sequence similarity), were finally polished with non-full-length
444 reads using Arrow (SMRT Link v7.0.0, https://www.pacb.com/wp-content/uploads/SMRT_Tools_Reference_Guide_v600.pdf).
445

446

447 *De novo assembly and scaffolding*

448 The aligner Minimap2 (v2.16-r922, Li et al. 2018) and assembler Miniasm (v0.2-r137-
449 dirty, Li et al. 2016) were used for raw data assembly generation. Racon (vv1.4.10,
450 Vaser et al. 2017) and Pilon (v1.22, Walker et al. 2014) were used for base-quality
451 improvement with raw ONT and Illumina read data. Chromosome-scale scaffolding
452 was performed by Phase Genomics (Seattle, Washington, USA) with Proximo Hi-C
453 (supplements 1.2). The resulting assembly was designated as 20-WGS-
454 PCE_<Avium|Fruticosa>2.0 _<Contig|Scaffold>.

455

456 *Correctness, completeness and contiguity of the Prunus cerasus genome sequence*
457 The BUSCO (Benchmark Universal Single-Copy Orthologs - Galaxy Version 4.1.4)
458 software was used for quantitative and quality assessment of the genome assemblies
459 based on near-universal single-copy orthologs. The long terminal repeat (LTR)
460 assembly Index (LAI) Ou et al. (2018) was calculated with LTR_retriever 2.9.0
461 (https://github.com/oushujun/LTR_retriever) to evaluate the assembly continuity
462 between the final genome sequence of *P. cerasus* 'Schattenmorelle' and *Prunus*
463 *fruticosa* ecotype Hármashatárhegy (Wöhner et al. 2021a, Pf_1.0), *P. avium* 'Tieton'
464 (Wang et al. 2020) and *P. persica* 'Lovell' (Verde et al. 2017) respectively. LTR_harvest
465 (genometools 1.6.1 implementation) was used to obtain LTR-RT candidates. The
466 genome size was also estimated by k-mer analysis (supplements 1.3) using Illumina
467 short read data. The merged datasets were subsequently used to generate a
468 histogram dataset representing the k-mers from all datasets. GenomeScope (Galaxy

469 Version 2.0, Ranallo-Benavidez 2020) was used to generate a histogram plot of k-mer
470 frequency of different coverage depths using the tetraploid ploidy level (k-mer length
471 19). Marker sequences and genetic positions from five available genetic sour cherry
472 maps (M172x25-F1, US-F1, 25x25-F1, Montx25-F1, RE-F1) and 14,644 SNP markers
473 (9+6k array) were downloaded from the Genome Database for Rosaceae (GDR,
474 <https://www.rosaceae.org/>). The marker sequences were mapped on the
475 chromosome sequences using the mapping software bowtie2 (Galaxy Version
476 2.5.0+galaxy0, Langmead et al. 2012) implementation on the Galaxy server
477 (<https://usegalaxy.org>) with standard settings.

478

479 *Structural and functional annotation*

480

481 For an inter species repeat comparison, a species-specific repeat library was
482 generated with RepeatModeler open-1.0.11, and the genome was subsequently
483 masked with RepeatMasker open-4.0.7. For structural genome annotation, another
484 species-specific repeat library for PCE_1.0 was generated with RepeatModeler2
485 (Flynn et al., 2020) version 2.0.2, and the genome was subsequently masked with
486 RepeatMasker 4.1.2. (Further details on the repeat masking software configuration
487 are available in supplements 1.4.1.).

488 To generate extrinsic evidence for structural annotation of protein coding genes, short
489 read RNA-Seq library SRR2290965 was aligned to the genome using HiSat2 version
490 2.1.0 (Kim et al., 2019). The output SAM file was converted to BAM format using
491 SAMtools (Li et al, 2009). The resulting alignment file was further used by both
492 BRAKER1 (Hoff et al., 2016; Hoff et al., 2019) and GeMoMa (Keilwagen et al. 2016).
493 Further, a custom protocol was used for integrating long read RNA-Seq data into
494 genome annotation (supplements 1.4.2). In short, protein coding genes were called in
495 Cupcake transcripts using GeneMarkS-T (Tang et al., 2015) and these predictions
496 were converted to hints for BRAKER1. In addition, intron coverage information from
497 long read to genome spliced alignment with Minimap2 (Li et al. 2018) was provided
498 to BRAKER1.

499 A combination of BRAKER1 (Hoff et al., 2016, Hoff et al., 2019), BRAKER2 (Bruna et
500 al., 2021), TSEBRA (Gabriel et al., 2021), and GeMoMa (Keilwagen et al. 2016) was

501 used for the final annotation of protein coding genes. BRAKER pipelines use a
502 combination of evidence-supported self-training GeneMark-ET/EP (Lomsadze et al.,
503 2014; Bruna et al., 2020) (here version 4.68) to generate a training gene set for the
504 gene prediction tool AUGUSTUS (Stanke et al., 2008; here version 3.3.2). BRAKER1
505 version 2.1.6 was here provided with BAM-files of from short and long read RNA-Seq
506 to genome alignments, and with gene structure information derived from Cupcake
507 transcripts using GeneMarkS-T. This generated a gene set that consists of ab initio
508 and evidence supported predictions. A separate gene set was generated with
509 BRAKER2, which uses protein to generate a gene set. We used the OrthoDB version
510 10 (Kriventseva, 2019) portion of plants in combination with the full protein sets of
511 *Prunus fruticosa* (Wöhner et al., 2021b), *Prunus armeniaca* (GCA 903112645.1),
512 *Prunus avium* (GCF_002207925.1), *Prunus dulcis* (GCF_902201215.1), *Prunus mume*
513 (GCF_000346735.1), *Prunus persica* (GCF_000346465.2) as input for BRAKER2. Both
514 the BRAKER1 and BRAKER2 AUGUSTUS gene sets were combined with a
515 GeneMarkS-T derived gene set using TSEBRA (Gabriel et al., 2021) from the
516 long_reads branch on GitHub with a custom configuration file (supplements 1.4.3.)
517 incorporating evidence from BRAKER1 and BRAKER2.
518 GeMoMa was run on the genome assembly of Schattenmorelle using 14 reference
519 species and experimental transcript evidence (supplements 1.4.4). GeMoMa gene
520 predictions of each reference species were combined with TSEBRA predictions using
521 the GeMoMa module GAF and subsequently, UTRs were predicted in a two-step
522 process based on mapped Iso-seq and RNA-seq data using the GeMoMa module
523 AnnotationFinalizer (supplements 1.4.5). First, UTRs were predicted based on Iso-seq
524 data. Second, UTRs were predicted based on RNA-seq data for gene predictions
525 without UTR prediction from the first step. An assembly hub for visualization of the
526 *Prunus cerasus* genome with structural annotation was generated using MakeHub
527 (Hoff, 2019; supplements 2). The functional annotation was performed with the Galaxy
528 Europe implementation of InterProScan (Galaxy Version 5.59-91.0+galaxy3, Jones et
529 al. 2014, Cock et al 2013, Zdobnov et al. 2001, Quevillon et al. 2005, Hunter et al.
530 2009). The chloroplast and mitochondria sequences were annotated with GeSeq
531 (Tillich et al. 2017, supplements 1.4.5).
532

533 *Identification of syntenic regions*

534

535 Structural comparison of orthologous loci between the subgenomes *Pce_S_a* and
536 *Pce_S_f* of *Prunus cerasus* and the two genotypes *Pa_T* and *Pf_{eH}* as representatives of
537 the two genome donor species *P. avium* and *P. fruticosa* was calculated with the final
538 annotations using SynMap2 (Haug-Baltzell et al. 2017) available at the CoGe platform
539 (<https://genomevolution.org/coge/>). Analysis on triplication events were performed
540 with standard settings and Last as Blast algorithm at a ratio coverage depth 3:3 in
541 SynMap2 (Haug-Baltzell et al. 2017).

542

543 *Identification of homoeologous exchange regions*

544

545 Homoeologous exchanges were identified on the amino acid, transcript and genomic
546 level.

547

548 *Calculation of amino acid identity*

549 Identity of amino acids (IAA) between all reference annotations homology-based gene
550 prediction was calculated by GeMoMa using the default parameters. Subsequently,
551 the *Pce_S* genome was divided into 250k windows, and the percentage of proteins
552 showing a higher IAA between *Pf_{eH}* (Wöhner et al. 2021b) and *Pa_T* (Wang et al. 2020)
553 to the respective subgenome (*Pce_S_a* and *Pce_S_f*) was determined. The percentage
554 of proteins in this window, which were more similar to *Pa_T* was finally subtracted from
555 the percentage of proteins which were more similar to *Pf_{eH}*. A proportion of transcripts
556 with higher intraspecific amino acid identity (between *Pa_T* and *Pce_S_a* or *Pf_{eH}* and
557 *Pce_S_f*) is expected compared to the proportion of transcripts with interspecific amino
558 acid identity match (*Pf_{eH}* and *Pce_S_a* or *Pa_T* and *Pce_S_f*). Opposite cases, indicate
559 potential translocations between the two subgenomes *Pce_S_a* and *Pce_S_f* and were
560 plotted into a circus plot (Figure S11).

561

562 *Read mapping and coverage analysis*

563 RNAseq raw data published by Bird et al. (2022) was obtained from NCBI sequence
564 read archive (SRA) for the following species: *P. cerasus* (SRX14816146,

565 SRX14816142, SRX14816138), *P. fruticosa* (SRX14816141), *P. avium* (SRX14816143),
566 *P. canescens* (SRX14816137), *P. serrulata* (SRX14816136), *P. mahaleb*
567 (SRX14816140), *P. pensylvanica* (SRX14816144), *P. maackii* (SRX14816139), *P.*
568 *subhirtella* (SRX14816145). Reads were adapter- and quality trimmed using the
569 software Trim Galore (version 0.6.3, parameters --quality 30 --length 50). Trimmed
570 reads were mapped against the *P. cerasus* subgenomes *Pce_s_a* and *Pce_s_f* using
571 STAR (version 2.7.8a, parameter --twopassMode Basic). The subsequent analysis
572 was performed in accordance to Keilwagen et al. (2022). The *Pce_s* genome was
573 divided into 250k windows. The percentage of covered bases using RNAseq data of
574 *P. cerasus* (SRX14816146, SRX14816142, SRX14816138) was estimated at a depth
575 of 1 for each window. The same was done with all other RNAseq data sets. The
576 percentage of covered bases from *P. avium* (SRX14816143) was subtracted from the
577 percentage of covered bases from *P. cerasus* (SRX14816146, SRX14816142,
578 SRX14816138). The same was done using the reads of *P. fruticosa* (SRX14816141).
579 For subgenome *Pce_s_a* it is expected that the intraspecific difference for transcripts
580 of data set *P. avium* (SRX14816143) is lower (close to 0) than the interspecific
581 difference for transcripts of data set *P. fruticosa* (SRX14816141) and vice versa.
582 Opposite cases indicate potential homoeologous exchanges between the two
583 subgenomes *Pce_s_a* and *Pce_s_f* and were plotted into a circos plot (Figure S11).
584 The nucleotide short reads from *Pce_s* were mapped against the genomes of the two
585 ancestral species *Pa_T* and *Pf_{eH}*. Subsequently, the mapped reads were filtered using
586 samtools for mapped reads in proper pair (-f 3) and primary alignments and not
587 supplementary alignment (-F 2304). Those reads were divided into four groups
588 according to the following criteria: 1. unique match to *Pa_T*; 2. unique match to *Pf_{eH}*; 3.
589 match to *Pa_T* and *Pf_{eH}*; 4. no match to *Pa_T* and *Pf_{eH}* (unique to *Pce_s*). The first two
590 separated read sets were then re-mapped against the subgenomes *Pce_s_a* and
591 *Pce_s_f*. The percentage of covered bases was calculated for a 100 k window. For the
592 subgenomes of *Pce_s*, the percentage of intraspecific covered bases (*Pce_s_a* to *Pa_T*,
593 *Pce_s_f* to *Pf_{eH}*) should be higher compared to the percentage of interspecific covered
594 bases (*Pce_s_a* to *Pf_{eH}*, *Pce_s_f* to *Pa_T*). The opposite case indicates possible
595 translocations and were plotted into a circos plot (Figure S11). Additionally, regions

596 of the Schattenmorelle genome assembly were determined that are uniquely covered
597 by *Pa_T* and *Pf_{eH}* filtered read sets.

598

599 *LTR insertion estimation*

600 The difference (identity) of left and right LTR was calculated using the script
601 EDTA_raw.pl from the software EDTA version 1.9 (<https://github.com/oushujun/EDTA>,
602 Ou et al. 2019). As input files we used the genome sequences of *P. cerasus* (*Pces_a*
603 and *Pces_f*), *Pa_T* (NCBI BioProject acc. no. PRJNA596862), *Pf_{eH}* (NCBI BioProject acc.
604 no. PRJNA727075) and a curated library of representative transposable elements
605 from *Viridiplantae* (<https://www.girinst.org/repbase/>). Because trees are not annual
606 plants, the identity obtained from the resulting .pass.list file was used for the
607 estimation of generation time after LTR insertion using the formula $T=K/2\mu$ (K is the
608 divergence of the LTR = 1 – identity) assuming a *Prunus* specific mutation frequency
609 of $\mu=7.7 \times 10^{-9}$ (Xie et al. 2016) per generation.

610

611 *Protein clustering, multiple sequence alignment and divergence of time estimation*

612 The protein datasets from *Pces_a* and *Pces_f*, *Pa_T*, *Pf_{eH}*, *Pp* (*Prunus persica* Whole
613 Genome Assembly v2.0, v2.0.a1), *Pm* (*Prunus mume* Tortuosa Genome v1.0), *Py*
614 (*Prunus yedoensis* var. *nudiflora* Genome v1.0), *Md* (*Malus x domestica* HFTH1 Whole
615 Genome v1.0) and *At* (TAIR10.1, RefSeq GCF_000001735.4) from the annotation step
616 were uploaded to Galaxy_Europe server as .fasta. The Proteinortho (Galaxy Version
617 6.0.32+galaxy0) was used to find orthologous proteins within the datasets. MAFFT
618 (Galaxy Version 7.505+galaxy0) was used to align the obtained single copy
619 orthogroups. The final alignments were merged with the Merge.files function (Galaxy
620 Version 1.39.5.0). Finally, the alignments were concatenated into a super protein and
621 the final sequences were aligned with MAFFT. A phylogenetic tree was reconstructed
622 with RAxML (maximum likelihood based inference of large phylogenetic trees, Galaxy
623 Version 8.2.4+galaxy3) and the obtained .nhx file was reformatted as .nwk file for
624 further processing using CLC Mainworkbench (21.0.1, QIAGEN Aarhus A/S).
625 Evolutionary analyses were conducted in MEGA X (Kumar et al. 2018). Estimation of
626 pairwise divergence time was performed according to Shirasawa et al. (2019) with a

627 divergence time from the reference species peach and apple (34-67 Mya,
628 www.timetree.org). Specific parameters for the calculation are listed in supplements.
629
630

631 **Data Availability**

632 Data supporting the findings of this study are deposited into the Open Agrar repository
633 (<https://doi.org/10.5073/20230324-105730-0>, Wöhner et al. (2023)) and on personal
634 request to the corresponding author. An assembly hub for genome and annotation
635 visualization is permanently hosted at [http://bioinf.uni-greifswald.de/private-
636 hubs/pcer/hub.txt](http://bioinf.uni-greifswald.de/private-hubs/pcer/hub.txt) .

637

638 **Author contributions**

639

640 TWW, OFE, HF wrote the manuscript. AHJW, KN, RPW and EJB performed DNA
641 isolation, sequencing, genome assembly and scaffolding. KJH, JK, LG, HT and TB
642 performed masking and annotation of the dataset. OA and LB performed the LTR-
643 dating and TW did the corresponding analysis and results compilation. JK performed
644 BUSCO analysis, read mapping and coverage analysis. TW the did the k-mer analysis,
645 interproscan, synteny analysis, phylogenetic analysis and circus plots. TW and JL
646 performed the Lai-Index analysis. HF, MS and AP conceived the study and made
647 substantial contributions to its design, acquisition, analysis and interpretation of data.
648 All authors contributed equally to the finalization of the manuscript.

649

650 **Acknowledgments**

651 We would acknowledge the Galaxy Europe, Galaxy USA server administration for
652 support and provision of resources. Special thanks to Eric Lyon for support with
653 Synmap2.

654

655 **References**

656 **Arang Rhie** (2020). Meryl. In GitHub repository. GitHub.
657 <https://github.com/marbl/meryl>.

658 **Arulsekar, S., and Parfitt, D.E.** (1986). Isozyme Analysis Procedures for Stone Fruits,
659 Almond, Grape, Walnut, Pistachio, and Fig. *HortSci* **21** (4): 928–933.

660 **Beaver, J.A., and Iezzoni, A.F.** (1993). Allozyme Inheritance in Tetraploid Sour Cherry
661 (*Prunus cerasus* L.). *jashs* **118** (6): 873–877.

662 **Bertioli, D.J., Jenkins, J., Clevenger, J., Dudchenko, O., Gao, D., Seijo, G., Leal-**
663 **Bertioli, S.C.M., Ren, L., Farmer, A.D., and Pandey, M.K., et al.** (2019). The
664 genome sequence of segmental allotetraploid peanut *Arachis hypogaea*. *Nat Genet* **51** (5): 877–884.

665 **Bird, K.A., Jacobs, M., Sebolt, A., Rhoades, K., Alger, E.I., Colle, M., Alekman,**
666 **M.L., Bies, P.K., Cario, A.J., and Chigurupati, R.S., et al.** (2022). Parental origins
667 of the cultivated tetraploid sour cherry (*Prunus cerasus* L.). *Plants People Planet* **4**
668 (5): 444–450.

669 **Bolger, A.M., Lohse, M., and Usadel, B.** (2014). Trimmomatic: a flexible trimmer for
670 Illumina sequence data. *Bioinformatics* **30** (15): 2114–2120.

671 **Brůna, T., Hoff, K.J., Lomsadze, A., Stanke, M., and Borodovsky, M.** (2021).
672 BRAKER2: automatic eukaryotic genome annotation with GeneMark-EP+ and
673 AUGUSTUS supported by a protein database. *NAR Genom Bioinform* **3** (1):
674 lqaa108.

675 **Brůna, T., Lomsadze, A., and Borodovsky, M.** (2020). GeneMark-EP+: eukaryotic
676 gene prediction with self-training in the space of genes and proteins. *NAR Genom*
677 *Bioinform* **2** (2): lqaa026.

678 **Cai, L., Stegmeir, T., Sebolt, A., Zheng, C., Bink, Marco C. A. M., and Iezzoni, A.**
679 (2018). Identification of bloom date QTLs and haplotype analysis in tetraploid sour
680 cherry (*Prunus cerasus*). *Tree Genetics & Genomes* **14** (2): 1–11.

681 **Chester, M., Gallagher, J.P., Symonds, V.V., Da Cruz Silva, A.V., Mavrodiev, E.V.,**
682 **Leitch, A.R., Soltis, P.S., and Soltis, D.E.** (2012). Extensive chromosomal
683 variation in a recently formed natural allopolyploid species, *Tragopogon miscellus*
684 (Asteraceae). *Proceedings of the National Academy of Sciences of the United*
685 *States of America* **109** (4): 1176–1181.

686

687 **Clarke, J.B., and Tobutt, K.R.** (2009). A standard set of accessions, microsatellites
688 and genotypes for harmonising the fingerprinting of cherry collections for the
689 ECPGR. *Acta Horticulturae* (814 (VOL 2)): 615–618.

690 **Clayton, J.W., and Tretiak, D.N.** (1972). Amine-Citrate Buffers for p H Control in
691 Starch Gel Electrophoresis. *J. Fish. Res. Bd. Can.* **29** (8): 1169–1172.

692 **Coart, E., Vekemans, X., Smulders, M.J.M., Wagner, I., van Huylenbroeck, J., van**
693 **Bockstaele, E., and Roldán-Ruiz, I.** (2003). Genetic variation in the endangered
694 wild apple (*Malus sylvestris* (L.) Mill.) in Belgium as revealed by amplified fragment
695 length polymorphism and microsatellite markers. *Molecular Ecology* **12** (4): 845–
696 857.

697 **Cock, P.J.A., Grüning, B.A., Paszkiewicz, K., and Pritchard, L.** (2013). Galaxy tools
698 and workflows for sequence analysis with applications in molecular plant
699 pathology. *PeerJ* **1**: e167.

700 **Dirlewanger, E., Claverie, J., Wünsch, A., and Iezzoni, A.F.** (2007). Cherry. In *Fruits*
701 and *Nuts* (Springer, Berlin, Heidelberg), pp. 103–118.

702 **Edger, P.P., Poorten, T.J., VanBuren, R., Hardigan, M.A., Colle, M., McKain, M.R.,**
703 **Smith, R.D., Teresi, S.J., Nelson, A.D.L., and Wai, C.M., et al.** (2019). Origin and
704 evolution of the octoploid strawberry genome. *Nat Genet* **51** (3): 541–547.

705 **Flynn, J.M., Hubley, R., Goubert, C., Rosen, J., Clark, A.G., Feschotte, C., and**
706 **Smit, A.F.** (2020). RepeatModeler2 for automated genomic discovery of
707 transposable element families. *Proc. Natl. Acad. Sci. U.S.A.* **117** (17): 9451–9457.

708 **Gabriel, L., Hoff, K.J., Brůna, T., Borodovsky, M., and Stanke, M.** (2021). TSEBRA:
709 transcript selector for BRAKER. *BMC Bioinformatics* **22** (1): 566.

710 **Hancock, J.F.** (2008). *Temperate Fruit Crop Breeding: Germplasm to Genomics*
711 (Dordrecht: Springer Science+Business Media B.V.).

712 **Hao, Y., Fleming, J., Petterson, J., Lyons, E., Edger, P.P., Pires, J.C., Thorne, J.L.,**
713 **and Conant, G.C.** (2022). Convergent evolution of polyploid genomes from across
714 the eukaryotic tree of life. *G3 (Bethesda, Md.)* **12** (6).

715 **Haug-Batzell, A., Stephens, S.A., Davey, S., Scheidegger, C.E., and Lyons, E.**
716 (2017). SynMap2 and SynMap3D: web-based whole-genome synteny browsers.
717 *Bioinformatics* **33** (14): 2197–2198.

718 **Hawkins, J.S., Proulx, S.R., Rapp, R.A., and Wendel, J.F.** (2009). Rapid DNA loss
719 as a counterbalance to genome expansion through retrotransposon proliferation
720 in plants. *Proceedings of the National Academy of Sciences of the United States
721 of America* **106** (42): 17811–17816.

722 **Höfer, M., Flachowsky, H., Schröpfer, S., and Peil, A.** (2021). Evaluation of Scab
723 and Mildew Resistance in the Gene Bank Collection of Apples in Dresden-Pillnitz.
724 *Plants* **10** (6): 1227.

725 **Hoff, K.J.** (2019). MakeHub: Fully automated generation of UCSC Genome Browser
726 Assembly Hubs (Cold Spring Harbor Laboratory).

727 **Hoff, K.J., Lange, S., Lomsadze, A., Borodovsky, M., and Stanke, M.** (2016).
728 BRAKER1: Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-
729 ET and AUGUSTUS. *Bioinformatics* **32** (5): 767–769.

730 **Hoff, K.J., Lomsadze, A., Borodovsky, M., and Stanke, M.** (2019). Whole-Genome
731 Annotation with BRAKER. In *Gene Prediction* (Humana, New York, NY), pp. 65–95.

732 **Hrotkó, K., Feng, Y., and Halász, J.** (2020). Spontaneous hybrids of *Prunus fruticosa*
733 Pall. in Hungary. *Genet Resour Crop Evol* **67** (2): 489–502.

734 **Hunter, S., Apweiler, R., Attwood, T.K., Bairoch, A., Bateman, A., Binns, D., Bork,
735 P., Das, U., Daugherty, L., and Duquenne, L., et al.** (2009). InterPro: the
736 integrative protein signature database. *Nucleic Acids Res* **37** (Database issue):
737 D211-5.

738 **Jakob, S.S., Meister, A., and Blattner, F.R.** (2004). The considerable genome size
739 variation of *Hordeum* species (poaceae) is linked to phylogeny, life form, ecology,
740 and speciation rates. *Mol Biol Evol* **21** (5): 860–869.

741 **Janick, J.** (2010). *Horticultural Reviews*, Volume 19 (New York, NY: John Wiley &
742 Sons).

743 **Jia, J., Xie, Y., Cheng, J., Kong, C., Wang, M., Gao, L., Zhao, F., Guo, J., Wang,
744 K., and Li, G., et al.** (2021). Homology-mediated inter-chromosomal interactions
745 in hexaploid wheat lead to specific subgenome territories following
746 polyploidization and introgression. *Genome Biol* **22** (1): 26.

747 **Jo, Y., Chu, H., Cho, J.K., Choi, H., Lian, S., and Cho, W.K.** (2015). De novo
748 transcriptome assembly of a sour cherry cultivar, Schattenmorelle. *Genomics Data*
749 **6**: 271–272.

750 **Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H.,**
751 **Maslen, J., Mitchell, A., and Nuka, G., et al. (2014). InterProScan 5: genome-**
752 **scale protein function classification. Bioinformatics** **30** (9): 1236–1240.

753 **Jones, D.T., Taylor, W.R., and Thornton, J.M. (1992). The rapid generation of**
754 **mutation data matrices from protein sequences. Computer applications in the**
755 **biosciences CABIOS** **8** (3): 275–282.

756 **José Quero-García, Amy Iezzoni, Gregorio López-Ortega, Cameron Peace,**
757 **Mathieu Fouché, Elisabeth Dirlewanger, and Mirko Schuster (2019). Advances**
758 **and challenges in cherry breeding. In Achieving sustainable cultivation of**
759 **temperate zone tree fruits and berries (Burleigh Dodds Science Publishing), pp.**
760 **55–88.**

761 **Keilwagen, J., Hartung, F., and Grau, J. (2019). GeMoMa: Homology-Based Gene**
762 **Prediction Utilizing Intron Position Conservation and RNA-seq Data. In Gene**
763 **Prediction (Humana, New York, NY), pp. 161–177.**

764 **Keilwagen, J., Wenk, M., Erickson, J.L., Schattat, M.H., Grau, J., and Hartung, F.**
765 **(2016). Using intron position conservation for homology-based gene prediction.**
766 **Nucleic Acids Res** **44** (9): e89.

767 **Kim, D., Paggi, J.M., Park, C., Bennett, C., and Salzberg, S.L. (2019). Graph-based**
768 **genome alignment and genotyping with HISAT2 and HISAT-genotype. Nat**
769 **Biotechnol** **37** (8): 907–915.

770 **Kim, S., Park, M., Yeom, S.-I., Kim, Y.-M., Lee, J.M., Lee, H.-A., Seo, E., Choi, J.,**
771 **Cheong, K., and Kim, K.-T., et al. (2014). Genome sequence of the hot pepper**
772 **provides insights into the evolution of pungency in Capsicum species. Nat Genet**
773 **46** (3): 270–278.

774 **Knight, C.A., and Beaulieu, J.M. (2008). Genome size scaling through phenotype**
775 **space. Ann Bot** **101** (6): 759–766.

776 **Kobel, F. (1927). Zytologische Untersuchungen an Prunoideen und Pomoideen. Art.**
777 **Institut Orell Füssli**

778 **Krebs, S.L., and Hancock, J.F. (1989). Tetrasomic inheritance of isoenzyme markers**
779 **in the highbush blueberry, Vaccinium corymbosum L. Heredity** **63** (1): 11–18.

780 **Kriventseva, E.V., Kuznetsov, D., Tegenfeldt, F., Manni, M., Dias, R., Simão, F.A.,**
781 **and Zdobnov, E.M. (2019). OrthoDB v10: sampling the diversity of animal, plant,**

782 fungal, protist, bacterial and viral genomes for evolutionary and functional
783 annotations of orthologs. *Nucleic Acids Res* **47** (D1): D807-D811.

784 **Kumar S., Stecher G., Li M., Knyaz C., and Tamura K.** (2018). MEGA X: Molecular
785 Evolutionary Genetics Analysis across computing platforms. *Molecular Biology*
786 and Evolution

787 **Langmead, B., and Salzberg, S.L.** (2012). Fast gapped-read alignment with Bowtie
788 2. *Nat Methods* **9** (4): 357–359.

789 **Lansari, A., and Iezzoni, A.** (1990). A Preliminary Analysis of Self-incompatibility in
790 Sour Cherry. *HortSci* **25** (12): 1636–1638.

791 **Le Dantec, L., Girollet, N., Gouzy, J., Sallet, E., Fouche, M., Quero-Garcia, J., &**
792 **Dirlewanger, E.** (2019). An Improved Assembly of the Diploid 'Regina' Sweet
793 Cherry Genome. In *Plant & Animal Genome Conference XXVII* (p. PO1193).

794 **Leitch, I.J., Hanson, L., Lim, K.Y., Kovarik, A., Chase, M.W., Clarkson, J.J., and**
795 **Leitch, A.R.** (2008). The ups and downs of genome size evolution in polyploid
796 species of *Nicotiana* (Solanaceae). *Ann Bot* **101** (6): 805–814.

797 **Li, H.** (2016). Minimap and miniasm: fast mapping and de novo assembly for noisy
798 long sequences. *Bioinformatics* **32** (14): 2103–2110.

799 **Li, H.** (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*
800 **34** (18): 3094–3100.

801 **Li, L., Chen, M., Zhang, X., and Jia, X.** (2022). Spatial Distribution Pattern of Root
802 Sprouts under the Canopy of *Malus sieversii* in a Typical River Valley on the
803 Northern Slopes of the Tianshan Mountain. *Forests* **13** (12): 2044.

804 **Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G.,**
805 **Abecasis, G., and Durbin, R.** (2009). The Sequence Alignment/Map format and
806 SAMtools. *Bioinformatics* **25** (16): 2078–2079.

807 **Lomsadze, A., Burns, P.D., and Borodovsky, M.** (2014). Integration of mapped
808 RNA-Seq reads into automatic training of eukaryotic gene finding algorithm.
809 *Nucleic Acids Res* **42** (15): e119.

810 **Maková, L., Vít, P., and Urfus, T.** (2018). Crop-to-wild hybridization in cherries-
811 Empirical evidence from *Prunus fruticosa*. *Evolutionary Applications* **11** (9): 1748–
812 1759.

813 **Mason, A.S., and Wendel, J.F.** (2020). Homoeologous Exchanges, Segmental
814 Allopolyploidy, and Polyploid Genome Evolution. *Front. Genet.* **11**: 1014.

815 Multivariate analysis of a sour cherry germplasm c germplasm collection. *J Amer Soc*
816 *Hort Sci* .

817 **OLDÉN, E.J., and NYBOM, N.** (1968). ON THE ORIGIN OF PRUNUS CERASUS L.
818 *Hereditas* **59** (2-3): 327–345.

819 **Ordidge, M., Lithauer, S., Venison, E., Blouin-Delmas, M., Fernandez-Fernandez,**
820 **F., Höfer, M., Kägi, C., Kellerhals, M., Marchese, A., and Mariette, S.** (2021).
821 Towards a Joint International Database: Alignment of SSR Marker Data for
822 European Collections of Cherry Germplasm. *Plants* **10** (6): 1243.

823 **Ou, S., Chen, J., and Jiang, N.** (2018). Assessing genome assembly quality using the
824 LTR Assembly Index (LAI). *Nucleic Acids Research* **46** (21): e126.

825 **Ou, S., Su, W., Liao, Y., Chougule, K., Agda, J.R.A., Hellinga, A.J., Lugo, C.S.B.,**
826 **Elliott, T.A., Ware, D., and Peterson, T.**, et al. (2019). Benchmarking transposable
827 element annotation methods for creation of a streamlined, comprehensive
828 pipeline. *Genome Biol* **20** (1): 275.

829 **Pinosio, S., Marroni, F., Zuccolo, A., Vitulo, N., Mariette, S., Sonnante, G.,**
830 **Aravanopoulos, F.A., Ganopoulos, I., Palasciano, M., and Vidotto, M.**, et al.
831 (2020). A draft genome of sweet cherry (*Prunus avium* L.) reveals genome-wide
832 and local effects of domestication. *The Plant Journal* **103** (4): 1420–1432.

833 **Pierce, B. A.** (2012). *Genetics: a conceptual approach*. Macmillan.

834 **Potter, D., Eriksson, T., Evans, R.C., Oh, S., Smedmark, J.E.E., Morgan, D.R.,**
835 **Kerr, M., Robertson, K.R., Arsenault, M., and Dickinson, T.A.**, et al. (2007).
836 Phylogeny and classification of Rosaceae. *Plant Syst. Evol.* **266** (1-2): 5–43.

837 **Quevillon, E., Silventoinen, V., Pillai, S., Harte, N., Mulder, N., Apweiler, R., and**
838 **Lopez, R.** (2005). InterProScan: protein domains identifier. *Nucleic Acids Research*
839 **33** (Web Server issue): W116-20.

840 **Ranallo-Benavidez, T.R., Jaron, K.S., and Schatz, M.C.** (2020). GenomeScope 2.0
841 and Smudgeplot for reference-free profiling of polyploid genomes. *Nat Commun*
842 **11** (1): 1432.

843 **Raptopoulos, T.** (1941). Chromosomes and fertility of cherries and their hybrids.
844 *Journal of Genetics*, 42(1), 91-114.

845 **Renny-Byfield, S., Chester, M., Kovařík, A., Le Comber, S.C., Grandbastien, M.-**
846 **A., Deloger, M., Nichols, R.A., Macas, J., Novák, P., and Chase, M.W., et al.**
847 (2011). Next generation sequencing reveals genome downsizing in allotetraploid
848 *Nicotiana tabacum*, predominantly through the elimination of paternally derived
849 repetitive DNAs. *Mol Biol Evol* **28** (10): 2843–2854.

850 **Rutland, C.A., Hall, N.D., and McElroy, J.S.** (2021). The Impact of Polyploidization
851 on the Evolution of Weed Species: Historical Understanding and Current
852 Limitations. *Front. Agron.* **3**: 5.

853 **Santi, F., and Lemoine, M.** (1990). Genetic markers for *Prunus avium* L. 2. Clonal
854 identifications and discrimination from *P cerasus* and *P cerasus* × *P avium*. *Ann.*
855 *For. Sci.* **47** (3): 219–227.

856 **Schiessl, S.-V., Katche, E., Ihien, E., Chawla, H.S., and Mason, A.S.** (2019). The
857 role of genomic structural variation in the genetic improvement of polyploid crops.
858 *The Crop Journal* **7** (2): 127–140.

859 **Schuster, M., Apostol, J., Iezzoni, A., Jensen, M., and Milatović, D.** (2017). Sour
860 cherry varieties and improvement. In *Cherries: botany, production and uses*
861 (Wallingford: CABI), pp. 95–116.

862 **Schuster, M., and Schreibner, H.** (2000). GENOME INVESTIGATION IN SOUR
863 CHERRY, *P. CERASUS* L. *Acta Hortic.* (538): 375–379.

864 **Schuster, M., and Wolfram, B.** (2008). NEW SOUR CHERRY CULTIVARS FROM
865 DRESDEN-PILLNITZ. *Acta Hortic.* (795): 83–86.

866 **Shirasawa, K., Isuzugawa, K., Ikenaga, M., Saito, Y., Yamamoto, T., Hirakawa,**
867 **H., and Isobe, S.** (2017). The genome sequence of sweet cherry (*Prunus avium*)
868 for use in genomics-assisted breeding. *DNA Res* **24** (5): 499–508.

869 **Shirasawa, K., Esumi, T., Hirakawa, H., Tanaka, H., Itai, A., Ghelfi, A., Nagasaki,**
870 **H., and Isobe, S.** (2019). Phased genome sequence of an interspecific hybrid
871 flowering cherry, 'Somei-Yoshino' (*Cerasus* × *yedoensis*). *DNA Res* **26** (5): 379–
872 389.

873 **Smit, A.F., Hubley, R., and Green, P.** (2015). RepeatModeler Open-1.0. 2008–2015.
874 Seattle, USA: Institute for Systems Biology. Available from: <http://www.repeatmasker.org>, Last Accessed May 1., 2018.

875 **Smit, A.F., Hubley, R., and Green, P.** (2015). RepeatMasker Open-4.0. 2013–2015.

877 **Soltis, D.E., Haufler, C.H., Darrow, D.C., and Gastony, G.J.** (1983). Starch Gel
878 Electrophoresis of Ferns: A Compilation of Grinding Buffers, Gel and Electrode
879 Buffers, and Staining Schedules. *American Fern Journal* **73** (1): 9.

880 **Soltis, D.E., and Rieseberg, L.H.** (1986). AUTOPOLYPLOIDY IN *TOLMIEA*
881 *MENZIESII* (SAXIFRAGACEAE): GENETIC INSIGHTS FROM ENZYME
882 ELECTROPHORESIS. *American Journal of Botany* **73** (2): 310–318.

883 **Stanke, M., Diekhans, M., Baertsch, R., and Haussler, D.** (2008). Using native and
884 syntenically mapped cDNA alignments to improve *de novo* gene finding.
885 *Bioinformatics* **24** (5): 637–644.

886 **STEBBINS, G.L.** (1947). Types of polyploids; their classification and significance.
887 *Advances in genetics* **1**: 403–429 (Elsevier).

888 **Tamura, K., Battistuzzi, F.U., Billing-Ross, P., Murillo, O., Filipski, A., and Kumar, S.** (2012). Estimating divergence times in large molecular phylogenies.
889 *Proceedings of the National Academy of Sciences of the United States of America*
890 **109** (47): 19333–19338.

892 **Tamura, K., Tao, Q., and Kumar, S.** (2018). Theoretical Foundation of the RelTime
893 Method for Estimating Divergence Times from Variable Evolutionary Rates. *Mol*
894 *Biol Evol* **35** (7): 1770–1782.

895 **Tang, S., Lomsadze, A., and Borodovsky, M.** (2015). Identification of protein coding
896 regions in RNA transcripts. *Nucleic Acids Res* **43** (12): e78.

897 **Tao, Q., Tamura, K., Mello, B., and Kumar, S.** (2020). Reliable Confidence Intervals
898 for RelTime Estimates of Evolutionary Divergence Times. *Mol Biol Evol* **37** (1): 280–
899 290.

900 **The Galaxy platform for accessible, reproducible and collaborative biomedical**
901 **analyses.** (2022). update. *Nucleic Acids Research*, 2022, 50. Jg., Nr. W1, S.
902 W345-W351.

903 **Tillich, M., Lehwerk, P., Pellizzer, T., Ulbricht-Jones, E.S., Fischer, A., Bock, R.,**
904 **and Greiner, S.** (2017). GeSeq - versatile and accurate annotation of organelle
905 genomes. *Nucleic Acids Res* **45** (W1): W6-W11.

906 **van de Peer, Y., Maere, S., and Meyer, A.** (2009). The evolutionary significance of
907 ancient genome duplications. *Nat Rev Genet* **10** (10): 725–732.

908 **Vaser, R., Sović, I., Nagarajan, N., and Šikić, M.** (2017). Fast and accurate de novo
909 genome assembly from long uncorrected reads. *Genome Res.* **27** (5): 737–746.

910 **Verde, I., Jenkins, J., Dondini, L., Micali, S., Pagliarani, G., Vendramin, E., Paris, R., Aramini, V., Gazza, L., and Rossini, L.**, et al. (2017). The Peach v2.0 release: high-resolution linkage mapping and deep resequencing improve chromosome-scale assembly and contiguity. *BMC Genomics* **18** (1): 225.

914 **Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C.A., Zeng, Q., Wortman, J., and Young, S.K.**, et al. (2014). Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLOS ONE* **9** (11): e112963.

918 **Wang, J., Liu, W., Zhu, D., Hong, P., Zhang, S., Xiao, S., Tan, Y., Chen, X., Xu, L., and Zong, X.**, et al. (2020). Chromosome-scale genome assembly of sweet cherry (*Prunus avium* L.) cv. Tieton obtained using long-read and Hi-C sequencing. *Hortic Res* **7** (1): 122.

922 **Wang, X., Liu, S., Zuo, H., Zheng, W., Zhang, S., Huang, Y., Pingcuo, G., Ying, H., Zhao, F., and Li, Y.**, et al. (2021). Genomic basis of high-altitude adaptation in Tibetan *Prunus* fruit trees. *Current biology CB* **31** (17): 3848-3860.e8.

925 **Wang, Z., Miao, H., Liu, J., Xu, B., Yao, X., Xu, C., Zhao, S., Fang, X., Jia, C., and Wang, J.**, et al. (2019). *Musa balbisiana* genome reveals subgenome evolution and functional divergence. *Nat. Plants* **5** (8): 810–821.

928 **Wang, X., Morton, J.A., Pellicer, J., Leitch, I.J., and Leitch, A.R.** (2021). Genome downsizing after polyploidy: mechanisms, rates and selection pressures. *The Plant Journal* **107** (4): 1003–1015.

931 **Wöhner, T.W., Emeriewen, O.F., Wittenberg, A.H.J., Schneiders, H., Vrijenhoek, I., Halász, J., Hrotkó, K., Hoff, K.J., Gabriel, L., and Lempe, J.**, et al. (2021a). The draft chromosome-level genome assembly of tetraploid ground cherry (*Prunus fruticosa* Pall.) from long reads. *Genomics* **113** (6): 4173–4183.

935 **Wöhner, T.W., Emeriewen, O.F., Wittenberg, A.H.J., Schneiders, H., Vrijenhoek, I., Halász, J., Hrotkó, K., Hoff, K.J., Gabriel, L., and Lempe, J.**, et al. (2021b). Supporting Materials for - The Draft Chromosome-level Genome Assembly of Tetraploid Ground Cherry (*Prunus fruticosa* Pall.) from Long Reads.

939 https://www.openagrар.de/receive/openagrар_mods_00070329 (2021) (accessed
940 1 June 2021)

941 **Wöhner, T., Emeriewen, Ofere, F., Wittenberg, Alexander, H.J., Nijbroek, K.,**
942 **Wang, Rui, Peng, Blom, E.-J., Keilwagen, J., Berner, T., Hoff, Katharina, J.,**
943 **and Gabriel, L., et al. (2023). Data set: The structure of the tetraploid sour cherry**
944 **'Schattenmorelle' (Prunus cerasus L.) genome reveals insights into its segmental**
945 **allopolyploid nature. <https://doi.org/10.5073/20230324-105730-0>**

946 **Xie, Z., Wang, L., Wang, L., Wang, Z., Lu, Z., Tian, D., Yang, S., and Hurst, L.D.**
947 (2016). Mutation rate analysis via parent-progeny sequencing of the perennial
948 peach. I. A low rate in woody perennials and a higher mutagenicity in hybrids.
949 Proceedings. Biological sciences **283** (1841).

950 **Yan, M., Zhang, X., Zhao, X., and Yuan, Z. (2019). The complete mitochondrial**
951 **genome sequence of sweet cherry (Prunus avium cv. 'summit'). Mitochondrial DNA**
952 **Part B** **4** (1): 1996–1997.

953 **Zdobnov, E. M., and Apweiler, R. (2001). InterProScan - an integration platform for**
954 **the signature-recognition methods in InterPro. Bioinformatics**, **17**(9), 847–848.
955 <https://doi.org/10.1093/bioinformatics/17.9.847>

956 **Zdobnov, E.M., and Apweiler R. (2001). R. InterProScan: protein domains identifier.**
957 **Bioinformatics (Oxford, England)** **17**: 847–848.

958 **Zenil-Ferguson, R., Ponciano, J.M., and Burleigh, J.G. (2016). Evaluating the role**
959 **of genome downsizing and size thresholds from genome size distributions in**
960 **angiosperms. American Journal of Botany** **103** (7): 1175–1186.

961 **Zhang, Z., Fu, T., Liu, Z., Wang, X., Xun, H., Li, G., Ding, B., Dong, Y., Lin, X., and**
962 **Sanguinet, K.A., et al. (2019). Extensive changes in gene expression and**
963 **alternative splicing due to homoeologous exchange in rice segmental**
964 **allopolyploids. Theor Appl Genet** **132** (8): 2295–2308.

965

966 **Figure legends**

967 **Figure 1 Morphology of *P. cerasus* L. 'Schattenmorelle'. (a) mature tree habitus, (b)**
968 **leaves, (c) inflorescence, (d) fruits.**

969

970 **Figure 2** The genome of *P. cerasus* 'Schattenmorelle'. Circos plot of 16
971 pseudomolecules of the subgenomes of *Pce_s_a* and *Pce_s_f*. (a) chromosome length
972 (Mb); (b) gene density in blocks of 250k; (c) distribution of repetitive sequences in
973 blocks of 250k (d) Gypsy elements in blocks of 250k; (e) Copia elements in block of
974 250k; (f) GC content in blocks of 1 Mb. (g) The inner ring shows markers from the 6+9k
975 SNP array located on both subgenomes.

976 **Figure 3** GenomeScope (Galaxy Version 2.0) estimation of the *P. cerasus* genome
977 size by k-mer counts obtained from the software Meryl (Galaxy Version 1.3+galaxy2).
978 Both programs are integrated on the GalaxyServerEurope. The k-mer-peaks indicate
979 that k-mers with a length of 19 bp occur in heterozygote (100x depth, 200x depth,
980 300x depth) and homozygote (400x depth) constitution within the genome. Coverage
981 depth of individual k-mers is assigned as coverage.

982

983 **Fig.4** Detected regions of homoeologous exchanges in the genome of *P. cerasus*
984 'Schattenmorelle'. Circos plot of 16 pseudomolecules of the subgenomes of *Pce_s_a*
985 and *Pce_s_f*. (a) chromosome length (Mb); (b) 16 in *Pce_s_a* and 12 in *Pce_s_f* detected
986 regions that match all three following analysis methods: (c) 1024 regions (100k
987 window) were intraspecific %-covered bases from mapped reads (*Pce_s_a* to *Pa_T*,
988 *Pce_s_f* to *Pf_{eH}*) was < than interspecific %-covered bases from mapped reads (*Pce_s_a*
989 to *Pf_{eH}*, *Pce_s_f* to *Pa_T*); (d) 148 regions were intraspecific difference of %-covered
990 bases from obtained RNAseq reads (*Pa* and *Pce_s_a*, *Pf* and *Pce_s_f*) > than interspecific
991 difference of %-covered bases from obtained RNAseq reads (*Pf* and *Pce_s_a*, *Pa* and
992 *Pce_s_f*); (e) 367 regions were the proportion of transcripts with intraspecific amino acid
993 identity (*Pa_T* and *Pce_s_a*, *Pf_{eH}* and *Pce_s_f*) < than the proportion of transcripts with
994 interspecific amino acid identity (*Pf_{eH}* and *Pce_s_a*, *Pa_T* and *Pce_s_f*).

995

996 **Fig. 5** Investigation on the evolution of the genome of *P. cerasus* 'Schattenmorelle'. **A**
997 Determination of insertion time from shared long terminal repeats (LTRs) in *P. cerasus*
998 subgenome *avium* (*Pce_s_a*) and *P. cerasus* subgenome *fruticosa* (*Pce_s_f*) compared
999 to *P. avium* 'Tieton' (*Pa_T*) and *P. fruticosa* ecotype Hármashatárhegy (*Pf_{eH}*). **B**
1000 Estimation of divergence of time (Mya) of *P. cerasus* subgenomes *Pce_s_a* and *Pce_s_f*
1001 compared to the donor species *P. avium* (*Pa*) and *P. fruticosa* (*Pf*). *Prunus yedonensis*

1002 (Pyn); *Prunus avium* (Pa); *Prunus persica* (Pp); *Prunus mume* (Pm); *Malus domestica*
1003 (Md); Paleocene (PAL); Eocene (EOC); Oligocene (OLI); Miocene (MIO); Pliocene (PII);
1004 Pleistocene (PLEI).

1005 **Tables**

1006 **Table 1** Characterization of repetitive sequences of *P. fruticosa* ecotype
1007 Hármashtárhegy (*Pf_{eH}*) compared to *P. avium* 'Tieton' (*Pa_T*), *P. persica* 'Lovell', and
1008 the two subgenomes of *P. cerasus* 'Schattenmorelle' *Pce_S_a* and *Pce_S_f*

1009 **Table 2** Comparison between the number of transcripts and %-IAA obtained from *P.*
1010 *fruticosa* ecotype Hármashtárhegy (*Pf_{eH}*) and *P. avium* cv 'Tieton' (*Pa_T*) representing
1011 the two ancestral species of *P. cerasus*

1012

1013 **Supplemental information (SI)**

1014

1015 **Document S1 - Supplemental material and methods**

1016

1017 **Supplemental Figures**

1018 **Figure S1** Hi-C heatmap post-scaffolding for the subgenomes *Pce_S_a* (left) and *Pce_S_f*
1019 (right) of *P. cerasus* cv 'Schattenmorelle'. The heatmap indicates the density of paired
1020 Hi-C reads which interact to each other in close proximity. High intense colour
1021 indicates high interaction.

1022

1023 **Figure S2** The chloroplast (a) and mitochondrial (b) sequences of *P. cerasus* L. cv
1024 'Schattenmorelle'.

1025

1026 **Figure S3** Analysis of completeness of the *P. cerasus* cv. Schattenmorelle
1027 subgenomes *P. cerasus* cv 'Schattenmorelle' subgenome *avium* (*Pce_S_a*) and *P.*
1028 *cerasus* cv 'Schattenmorelle' subgenome *fruticosa* (*Pce_S_f*) and combined datasets
1029 compared to *P. avium* cv. 'Tieton' (*Pa_T*) and *P. persica* cv. Lovell (*Pp_L*) by mapping of
1030 a set of universal single-copy orthologs using BUSCO. The bar charts indicate
1031 complete single copy (orange), complete duplicated (gray), fragmented (yellow) and
1032 missing (blue) genes. For evaluation the *embryophyta_odb10* BUSCO dataset
1033 (n=1614) was used. *P. cerasus* cv. Schattenmorelle show a 99 % completeness (S:

1034 16.7 %, D: 82.3 %, F: 0.4 %, M: 0.6 %, n: 1614) which reaches the completeness of
1035 *P. avium* cv. 'Tieton' (C: 98.3 %, S: 95.6 %, D: 2.7 %, F: 0.5 %, M: 1.5 %, n: 1614) and
1036 *P. persica* 'Lovell' (C: 99.3 %, S: 97.5 %, D: 1.8 %, F: 0.1 %, M: 0.6 %, n: 1614).

1037

1038 **Figure S4** Assessing the quality of repetitive sequences between the chromosome
1039 sequences of *P. avium* 'Tieton' (A), *P. fruticosa* ecotype Hármashtárhegy (B), *P.*
1040 *persica* 'Lovell', and (C) *P. cerasus* subgenome *avium* and *P. cerasus* subgenome
1041 *fruticosa* using the LAI index. The genomes *P. cerasus* [this study] and *P. avium* were
1042 sequenced with ONT 9.4.1 and Illumina (Wang et al. 2020), *P. fruticosa* with ONT
1043 10.3 (Wöhner et al. 2021) and *P. persica* with Illumina and Sanger sequencing of
1044 fosmid and BAC clones (Verde et al. 2017).

1045

1046 **Figure S5** Collinearity plots between the five published genetic maps of sour cherry
1047 (M172x25-F1, US-F1, 25x25-F1, Montx25-F1, RE-F1) and the *P. cerasus* cv
1048 'Schattenmorelle' subgenome *avium* (*Pce_S_a*) and *fruticosa* (*Pce_S_f*). X-axis represents
1049 the genetic position of a marker in the genetic linkage map given in centi Morgan (cM).
1050 Y-axis represents the physical position of a marker sequence within the genome
1051 sequence of the respective subgenome given in Mega base pairs (Mbp).

1052

1053 **Figure S6** Percentage of *P. cerasus* (*Pce*) proteins by IAA compared with 15 reference
1054 species. *P. fruticosa* ecotype Hármashtárhegy (*P_{fr}H*), *P. avium* 'Tieton' (*Pa_T*), *P.*
1055 *yedonensis* (*Pyed*), *P. domestica* (*Pd*), *P. armeniaca* (*Par*), *P. persica* (*Pp*), *Pyrus*
1056 *communis* (*Pyrco*), *Populus trichocarpa* (*Poptri*), *Vitis vinifera* (*Vv*), *Arabidopsis*
1057 *thaliana* (*At*), *Malus domestica* (*Md*).

1058

1059 **Figure S7** Synmap2 plots of self-self comparisons between (A) *Prunus persica*
1060 'Lovell' (*Pp_L*) and *P. avium* 'Sato Nishiki' (*Pa_S*), (B) *P. fruticosa* 'Hármashtárhegy'
1061 (*Pf_{EH}*), (C) *P. cerasus* *avium* 'Schattenmorelle' (*Pce_S_a*), (D) *P. cerasus* *fruticosa*
1062 'Schattenmorelle' (*Pce_S_f*) for the identification of triplicated regions (TR) 1-7.

1063

1064 **Figure S8** Positional co-linearity comparison between the two subgenomes *P.*
1065 *cerasus* *avium* 'Schattenmorelle' (*Pce_S_a*), *P. cerasus* *fruticosa* 'Schattenmorelle'

1066 (*Pce_S_f*) and *P. avium* 'Tieton' (*Pa_T*), *P. fruticosa* 'Hármashatárhegy' (*Pf_{eH}*), using the
1067 molecular markers from the 9+6k SNP array. The plots were generated using the R-
1068 software package chromoMap v0.4.1.

1069

1070 **Figure S9** Synteny between (a) *P. cerasus* subgenome _*avium* (*Pce_S_a*) and *P. cerasus*
1071 subgenome _*fruticosa* (*Pce_S_f*), and (b) the subgenomes and the genotypes *Prunus*
1072 *avium* 'Tieton' (*Pa_T*) and *P. fruticosa* (*Pf_{eH}*) of the ancestral species *P. avium* and *P.*
1073 *fruticosa*. Blue arrows indicate positions where inversions occurred.

1074

1075 **Figure S10** Detected regions of homoeologous exchanges in the genome of *P.*
1076 *cerasus* 'Schattenmorelle'. Circos plot of 16 single pseudomolecules 1 to 8 (see
1077 continuing plots) of the subgenomes of *Pce_S_a* and *Pce_S_f*. (a) chromosome length
1078 (Mb); (b) region in *Pce_S_a* and in *Pce_S_f* detected that match all three following analysis
1079 methods: (c) regions (100k window) were intraspecific %-covered bases from
1080 mapped reads (*Pce_S_a* to *Pa_T*, *Pce_S_f* to *Pf_{eH}*) was < than interspecific %-covered
1081 bases from mapped reads (*Pce_S_a* to *Pf_{eH}*, *Pce_S_f* to *Pa_T*); (d) regions were intraspecific
1082 difference of %-covered bases from obtained RNAseq reads (*Pa* and *Pce_S_a*, *Pf* and
1083 *Pce_S_f*) > than interspecific difference of %-covered bases from obtained RNAseq
1084 reads (*Pf* and *Pce_S_a*, *Pa* and *Pce_S_f*); (e) regions were the proportion of transcripts
1085 with intraspecific amino acid identity (*Pa_T* and *Pce_S_a*, *Pf_{eH}* and *Pce_S_f*) < than the
1086 proportion of transcripts with interspecific amino acid identity (*Pf_{eH}* and *Pce_S_a*, *Pa_T*
1087 and *Pce_S_f*).

1088

1089 **Supplemental Tables**

1090

1091 **Table S1** Statistics of different assemblies for *P. cerasus* cv 'Schattenmorelle' (*Pce_S*)
1092 and the subgenomes *Pce_S_a* and *Pce_S_f*

1093

1094 **Table S2** Pseudomolecule statistics for *Pce_S*

1095

1096 **Table S3** Iso-Seq results

1097

1098 **Table S4** Functional annotation results generated by interproscan using BRAKER &
1099 GeMoMa combination of ab-initio and homology-based structural gene annotation
1100 and statistics

1101

1102 **Table S5** Mapping of marker sequences to *Prunus cerasus* 'Schattenmorelle' (*Pce_S*)
1103 genome

1104

1105 **Table S6** Position of inversions within the subgenomes of *Pce_S* indicated by
1106 collinearity of marker positions.

1107

1108 **Supplemental Notes**

1109

1110 **Note S1** Access to the assembly hub for genome and annotation visualization at
1111 UCSC Genome Browser.

1112

1113 **Note S2** Calculation of the 70% quantile from IAA.

1114 **Note S3** Calculation of %-covered bases obtained from RNAseq data.