

1 **Chromosome-level and haplotype-resolved genome assembly**
2 **enabled by high-throughput single-cell sequencing of gamete**
3 **genomes**

4
5 José A. Campoy^{1,∇}, Hequan Sun^{1,2,∇}, Manish Goel¹, Wen-Biao Jiao¹, Kat Folz-Donahue³,
6 Nan Wang⁴, Manuel Rubio⁵, Chang Liu^{4,6}, Christian Kukat³, David Ruiz⁵, Bruno Huettel⁷ and
7 Korbinian Schneeberger^{1,2,*}

8
9 [∇]These authors contributed equally.

10
11 ¹Department of Chromosome Biology, Max Planck Institute for Plant Breeding Research,
12 Carl-von-Linné-Weg 10, 50829 Cologne, Germany ²Faculty of Biology, LMU Munich,
13 Großhaderner Str. 2, 82152 Planegg-Martinsried, Germany ³FACS & Imaging Core Facility,
14 Max Planck Institute for Biology of Ageing, 50931 Cologne, Germany ⁴Center for Plant
15 Molecular Biology (ZMBP), University of Tübingen, Auf der Morgenstelle 32, 72076
16 Tübingen, Germany ⁵Departament of Plant Breeding, CEBAS-CSIC, PO Box 164, E-30100,
17 Espinardo, Murcia, Spain ⁶Institute of Biology, University of Hohenheim, Garbenstraße 30,
18 70599 Stuttgart, Germany ⁷Max Planck-Genome-center Cologne, Carl-von-Linné-Weg 10,
19 50829 Cologne, Germany

20
21 *Correspondence: Korbinian Schneeberger (schneeberger@mpipz.mpg.de)

22
23 Key words: single-cell sequencing, haplotype-resolved assembly, haplotyping, phasing, *de*
24 *novo* assembly

25 **Generating chromosome-level, haplotype-resolved assemblies of heterozygous**
26 **genomes remains challenging. To address this, we developed gamete binning, a**
27 **method based on single-cell sequencing of haploid gametes enabling separation of**
28 **the whole-genome sequencing reads into haplotype-specific reads sets. After**
29 **assembling the reads of each haplotype, the contigs are scaffolded to chromosome-**
30 **level using a genetic map derived from the gametes. As a proof-of-concept, we**
31 **assembled the two genomes of a diploid apricot tree based on whole-genome**
32 **sequencing of 445 individual pollen grains. The two haplotype assemblies (N50: 25.5**
33 **and 25.8 Mb) featured a haplotyping precision of >99% and were accurately scaffolded**
34 **to chromosome-level.**

35 **Introduction**

36 Currently, most diploid genome assemblies ignore the differences between the
37 homologous chromosomes and assemble the genomes into one pseudo-haploid sequence,
38 which is an artificial consensus of both haplotypes. Such an artificial consensus can result in
39 imprecise gene annotation and misleading biological interpretation^{1,2}. To avoid these
40 problems, it is a common strategy to inbreed or to generate double-haploid genotypes to
41 enable the assembly of homozygous genomes.

42 Recent alternatives allowing for the assembly of both haplotypes include
43 chromosome sorting³, Strand-seq⁴⁻⁶ and high-throughput chromosome conformation capture
44 (Hi-C)⁷⁻¹³ sequencing. Chromosome sorting separates individual chromosomes before
45 sequencing, and thus enables the sequencing and assembly of individual haplotypes.
46 However, sorting of particular chromosomes may not always be possible if they cannot be
47 discriminated based on their fluorescence intensity or light scatter¹⁴ and may need tedious
48 generation of specific lines for sorting¹⁵. The more recent method Strand-seq is a single-cell
49 technique that requires neither parents nor gametes which can be potentially used to cluster
50 long sequencing reads by chromosome, phase haplotypes, and scaffold using genetic map
51 techniques, however, the difficulty for generating Strand-seq data has limited its application

52 to a narrow number of model species. In contrast, the analysis of the chromosome
53 conformation, including Hi-C technologies which enable the detection of chromatin
54 interactions at an unprecedented scale, has been successfully applied for haplotype phasing
55 and genome scaffolding for a wide range of species^{7,9-12,16}. However, despite its simple
56 application, Hi-C-based phasing can be error prone due to some weaknesses in defining the
57 alleles that distinguish haplotypes, which in turn can lead to haplotype switch errors¹¹ and
58 result in mis-scaffolding of small contigs due to the lack of sufficient informative connections
59 to other contigs^{8,12,13}. Also the reconstruction of whole chromosomes structures can be error-
60 prone as already one local mis-scaffolding is sufficient to introduce severe mis-assemblies
61 like falsely joining chromosome arms¹⁰. It is therefore necessary to carefully inspect
62 assemblies that rely on Hi-C for phasing or scaffolding to identify errors, which in turn require
63 correction based on additional evidence including, for example, the integration of genetic
64 maps^{8,10,17}.

65 An elegant alternative for haplotype phasing, called trio binning, is based on the
66 separation of whole-genome sequencing reads into haplotype-specific read sets before
67 assembly using the genomic differences between the parental genomes². While this is a
68 powerful method, it can be limiting if the parents are not available or are unknown¹⁸. A
69 solution for this is the sequencing of a few gamete genomes (derived from the focal
70 individual), which is sufficient for the inference of genome-wide haplotypes, but relies on
71 existing long-contiguity reference sequences¹⁹⁻²².

72 In addition to resolving haplotypes, the generation of chromosome-level assemblies,
73 which are necessary to understand the full complexity of genomic differences including all
74 kinds of structural rearrangements, is similarly challenging^{23,24}. While recent improvements in
75 long DNA molecule sequencing²⁵ and as mentioned above in Hi-C data generation promise
76 the assembly of telomere-to-telomere contigs, genetic maps can reliably help to resolve mis-
77 assemblies and guide chromosome-level scaffolding¹⁰. The generation of genetic maps,
78 however, relies on a substantial amount of meiotic recombination which usually implies the
79 genotyping of hundreds of recombinant genomes^{26,27}. Creating and genotyping sufficiently

80 large populations is not possible in some species (like many of the mammals including
81 humans), and for those species for which it is possible it can be time-consuming and costly,
82 and may post great challenges if the individuals show long juvenility or sterility¹⁶.

83 To address all these challenges, we present gamete binning, a method for
84 chromosome-level, haplotype-resolved genome assembly - independent of parental
85 genomes or recombinant progenies (Fig. 1). The method starts by isolating gamete nuclei
86 from the focal individual followed by high-throughput single-cell sequencing of hundreds of
87 the haploid gamete genomes. (For clarification, we collectively refer to both gametophytes in
88 plants and gametes in animals collectively as gametes, as both have haploid genomes.) The
89 segregation of sequence variation in the gamete genomes enables a straightforward phasing
90 of all variants into two haplotypes, which subsequently allows for genetic mapping and
91 sorting of whole-genome sequencing reads into distinct read sets - each representing a
92 different haplotype. Assembling these independent read sets leads to haplotype-resolved
93 genome assemblies, which can be scaffolded to chromosome-level using a gamete genome-
94 derived genetic map.

95 **Results**

96 **Preliminary diploid-genome assembly**

97 We used gamete binning to assemble the two haploid genomes of a specific, diploid
98 apricot tree (*Prunus armeniaca*; cultivar 'Rojo Pasion'²⁸), which grows in Murcia,
99 southeastern Spain (Supplementary Figure 1). We first performed a preliminary *de novo*
100 genome assembly using *Canu*²⁹ with 19.9 Gb long reads (PacBio, Supplementary Figure 2)
101 derived from DNA extracted from fruits and corresponding to 82x genome coverage
102 according to a genome size of ~242.5 Mb estimated by *findGSE*³⁰ (Methods; Supplementary
103 Figure 3). After purging haplotype-specific contigs, the curated assembly consisted of 939
104 contigs with a combined length of 230.9 Mb and an N50 of 563.8 kb, which represents a
105 haploid, but mosaic assembly of the apricot genome (Methods).

106 **High-throughput single-cell sequencing of pollen**

107 To advance this assembly, we isolated pollen grains from ten closed flowers (to avoid
108 contamination of foreign pollen) and released their nuclei following a protocol based on pre-
109 filtering followed by bursting³¹ (Fig. 1a; Methods). The nuclei mixture was cleaned up using
110 propidium iodide staining plus sorting by flow cytometry, leading to a solution with 12,600
111 nuclei that were loaded into a 10x Chromium Controller in two batches - each with 6,300
112 nuclei (Supplementary Figures 1a-d; Supplementary Figure 4; Methods). With this we
113 generated two 10x single-cell genome (CNV) sequencing libraries, which were sequenced
114 with 95 and 124 million 151 bp paired-end reads (Illumina). By exploring the *cellranger*-
115 corrected cell barcodes within the read data of both libraries, we extracted 691 read sets -
116 each with a minimum of 5,000 read pairs (Methods; Fig. 2a).

117 Aligning the reads of each pollen genome to the curated assembly, we found that the
118 reads of 246 sets featured high similarity to thrip genomes or included more than one haploid
119 genome, possibly due to random attachment of multiple nuclei during 10x Genomics library
120 preparation or the uncompleted separation of pollen nuclei during pollen maturation³²
121 (Supplementary Figure 6a-c; Methods). Thus, we selected a set of 445 haploid pollen
122 genomes. In general, the short-read alignments did not show any biases or preferences for
123 specific regions of the genome as reported for some of the single-cell genome amplification
124 kits, but covered nearly all regions (99.1%) of the curated assembly (Fig. 2b; Supplementary
125 Figure 6d).

126 **Haplotype phasing and genetic mapping**

127 With short read alignments, we identified 578,209 heterozygous SNPs on 702 contigs
128 with a total length of 218.0 Mb (Fig. 2b; Methods). Even though this implied 1 SNP marker
129 every 377 bp on average, we observed that the distances between some of the SNP markers
130 were larger than the usual long reads, which would hamper the haplotype assignment of
131 reads whenever they aligned to such regions. Overall, we observed 10,452 regions larger
132 than 2 kb without markers (110.9 Mb) including 237 regions (12.5 Mb), that spanned entire

133 contigs. Regions without markers occur if the two haplotypes are identical (which is a
134 common phenomenon in domesticated genomes) or if a region exists only in one of the
135 haplotypes (e.g. a large indel). We distinguished these two cases using the short-read
136 coverage of the combined pollen read sets, assuming that the regions that are only present
137 in one haplotype are supported by only approximately half of the reads (Methods). While
138 7,199 regions (74.5 Mb) were shared between the haplotypes (and were labelled as
139 conserved), we found that 3,253 regions (36.4 Mb) were specific to one of the haplotypes
140 (i.e. deletions; Fig. 2b). Such regions (i.e. deletions) which are specific to one haplotype can
141 also be used as markers. If such deletions were linked to nearby SNP markers, we phased
142 them according to their linked alleles. For deletions on contigs without additional markers, we
143 used the absence and presence of read alignments in the pollen to assign genotypes.

144 The haploid nature of the 445 selected individual pollen genomes allowed us to phase
145 all SNP and deletion markers into two haplotypes simply by using the linkage within the
146 pollen genomes (Fig. 2c-d). To phase the haplotypes across the contigs, we generated two
147 virtual markers for each contig representing the (imputed) alleles at both ends of the contig.
148 The markers were grouped into a genetic map with eight linkage groups (corresponding to
149 the eight homologous chromosome pairs) including 891 contigs with a total length of 228.0
150 Mb (corresponding to about 99% of the complete assembly) using *JoinMap* 4.0³³ (Fig. 2e;
151 Fig. 3a) (Methods).

152 **Haplotype-specific long read separation and chromosome-level assembly**

153 After this, we aligned the PacBio reads to the curated assembly. Using the phased
154 alleles (of the SNP and deletion markers) within each of the individual PacBio read
155 alignments, we separated 93.4% of the reads into one of 16 haplotype-specific clusters
156 representing the two haplotypes of each of the eight linkage groups. Reads that aligned in
157 regions that were conserved between the two haplotypes were randomly assigned to one of
158 the two haplotype-specific clusters (Fig. 3a; Methods). Similarity analyses revealed that most
159 of the remaining 6.6% reads were related to organellar genomes or repetitive sequences.

160 The 16 haplotype-specific read sets were independently assembled using *Flye*³⁴,
161 which led to 16 haplotype-specific chromosome assemblies with average N50 values ranging
162 from 662.3 to 664.6 kb (Table 1; Methods). Using the genetic map, we combined the contigs
163 of each assembly into a pseudo-molecule. This led to two haplotype-resolved chromosome-
164 level assemblies, both with N50 above 25.0 Mb (Fig. 3a-b; Methods).

165 To assess haplotype accuracy, we additionally whole-genome sequenced the
166 parental cultivars of 'Rojo Pasion' known as 'Currot' and 'Orange Red'. Using Illumina
167 sequencing technology, we generated 15.7 and 16.2 Gb short reads of each of the diploid
168 parental genomes, respectively. Overall, we found that ~99.1% of the *k*-mers that were
169 specific to one of the haplotype assemblies could be found in the corresponding parental
170 genome illustrating that almost all of the variation was correctly assigned to haplotypes (Fig.
171 3c; Table 1; Methods). Having proved the haplotype accuracy, the assemblies were polished
172 resulting in final haplotype assemblies. The final haplotype assembly sizes were 216.0 and
173 215.2 Mb for 'Currot'-genotype (8 scaffolds, N50: 25.8 Mb) and 'Orange Red'-genotype (8
174 scaffolds, N50: 25.5 Mb), respectively (Table 1).

175 We estimated the overall assembly quality by comparing the *k*-mer distributions of the
176 assemblies and the Illumina short read sets of the focal and parental using KAT³⁵ and
177 Merqury³⁶. Both haplotype genome assembly showed very high quality values ($QV > 36$) and
178 the absence of allelic duplications between the haplotypes, though a fraction of ~7% of the
179 heterozygous *k*-mers in the reads was missing in the assemblies (Supplementary Figures 7c,
180 8).

181 To further assess the overall structures of the assembled chromosomes, we
182 compared them to recently assembled chromosomes of very closely-related species such as
183 the heterozygous 'Chuanzhong' apricot (*Prunus armeniaca*)³⁷, the Japanese apricot
184 (*Prunus mume*)³⁸, and a more distantly-related species, peach (*Prunus persica*: doubled-
185 haploid genome)³⁹ using SyRI²³ (a tool designed for the comparison of chromosome-level
186 assemblies). Our assemblies showed high consistency in the synteny to these assemblies

187 across entire chromosomes, reflecting the reliability of the genetic map and the assembled
188 genome structures (Fig. 3d; Supplementary Figure 6).

189 As yet another way to assess the quality of the genome we generated two Hi-C
190 libraries from DNA extracted from leaves of Rojo Pasion and sequenced them totaling in
191 191.2 million read pairs (or ~240x genome coverage). We created Hi-C contact maps using
192 each of the homologous chromosome pairs separately as well as using the entire genome
193 (Fig. 3e; Fig. 4a; Supplementary Figures 9-15). In general, the contiguity of contact signals
194 surrounding the main diagonal of the map again demonstrated the high quality of the
195 structure of the assemblies.

196 **Comparing gamete binning with Hi-C based phasing and genome scaffolding**

197 However, the perhaps more interesting way to use the Hi-C data is its application for
198 genome phasing and scaffolding and the comparison of its assembly performance to that of
199 gamete binning.

200 Applying *ALLHiC*⁸ to the Hi-C reads sets generated 16 scaffolds (representing the 16
201 haploid chromosomes), with sizes ranging from 11.2 to 51.1 Mb (Methods). (Using a different
202 Hi-C-based phasing and scaffolding tool, *SALSA2*⁴⁰, did not lead to comparable results, thus
203 not compared further.). For comparison, we also generated Hi-C contact maps for *ALLHiC*
204 based assemblies (Fig. 4b). Interestingly, the contact maps of the gamete binning and
205 *ALLHiC* based assemblies were strikingly different. Only the gamete binning assembly
206 showed (beside the contact within the haplotypes) the expected contact signals between two
207 different haplotypes, which also were reported for other species^{8,41}. The absence of these
208 signals in the Hi-C based assembly suggests that the assembly was falsely merging
209 sequences from different haplotypes and the contigs were likely to be scaffolded in the wrong
210 order.

211 To test if the Hi-C based assemblies were truly a mixture of the two haplotypes, we
212 checked the presence of parental-specific *k*-mers within each of the 16 haplotype-specific
213 chromosome-level assemblies (Fig. 4c). This revealed that the majority of the haplotype
214 specific assemblies were in fact mixtures of the two haplotypes, which is in great contrast

215 with the high haplotyping accuracy of gamete binning. Finally, a whole-genome alignment of
216 the Hi-C based assembly to the genetic map based assembly of gamete binning revealed
217 many ambiguities between the genetic maps and the Hi-C based assembly within essentially
218 all haplotype-specific chromosome assemblies (Fig. 4d).

219 Taken together, besides its broad application, Hi-C-based phasing and scaffolding
220 was far from being error-free. Some of the errors combined large pieces from different
221 haplotypes, which resulted in falsely arranged chromosomes and severe phasing errors.
222 Though, gamete-binning might be more tedious in its experimental requirements, the
223 improved assembly quality might justify the additional effort.

224 **Haplotype diversity and (non-allelic) meiotic recombination**

225 In contrast to conventional diploid genome assemblies where the two haplotypes are
226 merged into one artificial consensus sequence, separate haploid assemblies allow for the
227 analysis of haplotype diversity. Comparing the two haplotype assemblies of 'Rojo Pasion'
228 using *SyRI*²³ allowed us to gain first insights into the haplotype diversity within an individual
229 apricot tree. Despite high levels of synteny, the two assemblies revealed large-scale
230 rearrangements (23 inversions, 1,132 translocation/transpositions and 2,477 distal
231 duplications) between the haplotypes making up more than 15% of the assembled sequence
232 (38.3 and 46.2 Mb in each of assemblies; Supplementary Table 1). Using the Hi-C contact
233 maps (Fig. 3e; Supplementary Figures 9-15), we validated the 17 largest rearrangements (>
234 500 kb) between the haplotype assemblies. Using a comprehensive RNA-seq dataset
235 sequenced from multiple tissues of 'Rojo Pasion' including reproductive buds, vegetative
236 buds, flowers, leaves, fruits (seeds removed) and barks as well as a published apricot RNA-
237 seq dataset³⁷, we predicted 30,378 and 30,661 protein-coding genes within each of the
238 haplotypes (with an annotation completeness of 96.4% according to a BUSCO⁴² analysis).
239 Mirroring the huge differences in the sequences, we found the vast amount of 942 and 865
240 expressed, haplotype-specific genes in each of the haplotypes (Methods; Supplementary
241 Tables 2-3). Such deep insights into the differences between the haplotypes, which are only

242 enabled by chromosome-level and haplotype-resolved assemblies, will generally be of high
243 value for the analysis of agronomically relevant variation.

244 Moreover, the chromosome-level assemblies also allow for fine-grained analyses of
245 the haploid pollen genomes, which have already undergone recombination during meiosis.
246 Meiotic recombination is the major mechanism to generate novel variation in offspring
247 genomes. During meiosis new haplotypes are formed by sequence exchanges between two
248 homologous chromosomes. To keep chromosome structures intact during such exchanges, it
249 is essential that recombination only occurs in syntenic regions as otherwise large parts of the
250 chromosome can be lost or duplicated in the newly formed molecules. Re-analyzing the 445
251 pollen nuclei genomes using one of the chromosome-level assemblies as reference, we
252 detected 2,638 meiotic crossover (CO) events (Methods). To improve the resolution of the
253 predicted CO events (6.1 kb), we selected 2,236 CO events detected in 369 nuclei with a
254 sequencing depth above 0.1x genome coverage (Supplementary Table 4). Along the
255 chromosomes, CO events were broadly and positively correlated with the density of protein-
256 coding genes and were almost completely absent in rearranged regions as expected (Fig. 5;
257 Methods). By investigating the fine-scale pattern of short read alignment of each nuclei, we
258 identified six CO events located in rearranged regions (0.3% of 2,236 CO events found in
259 1.6% of the pollen genomes), which led to stark chromosomal rearrangements. In each of
260 the six chromosomes we found duplicated read coverage and pseudo-heterozygous variation
261 in the regions that were involved in the chromosome rearrangements as induced by the non-
262 allelic CO (Fig. 6). This evidences the existence of non-allelic recombination in pollen
263 genomes and might open up a more detailed view on the actual meiotic recombination
264 patterns as compared to what could be observed in offspring individuals.

265 **Conclusion**

266 Taken together, following the elegant rationale of haplotype-based read separation
267 before genome assembly introduced by trio binning², we present gamete binning. In contrast
268 to trio binning, gamete binning does not rely on paternal genomes, but instead uses the

269 genomes of individual gametes to resolve haplotypes. In addition, the recombination patterns
270 in these gamete genomes can be used to calculate a genetic map, which in turn enables the
271 generation of chromosome-level assemblies. High-throughput analysis of gamete genomes
272 avoids tedious generation of offspring progeny and allows to sample the required material in
273 its ecological context, which makes it possible to analyze meiotic recombination as it occurs
274 in natural environments. As a result, gamete binning can efficiently and effectively enable
275 haplotype-resolved and chromosome-level genome assembly of any heterozygous individual
276 with accessible gametes.

277 **Online Methods**

278 **DNA extraction, Illumina/PacBio library preparation and sequencing**

279 Fresh developing fruits of 'Rojo Pasion' were frozen in liquid nitrogen immediately
280 after being sampled in Murcia, Spain. After being shipped to the Max Planck Institute for
281 Plant Breeding Research (MPIPZ, Cologne, Germany), DNA was extracted from the
282 mesocarp and exocarp of the fruits using the Plant DNA Kit of Macherey-Nagel™ to create a
283 PacBio sequencing library. Meanwhile, fresh leaves were sampled from the parental cultivars
284 ('Currot' and 'Orange Red') at the experimental field of CEBAS-CSIC in Murcia, Spain, and
285 Illumina short read libraries were prepared after DNA extraction using the Plant DNA Kit of
286 Macherey-Nagel™.

287 All libraries were sequenced with the respective sequencing machines (Illumina
288 HiSeq 3000 and PacBio Sequel I) at Max Planck Genome-centre Cologne (MP-GC), which
289 led to 19.9 Gb long reads for 'Rojo Pasion' (PacBio; Supplementary Figure 2) and 15.7 and
290 16.2 Gb short reads for the parental cultivars (Illumina). Note that the parental WGS data
291 were only used for haplotype validation and for sorting the individual chromosome
292 assemblies to two sets of eight chromosomes to match the inheritance of the chromosomes.

293 **Pollen nuclei DNA extraction, 10x sc-CNV library preparation and sequencing**

294 Dormant shoots of 'Rojo Pasion' bearing developed flower buds were collected in
295 Murcia, Spain. Then, the shoots were shipped at 4 °C to MPIPZ (Cologne, Germany) and
296 were grown in long-day conditions in the greenhouse. Flowers at the pre-anthesis stage were
297 frozen in liquid nitrogen. Anthers from ten 'Rojo Pasion'²⁸ flowers were extracted with forceps
298 and submerged in woody pollen buffer (WPB)⁴³. Around 500,000 pollen grains were
299 extracted from anthers by vortexing them in WPB. The nuclei were isolated from the pollen
300 using a modified bursting method³¹. Isolated pollen was prefiltered (100µm) and bursted
301 (30um) using Celltrics™ sieves and woody pollen buffer. The nuclei were then stained with
302 propidium iodide (PI) at 50 µg/mL just before sorting and counting by flow cytometry to
303 remove pollen grain debris using a BD FACSAria Fusion™ with high-speed sort settings (70

304 μm nozzle and 70 PSI sheath pressure) and 0.9% NaCl as sheath fluid. The nuclei were
305 identified by PI fluorescence, light scattering, and autofluorescence characteristics
306 (Supplementary Figure 4). A total of 12,600 nuclei were counted and collected in a solution
307 of 4.2 μL phosphate-buffered saline with 0.1% bovine serum albumin.

308 According to manufacturer's instructions, the nuclei were loaded into a 10xTM
309 Chromium controller in two batches with 6,300 nuclei each, i.e., two 10x sc-CNV libraries
310 were prepared. In each library, DNA fragments from the same nucleus were ligated with a
311 unique 16-bp barcode sequence (of A/C/G/T). Both libraries were sequenced using Illumina
312 HiSeq3000 in the 2x151 bp paired-end mode, totaling 95 and 124 million read pairs,
313 respectively (61.7 Gb).

314 **Hi-C library preparation and sequencing**

315 Approximately 0.5 grams of flash-frozen leaf samples of 'Rojo Pasion', which were
316 collected from the field, were thawed and fixed with 1% formaldehyde for 30 min at room
317 temperature under vacuum. Subsequently, the in situ Hi-C library preparation was performed
318 according to a protocol established for rice seedlings⁴⁴. The libraries were sequenced on an
319 Illumina HiSeq3000 instrument; in total, around 191.2 million pair-end reads were obtained.

320 **RNA extraction and sequencing**

321 Fruits tissue was collected in the same way for the PacBio sequencing library. Tissue
322 from reproductive buds, vegetative buds, flowers, leaves, bark tissues were collected from
323 the same shoots used for pollen nuclei isolation. RNA was extracted from these tissues using
324 the NucleoSpin® RNA Plant of Macherey-NagelTM to prepare Illumina libraries.

325 All libraries were sequenced with Illumina HiSeq 3000 at Max Planck Genome-centre
326 Cologne (MP-GC) in the 150 bp single-end mode, which respectively led to 32.8
327 (reproductive buds), 28.9 (vegetative buds), 30.2 (flowers), 23.8 (leaves), 18.6 (fruit) and
328 26.1 (bark) million reads, totaling 24.1 Gb.

329 **Genome size estimation**

330 After trimming off 10x Genomics barcodes and hexamers from the 61.7 Gb reads of
331 the two 10x sc-CNV libraries, *k*-mer counting ($k=21$) was performed with *Jellyfish*⁴⁵. The *k*-
332 mer histogram was provided to *findGSE*³⁰ to estimate the size of the ‘Rojo Pasion’ genome
333 under the heterozygous mode (with ‘*exp_hom*=200’; Supplementary Figure 3).

334 **Preliminary diploid-genome assembly and curation**

335 With the 19.9 Gb raw PacBio reads of ‘Rojo Pasion’ (Supplementary Figure 2), a
336 preliminary diploid assembly was constructed using *canu*²⁹ (with options
337 ‘*genomeSize*=242500000 *corMhapSensitivity*=high *corMinCoverage*=0 *corOutCoverage*=100
338 *correctedErrorRate*=0.105’).

339 All raw Illumina reads from the 10x libraries were firstly aligned to the initial assembly
340 using *bowtie2*⁴⁶. Then the *purge haplotigs* pipeline was then used to remove haplotigs (i.e.,
341 haplotype-specific contigs inflating the true haploid genome) based on statistical analysis of
342 sequencing depth, and identify primary contigs to build up a curated haploid assembly⁴⁷. To
343 reduce the false positive rate in defining haplotigs, each haplotig was blasted to the curated
344 assembly; if over 50% of the haplotig could not be covered by any primary contigs, it was re-
345 collected as a primary contig.

346 **SNP marker selection**

347 After trimming off 10x barcodes and hexamers, all pooled Illumina reads from the 10x
348 sc-CNV libraries (61.7 Gb) were re-aligned to the curated haploid assembly using *bowtie2*⁴⁶.
349 With 87.2% reads aligned, 989,132 raw SNPs were called with *samtools* and *bcftools*⁴⁸.
350 Three criteria were used to select potential allelic SNPs (578,209), including i) the alternative
351 allele frequency must be between 0.38 to 0.62; ii) the alternative allele must be carried by 60-
352 140 reads; iii) the total sequencing depth at a SNP must be between 120-280x (as compared
353 with genome-wide mode depth of 208x; Fig. 2b).

354 **Deletion marker selection and genotyping**

355 The assemblies included 10,452 regions of over 2 kb without SNP marker (total:
356 110.9 Mb). If the average sequencing depth of such regions was less than or equal to 146x
357 (i.e., the value at the valley between middle and right-most peaks in the sequencing depth
358 distribution; Fig. 2b), it was selected as a deletion-like marker. This revealed 3,253 deletion
359 markers (36.4 Mb), including 237 on contigs without a single SNP marker. The remaining
360 7,199 regions (74.5 Mb) were defined as conserved (homozygous regions) between two
361 haplotypes (Fig. 2b). For each deletion marker and gamete genome, we assessed the
362 normalized read (*RPKM* value) could of the reads aligned within the deletion using
363 *bedtools*⁴⁹. The genotype at such a deletion marker was initialized as *a* or *n*, where *a* refers
364 to the presence of reads (and therefore relates to the haplotype without the deletion) and *n*
365 refers to the absence of reads (either the deletion haplotype or not having enough
366 information).

367 **Haplotype phasing and CO identification**

368 Barcodes in the raw reads were corrected using *cellranger*, with which 182.1 million
369 read pairs (51.0 Gb) were clustered into 691 read sets. Reads of each read set were aligned
370 to the curated assembly using *bowtie2*⁴⁶, bases were called using *bcftools*⁵⁰, and a simple bi-
371 marker majority voting strategy was applied to phase the SNPs along each contig (Fig. 2c).
372 After phasing, we identified COs as consistent switches between the haplotypes.

373 **Ploidy evaluation of single-cell sequencing**

374 For each nucleus, with short read alignment and base calling to the curated
375 assembly, we counted the number of inter-genotype transitions (genotype *a* to *b* and *b* to *a*)
376 at phased SNP markers over all contigs. Correlating this to the number of covered markers
377 revealed two clusters of nuclei (Supplementary Figure 6c). One cluster with 217 nuclei
378 showed that inter-genotype transitions increased linearly with the number of covered markers
379 (while there were high ratios of more than 5 transitions in every 100 markers), which
380 indicated the sequencing data were mixed from more than one nuclei. The other cluster of

381 445 nuclei (31.2 Gb with 111.4 million read pairs) showed a limited increase (probably due to
382 sequencing errors or markers from repetitive regions), which supported the expected haploid
383 status.

384 **Imputation of virtual markers at ends of contigs**

385 Let a and b denote the parental genotypes. The genotype of a nucleus at both ends
386 of a contig (referred to as virtual markers) can be represented by aa , bb or ab (or ba) where
387 aa/bb indicates an identical genotype along the contig while ab (or ba) indicates a CO event
388 in the regions of contig. Then we can build up genotype sequences at the two ends of all
389 contigs (with SNP markers) by imputing at all nuclei. For example, given a contig, sequences
390 of $aaaaa**ab**bbbbbbb$ (marker 1) and $aaaaaa**ab**bbbbbbb$ (marker 2) means there is a CO (in
391 bold) at the 7th (of 15) nuclei (Fig. 2c).

392 **Linkage grouping and genetic mapping**

393 All virtual markers (defined using SNP markers along contigs) were classified into 8
394 linkage groups (653 contigs: 212.9 Mb) after pairwise comparison of their genotype
395 sequences using *JoinMap4.0*³³ (with haploid population type: HAP; and logarithm of the odds
396 (LOD) values larger than 3.0).

397 After filtering out contigs with >10% missing nuclei information or nuclei with >10%
398 missing contigs, a high-quality genetic map consisting of 216 contigs (147.3 Mb,
399 corresponding to 622.0 cM; Fig. 3a) was first obtained using regression mapping in *JoinMap*
400 4.0® with the following settings: LOD larger than 3.0, a “goodness-of-fit jump” threshold of
401 5.0 for removal of loci and a “two rounds” mapping strategy³³. Genotype sequences imputed
402 at contig ends or deletions (i.e., respective virtual markers) were used to integrate the
403 remaining 723 contigs into the genetic map. For example, given a deletion marker (e.g., p
404 and q in Fig. 2c-e), if the respective contig had already existed in the genetic map, phasing
405 was only performed at the deletion (according to surrounding phased SNPs); otherwise,
406 phasing plus positioning to the genetic map would be applied. Both operations were based
407 on finding the minimum divergence of the genotype sequence of the marker to that of the

408 other contigs (in the corresponding genetic map). The final genetic map was completed as
409 891 contigs of 228.0 Mb.

410 **Haplotype-specific PacBio read separation**

411 PacBio reads (19.9 Gb) were classified based on three major cases after being
412 aligned to the curated assembly using *minimap2*⁵¹. First, a read covering phased SNP
413 markers was directly clustered into the haplotype supported by the respective alleles in the
414 read. Second, a read covering no SNP markers but overlapping a deletion marker was
415 clustered into the respective genotype based on its phasing with neighboring imputed
416 markers in genetic map. Third, a read in a conserved region was assigned to one of the
417 haplotypes randomly.

418 **Haplotype assembly and chromosome-level scaffolding**

419 Independent assemblies were performed with sixteen sets of reads, i.e., for every two
420 haplotypes in each of the eight linkage groups using *Flye*³⁴ with the default settings.

421 Using the 891 contigs of the curated assembled that were assigned to chromosomal
422 positions with the genetic mapping, we created a pseudo reference genome, with which the
423 newly assembled contigs were scaffolded using *RAGOO*⁵², leading to chromosome-level
424 assemblies (i.e., those labeled with 'scaf' in Fig. 3b).

425 **Haplotype evaluation**

426 The genotypes of the sixteen assemblies were firstly identified by comparing *k*-mers
427 in each assembly with Illumina WGS of the parental cultivar ($k=21$; Fig. 3c). Although
428 evaluation can always be performed in each linkage group, we combined the eight linkage-
429 group-wise assemblies for 'Currot'-genotype and the other eight for 'Orange Red'-genotype,
430 respectively.

431 After polishing the assemblies respectively with the 'Currot'-genotype and 'Orange
432 Red'-genotype PacBio reads using *apollo*⁵³, we built up two sets of haplotype-specific *k*-mers
433 from the assemblies, r_C and r_O . Correspondingly, a set of 'Currot'-specific *k*-mers (with

434 coverage from 10 to 60x), p_C , was selected from the parental Illumina WGS that did not exist
435 in ‘Orange Red’ short reads (coverage over 1x) but in ‘Rojo Pasion’ pollen short reads
436 (coverage from 10 to 300x); similarly, a set of ‘Orange Red’-specific k -mers, p_O , was also
437 collected. Then we intersected r_C and r_O with p_C and p_O respectively, leading to four subsets
438 $r_C \cap p_C$, $r_C \cap p_O$, $r_O \cap p_C$, and $r_O \cap p_O$, which were used to calculate average haplotyping accuracy.
439 All k -mer processing (counting, intersecting and difference finding) were performed with
440 *KMC*⁵⁴. After haplotype validation, the assemblies were further polished with the respective
441 parental short read alignment using *pilon*⁵⁵ (with options ‘--fix bases --mindepth 0.85’)
442 generating v1.0 of the assemblies. Manual correction of the v.1.0 assemblies was performed
443 according to focal and parental reads to generate assembly v1.1. Finally, k -mer-based
444 assembly validation was performed with KAT³⁵ and Merqury³⁶.

445 **Genome annotation**

446 We annotated protein-coding genes for each haplotype assembly (v1.0) by integrating
447 evidences from *ab initio* gene predictions (using three tools *Augustus*⁵⁶, *GlimmerHMM*⁵⁷ and
448 *SNAP*⁵⁸), RNA-seq read assembled transcripts and homologous protein sequence
449 alignments. We aligned protein sequences from the database UniProt/Swiss-Prot,
450 *Arabidopsis thaliana* and *Prunus persica* to each haplotype assembly using the tool
451 *Exonerate*⁵⁹ with the options “--percent 60 --minintron 10 --maxintron 60000”. We mapped
452 RNA-seq reads from reproductive buds, vegetative buds, flowers, leaves, fruits (except
453 seeds) and bark tissues, as well as a published Apricot RNA-seq dataset³⁷, using HISAT⁶⁰,
454 and we assembled the transcripts using *StringTie*⁶¹. Finally, we used the tool
455 *EvidenceModeler*⁶² to integrate the above evidence in order to generate consensus gene
456 models for each haplotype assembly.

457 We annotated the transposon elements (TE) using the tools *RepeatModeler* and
458 *RepeatMasker* (<http://www.repeatmasker.org>). We filtered the TE related genes based on
459 their coordinates overlapping with TEs (overlapping percent > 30%), sequence alignment

460 with TE-related protein sequences and *A. thaliana* TE related gene sequences (both
461 requiring *blastn* alignment identity and coverage both larger than 30%).

462 We improved the resulting gene models using in-house scripts. Firstly, we ran a
463 primary gene family clustering using *orthoFinder*⁶³ based on the resulting gene models from
464 each haplotype to find haplotype-specific genes. We then aligned these specific gene
465 sequences to the other haplotype using *blastn*⁶⁴ to check whether it was specific because the
466 ortholog was unannotated in the other haplotype. For these potentially unannotated genes
467 (*blastn* identity > 60% and *blastn* coverage > 60%), we checked the gene models from *ab*
468 *initio* prediction around the aligned regions to add the unannotated gene if both the gene
469 model and the aligned region had an overlapping rate larger than 80%. We also directly
470 generated new gene models based on the *Scipio*⁶⁵ alignment after confirming the existence
471 of start codon, stop codon and splicing site. Finally, the completeness of assembly and
472 annotation were evaluated by the *BUSCO*⁴² v4 tool based on 2,326 eudicots single-copy
473 orthologs from OrthoDB v10⁶⁶. A similar process was used to filter for haplotype-specific
474 genes (Supplementary Tables 2-3). Finally, a genome annotation lift-over was performed
475 from v1.0 to v1.1 using *liftoff*⁶⁷ with default parameters.

476 **Genome assembly comparison**

477 All genome assemblies, including 'Rojo Pasion' haplotypes, 'Chuanzhihong' apricot
478 (*Prunus armeniaca*)³⁷, Japanese apricot (*Prunus mume*)³⁸ and 'Lovell' peach (*Prunus*
479 *persica*)³⁹, were aligned to each other using *nucmer* from the *MUMmer4*⁶⁸ toolbox with
480 parameters '-max -l 40 -g 90 -b 100 -c 200'. The alignments were further filtered for
481 alignment length (>100 bp) and identity (>90%), with which structural rearrangements and
482 local variations were identified using *SyR*²³. To follow the nomenclature of the *Prunus*
483 community, the 'Rojo Pasion' chromosomes were numbered according to the numbering in
484 'Lovell' peach³⁹.

485 **Hi-C data analysis**

486 We used *ALLHiC*⁸ and *SALSA2*⁴⁰ for phasing and scaffolding. All 191.2 million Hi-C
487 read pairs were aligned (using *BWA* version 0.7.15-r1140) to the haplotype-resolved unitigs
488 assembled by *Canu*. Only uniquely mapped read pairs were selected using
489 *filterBAM_forHiC.pl* from the *ALLHiC* package. The selected alignments were used as input
490 for *ALLHiC_partition* (“*ALLHiC_partition* -b clean.bam -r unitigs.fa -e GATC -k 19”) and
491 *SALSA2* (“python run_pipeline.py -a unitigs.fa -l unitigs.fa.fai -g unitigs.gfa -m yes -b
492 alignment.bed -e GATC -o SALSA2_out -i 8”, where the file alignment.bed was generated
493 and sorted from clean.bam using *bedtools bamtobed* (version v2.29.0) and unitigs.gfa was
494 collected from the *Canu* output). For *ALLHiC*, we had to set group number as 19 to get 16
495 linkage groups (of chromosome-level size), and 3 smaller groups below 2.5 Mb, which were
496 not considered further. We continued with *ALLHiC* pipeline as it provided more accurate than
497 those from *SALSA2*. The subsequent pipeline of *ALLHiC* were run by default except for using
498 “-RE GATC” in the “allhic extract” command. For comparison, we also aligned all raw Hi-C
499 reads to haploid assemblies generated by gamete binning, and selected the uniquely
500 mapped read pairs as described above. Hi-C maps were visualized using *ALLHiC_plot* at
501 300-500 kb resolution. Alignments of *ALLHiC* and gamete binning based assemblies were
502 obtained using *minimap2* and dot plot was drawn with script *paFCoordsDotPlotly.R* at
503 <https://github.com/tpoorten/dotPlotly>.

504 **Crossover identification**

505 All 220 million pollen nuclei-derived short read pairs were pooled and aligned to the
506 ‘Currot’-genotype assembly, from which 739,342 SNP markers were defined with an
507 alternative allele frequency distribution of 0.38 to 0.62 and alternative allele coverage of 50 to
508 150x. Then, short reads of 445 nuclei were independently aligned to the ‘Currot’-genotype
509 assembly using *bowtie2*⁴⁶ and bases were called using *bcftools*⁵⁰. Finally, *TIGER*⁶⁹ was used
510 to identify COs. The landscape of COs from 369 nuclei with a sequencing depth over 0.1x
511 was calculated within 500 kb sliding windows along each chromosome at a step of 50 kb

512 (Fig. 5), where for each window, the recombination frequency (cM/Mb) was defined as
513 $C/n/(w/10^6) * 100\%$, where C is the number of recombinant nuclei in that window, n is the
514 total number of nuclei (369) and w is the window size. SNP/Mb and $gene/Mb$ were calculated
515 for the same windows as $x/(w/10^6)$, where x was the count of the feature in the respective
516 window.

517 **Acknowledgements**

518 The authors would like to thank Antonio Molina and José Egea for kindly providing
519 plant material, Saurabh Pophaly for help in transferring the read data to public servers, Detlef
520 Weigel for supportive guidance, and Kristin Krause and Vidya Oruganti for useful discussions
521 and comments for improving the manuscript. This work was funded by the “Humboldt
522 Research Fellowship for Experienced Researchers” (Alexander von Humboldt Foundation)
523 (J.A.C.), the Marie Skłodowska-Curie Individual Fellowship PrunMut (789673) (J.A.C.), the
524 Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s
525 Excellence Strategy – EXC 2048/1– 390686111 (K.S.), and the European Research Council
526 (ERC) Grant “INTERACT” (802629) (K.S.). C.K. acknowledges the ISAC SRL Emerging
527 Leaders Program.

528 **Author contributions**

529 J.A.C, H.S. and K.S. designed the project. J.A.C., B.H., K. F.-D., C.K., D.R., M.R.,
530 N.W. and C.L. performed wet-lab experiments. H.S., J.A.C, M.G., and W-B.J. performed all
531 data analysis. J.A.C., H.S. and K.S. wrote the manuscript with input from all authors. All
532 authors read and approved the final manuscript.

533 **Competing interests**

534 The authors declare no competing interests.

535 **Data availability**

536 Data supporting the findings of this work are available within the paper and its
537 Supplementary Information files. Read data sequenced from two 10x sc-CNV libraries, two
538 Hi-C libraries, one PacBio library from ‘Rojo Pasion’, two Illumina libraries for ‘Currot’ and
539 ‘Orange Red’ that support the work in this study as well as the haploid assemblies and
540 annotations generated are available in European Nucleotide Archive (ENA) under accession
541 number “PRJEB37669”. Data was uploaded to ENA using EMBLmyGFF⁷⁰. All other relevant
542 data are available upon request.

543 **Code availability**

544 Customs scripts supporting this work are available at github.com/schneeberger-

545 [lab/GameteBinning](#).

546 REFERENCES

- 547 1. Korfach, J. *et al.* De novo PacBio long-read and phased avian genome assemblies
548 correct and add to reference genes generated with intermediate and short reads.
549 *Gigascience* **6**, 1–16 (2017).
- 550 2. Koren, S. *et al.* De novo assembly of haplotype-resolved genomes with trio binning.
551 *Nat. Biotechnol.* **36**, 1174–1182 (2018).
- 552 3. Yang, H., Chen, X. & Wong, W. H. Completely phased genome sequencing through
553 chromosome sorting. *Proc. Natl. Acad. Sci. U. S. A.* **108**, 12–17 (2011).
- 554 4. Falconer, E. & Lansdorp, P. M. Strand-seq: A unifying tool for studies of chromosome
555 segregation. *Semin. Cell Dev. Biol.* **24**, 643–652 (2013).
- 556 5. Hills, M. *et al.* Construction of whole genomes from scaffolds using single cell strand-
557 seq data. *bioRxiv* (2018). doi:<https://doi.org/10.1101/271510>
- 558 6. Chaisson, M. J. P. *et al.* Multi-platform discovery of haplotype-resolved structural
559 variation in human genomes. *Nat. Commun.* **10**, 1–16 (2019).
- 560 7. Selvaraj, S., Dixon, J. R., Bansal, V. & Ren, B. Whole-genome haplotype
561 reconstruction using proximity-ligation and shotgun sequencing. *Nat. Biotechnol.* **31**,
562 1111–1118 (2013).
- 563 8. Zhang, X., Zhang, S., Zhao, Q., Ming, R. & Tang, H. Assembly of allele-aware,
564 chromosomal-scale autopolyploid genomes based on Hi-C data. *Nat. Plants* **5**, 833–
565 845 (2019).
- 566 9. Garg, S. *et al.* Efficient chromosome-scale haplotype-resolved assembly of human
567 genomes. *bioRxiv* 810341 (2019). doi:10.1101/810341
- 568 10. Linsmith, G. *et al.* Pseudo-chromosome-length genome assembly of a double haploid
569 ‘bartlett’ pear (*Pyrus communis* L.). *Gigascience* **8**, 1–17 (2019).
- 570 11. Edger, P. P. *et al.* Origin and evolution of the octoploid strawberry genome. *Nat.*
571 *Genet.* **51**, 541–547 (2019).
- 572 12. Bickhart, D. M. *et al.* Single-molecule sequencing and chromatin conformation capture
573 enable de novo reference assembly of the domestic goat genome. *Nat. Genet.* **49**,

- 574 643–650 (2017).
- 575 13. Wallberg, A. *et al.* A hybrid de novo genome assembly of the honeybee, *Apis*
576 *mellifera*, with chromosome-length scaffolds. *BMC Genomics* **20**, 1–19 (2019).
- 577 14. Doležel, J. *et al.* Advances in plant chromosome genomics. *Biotechnol. Adv.* **32**, 122–
578 136 (2014).
- 579 15. The International Wheat Genome Sequencing Consortium. A chromosome-based
580 draft sequence of the hexaploid bread wheat (*Triticum aestivum*) genome. *Science*
581 (80-). **345**, (2014).
- 582 16. Zhang, J. *et al.* Allele-defined genome of the autopolyploid sugarcane *Saccharum*
583 *spontaneum* L. *Nat. Genet.* **50**, 1565–1573 (2018).
- 584 17. Garg, S. *et al.* Accurate chromosome-scale haplotype-resolved assembly of human
585 genomes. (2020). doi:10.1101/810341
- 586 18. Zhang, X., Wu, R., Wang, Y., Yu, J. & Tang, H. Unzipping haplotypes in diploid and
587 polyploid genomes. *Comput. Struct. Biotechnol. J.* **18**, 66–72 (2020).
- 588 19. Li, R. *et al.* Inference of Chromosome-length Haplotypes Using Genomic Data of
589 Three to Five Single Gametes. *bioRxiv* 361873 (2018). doi:10.1101/361873
- 590 20. Kirkness, E. F. *et al.* Sequencing of isolated sperm cells for direct haplotyping of a
591 human genome. *Genome Res.* **23**, 826–832 (2013).
- 592 21. Shi, D. *et al.* Single-pollen-cell sequencing for gamete-based phased diploid genome
593 assembly in plants. *Genome Res.* 1–11 (2019).
- 594 22. Wu, J. *et al.* The genome of the pear (*Pyrus bretschneideri* Rehd.). *Genome Res.* **23**,
595 396–408 (2013).
- 596 23. Goel, M., Sun, H., Jiao, W. B. & Schneeberger, K. SyRI: finding genomic
597 rearrangements and local sequence differences from whole-genome assemblies.
598 *Genome Biol.* **20**, 1–13 (2019).
- 599 24. Jiao, W. B. & Schneeberger, K. Chromosome-level assemblies of multiple *Arabidopsis*
600 genomes reveal hotspots of rearrangements with altered evolutionary dynamics. *Nat.*
601 *Commun.* **11**, 1–10 (2020).

- 602 25. Amarasinghe, S. L. *et al.* Opportunities and challenges in long-read sequencing data
603 analysis. *Genome Biol.* **21**, 1–16 (2020).
- 604 26. Sun, H. *et al.* Linked-read sequencing of gametes allows efficient genome-wide
605 analysis of meiotic recombination. *Nat. Commun.* **10**, 1–9 (2019).
- 606 27. Dréau, A., Venu, V., Avdievich, E., Gaspar, L. & Jones, F. C. Genome-wide
607 recombination map construction from single individuals using linked-read sequencing.
608 *Nat. Commun.* **10**, (2019).
- 609 28. Egea, J., Dicenta, F. & Burgos, L. 'Rojo Pasion' apricot. *Hortscience* **39**, 1490–1491
610 (2004).
- 611 29. Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer
612 weighting and repeat separation. *Genome Res.* **27**, 722–736 (2017).
- 613 30. Sun, H., Ding, J., Piednoël, M. & Schneeberger, K. FindGSE: Estimating genome size
614 variation within human and Arabidopsis using k -mer frequencies. *Bioinformatics* **34**,
615 550–557 (2018).
- 616 31. Kron, P. & Husband, B. C. Using flow cytometry to estimate pollen DNA content:
617 Improved methodology and applications. *Ann. Bot.* **110**, 1067–1078 (2012).
- 618 32. Julian, C., Rodrigo, J. & Herrero, M. Stamen development and winter dormancy in
619 apricot (*Prunus armeniaca*). *Ann. Bot.* **108**, 617–625 (2011).
- 620 33. van Ooijen, J. W. JoinMap® 4, Software for the calculation of genetic linkage maps in
621 experimental populations. Wageningen, Netherlands: Kyazma B.V. (2006).
- 622 34. Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. a. Assembly of long, error-prone
623 reads using repeat graphs. *Nat. Biotechnol.* **37**, 540–546 (2019).
- 624 35. Mapleson, D., Garcia Accinelli, G., Kettleborough, G., Wright, J. & Clavijo, B. J. KAT: a
625 K-mer analysis toolkit to quality control NGS datasets and genome assemblies.
626 *Bioinformatics* **33**, 574–576 (2017).
- 627 36. Rhie, A., Walenz, B., Koren, S. & Phillippy, A. Merqury: reference-free quality,
628 completeness, and phasing assessment for genome assemblies. *bioRxiv* (2020).
629 doi:10.1101/2020.03.15.992941

- 630 37. Jiang, F. *et al.* The apricot (*Prunus armeniaca* L.) genome elucidates Rosaceae
631 evolution and beta-carotenoid synthesis. *Hortic. Res.* **6**, 1–12 (2019).
- 632 38. Zhang, Q. *et al.* The genome of *Prunus mume*. *Nat. Commun.* **3**, 1–8 (2012).
- 633 39. Verde, I. *et al.* The high-quality draft genome of peach (*Prunus persica*) identifies
634 unique patterns of genetic diversity, domestication and genome evolution. *Nat. Genet.*
635 **45**, 487–494 (2013).
- 636 40. Ghurye, J. *et al.* Integrating Hi-C links with assembly graphs for chromosome-scale
637 assembly. *PLoS Comput. Biol.* **15**, 1–19 (2019).
- 638 41. Chen, H. *et al.* Allele-aware chromosome-level genome assembly and efficient
639 transgene-free genome editing for the autotetraploid cultivated alfalfa. *Nat. Commun.*
640 **11**, (2020).
- 641 42. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M.
642 BUSCO: Assessing genome assembly and annotation completeness with single-copy
643 orthologs. *Bioinformatics* **31**, 3210–3212 (2015).
- 644 43. Loureiro, J. Two new nuclear isolation buffers for plant DNA flow cytometry: A test with
645 37 species. *Ann. Bot.* 875–888 (2007).
- 646 44. Liu, C., Cheng, Y. J., Wang, J. W. & Weigel, D. Prominent topologically associated
647 domains differentiate global chromatin packing in rice from *Arabidopsis*. *Nat. Plants* **3**,
648 742–748 (2017).
- 649 45. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of
650 occurrences of k-mers. *Bioinformatics* **27**, 764–770 (2011).
- 651 46. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient
652 alignment of short DNA sequences to the human genome. *Genome Biol.* **10**, 1–10
653 (2009).
- 654 47. Roach, M. J., Schmidt, S. a. & Borneman, A. R. Purge Haplotigs: Allelic contig
655 reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics* **19**, 1–10
656 (2018).
- 657 48. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,

- 658 2078–2079 (2009).
- 659 49. Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for comparing
660 genomic features. *Bioinformatics* **26**, 841–842 (2010).
- 661 50. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping
662 and population genetical parameter estimation from sequencing data. *Bioinformatics*
663 **27**, 2987–2993 (2011).
- 664 51. Li, H. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* **34**,
665 3094–3100 (2018).
- 666 52. Alonge, M. *et al.* RaGOO: Fast and accurate reference-guided scaffolding of draft
667 genomes. *Genome Biol.* **20**, 1–17 (2019).
- 668 53. Firtina, C. *et al.* Apollo: A Sequencing-Technology-Independent, Scalable, and
669 Accurate Assembly Polishing Algorithm. *Bioinformatics* 1–10 (2020).
- 670 54. Kokot, M., Dlugosz, M. & Deorowicz, S. KMC 3: counting and manipulating k-mer
671 statistics. *Bioinformatics* **33**, 2759–2761 (2017).
- 672 55. Walker, B. J. *et al.* Pilon: An integrated tool for comprehensive microbial variant
673 detection and genome assembly improvement. *PLoS One* **9**, 1–14 (2014).
- 674 56. Stanke, M. *et al.* AUGUSTUS: Ab initio prediction of alternative transcripts. *Nucleic*
675 *Acids Res.* **34**, 435–439 (2006).
- 676 57. Majoros, W. H., Pertea, M. & Salzberg, S. L. TigrScan and GlimmerHMM: Two open
677 source ab initio eukaryotic gene-finders. *Bioinformatics* **20**, 2878–2879 (2004).
- 678 58. Johnson, A. D. *et al.* SNAP: A web-based tool for identification and annotation of proxy
679 SNPs using HapMap. *Bioinformatics* **24**, 2938–2939 (2008).
- 680 59. Slater, G. S. C. & Birney, E. Automated generation of heuristics for biological
681 sequence comparison. *BMC Bioinformatics* **6**, 1–11 (2005).
- 682 60. Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low
683 memory requirements Daehwan HHS Public Access. *Nat. Methods* **12**, 357–360
684 (2015).
- 685 61. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from

- 686 RNA-seq reads. *Nat. Biotechnol.* **33**, 290–295 (2015).
- 687 62. Haas, B. J. *et al.* Automated eukaryotic gene structure annotation using
688 EVIDENCEModeler and the Program to Assemble Spliced Alignments. *Genome Biol.* **9**,
689 1–22 (2008).
- 690 63. Emms, D. M. & Kelly, S. OrthoFinder: Phylogenetic orthology inference for
691 comparative genomics. *Genome Biol.* **20**, 1–14 (2019).
- 692 64. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
693 alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 694 65. Keller, O., Odrionitz, F., Stanke, M., Kollmar, M. & Waack, S. Scipio: Using protein
695 sequences to determine the precise exon/intron structures of genes and their
696 orthologs in closely related species. *BMC Bioinformatics* **9**, 1–12 (2008).
- 697 66. Kriventseva, E. V. *et al.* OrthoDB v10: Sampling the diversity of animal, plant, fungal,
698 protist, bacterial and viral genomes for evolutionary and functional annotations of
699 orthologs. *Nucleic Acids Res.* **47**, D807–D811 (2019).
- 700 67. Shumate, A. & Salzberg, S. L. Liftoff: an accurate gene annotation mapping tool.
701 *bioRxiv* (2020). doi:10.1101/2020.06.24.169680
- 702 68. Marçais, G. *et al.* MUMmer4: A fast and versatile genome alignment system. *PLoS*
703 *Comput. Biol.* **14**, 1–14 (2018).
- 704 69. Rowan, B. A., Patel, V., Weigel, D. & Schneeberger, K. Rapid and inexpensive whole-
705 genome genotyping-by-sequencing for crossover localization and fine-scale genetic
706 mapping. *G3 Genes, Genomes, Genet.* **5**, 385–398 (2015).
- 707 70. Norling, M., Jareborg, N. & Dainat, J. EMBLmyGFF3: A converter facilitating genome
708 annotation submission to European Nucleotide Archive. *BMC Res. Notes* **11**, 1–5
709 (2018).
- 710

711 **Figure legends**

712 **Figure 1. Overview of gamete binning.** **a.** Extraction of gamete nuclei. **b.** Single-cell
713 genome sequencing of haploid gametes and haplotype phasing. **c.** Genetic map
714 construction based on the recombination patterns in the gamete genomes. **d.** Long-
715 read sequencing of somatic material. **e.** Separation of long reads based on genetic
716 linkage groups using phased alleles. **f.** Independent assembly of each haplotype of
717 each linkage group. **g.** Scaffolding assemblies to chromosome-level using the gamete-
718 derived genetic map.

719 **Figure 2. Single-pollen nuclei sequencing, variant phasing and genetic mapping.** **a.**
720 Sequencing depths of 691 pollen nuclei. **b.** Sequencing depth histogram of pooled
721 pollen short reads. The left-most peak revealed 0.9% of the genome that were not well
722 covered in the pollen read sets (i.e., $\leq 5x$). The middle peak indicated regions covered
723 only by half of the genomes and present in only one of the haplotypes, and the right-
724 most peak indicated regions, which were present in both haplotypes and showed the
725 expected coverage. In regions represented in both haplotypes, 578,209 SNPs were
726 defined. Regions without SNP markers were classified into 3,253 deletions and 7,199
727 conserved regions (Methods). **c.** SNP phasing along contigs. Genotyping was first
728 performed for each individual nuclei at each SNP marker. As shown, both genotypes
729 (in red and blue) were mixed in the curated but mosaic assembly. After phasing, 8 and
730 7 nuclei were respectively clustered for genotype A and B, and crossover could be
731 identified. With this, representative markers were imputed at ends of contigs. **d.**
732 Imputation of markers at deletions by genotyping using normalized read count. Two
733 cases were considered for phasing (and positioning) a deletion marker (in the genetic
734 map). If it was linked with surrounding SNP alleles, it could be phased accordingly;
735 otherwise, comparison its genotype sequence to genotype sequences of all other
736 markers (including SNP-derived markers at ends of contigs) would be performed to find
737 its value of phase (and positioning). **e.** Linkage group and genetic map construction

738 using the set of imputed markers (SNP-derived markers labeled as 1-8 and deletion
739 markers as p and q). For example, the genotype sequences of 6, 8 and q needed to be
740 flipped (i.e., phase values were 1 - contig phasing). Further ordering of the markers
741 (using *JoinMap*) led to linkage group-wise genetic maps.

742 **Figure 3. Genetic mapping, haplotype-specific assembly and validation.** **a.** Top: Genetic
743 map with a total genetic length of 622.0 cM (Methods). Middle: up to 2 Gb reads were
744 assigned to one of the two haplotypes of each linkage group. Bottom: a combination of
745 haplotype-A/B linkage groups led to two assemblies with 214.6 and 215.3 Mb. **b.**
746 Contig size distributions before (ctg-A, ctg-B) and after scaffolding (scaf-CU for the
747 assembly with sequence from 'Currot'; and scaff-OR for the assembly with sequence
748 from 'Orange Red'). After scaffolding, eight chromosome-scale pseudo-molecules were
749 obtained for each haplotype as labeled by "Chrs". **c.** Haplotype validation for the two
750 assemblies of each linkage group (LG-1-8) using parent-specific k -mers (of 'Orange
751 Red' and 'Currot'). With each linkage group, the two assemblies could be clearly
752 identified as either 'Currot'-haplotype or 'Orange Red'-haplotype using parental k -mers.
753 After combining the 'Currot'-related assemblies and 'Orange Red'-related assemblies
754 to genome-level, k -mer comparison revealed a haplotype accuracy of 99.1%. **d.** Using
755 the 'Currot'-haplotype as representative and comparing it to the assembly of the double
756 haploid *Prunus* ssp. reference genome (*Prunus persica*, and other closely-related
757 species; Supplementary Figure 6) revealed high levels of synteny and thus implies high
758 accuracy of the genetic map and chromosome-level scaffolding. **e.** Hi-C contact map
759 based on the assemblies of the two haplotypes of chromosome 1 (Currot (CU) and
760 Orange Red (OR)) at a resolution of 300 kb. The contact signal showed a high
761 contiguity within the haplotypes (main diagonal line) and confirmed two large inversions
762 (v_{11} and v_{12}) which we observed in the assembled sequence of the two haplotypes.
763 (See Supplementary figures 9-15 for chromosomes 2-8).

764 **Figure 4. Comparison of Hi-C based phasing and scaffolding with gamete binning. a.**
765 Hi-C contact map based on all 16 haplotype assemblies generated with gamete binning
766 (Currot (CU) and Orange Red (OR)). **b.** Hi-C contact map based on all 16 haplotype
767 assemblies generated with Hi-C data (Currot (CU) and Orange Red (OR)). Note the
768 contact signal along the main diagonal was much weaker as compared to the signal
769 based on the gamete binning assembly, and virtually no contact signals between two
770 different haplotypes could be identified. **c.** Haplotype validation of each linkage group
771 (LG-1-8) of the Hi-C assembly using parent-specific *k*-mers. Almost all haplotype-
772 specific assemblies included *k*-mers specific to both of the parental alleles indicating
773 severe errors in the phasing. **d.** Alignment of the Hi-C based assembly to genetic map
774 derived assembly (i.e. gamete binning derived assembly) revealed mis-joining or
775 splitting of linkage groups within the Hi-C assembly. For example, CUR1G was split
776 into Hi-C:19g1, 19g4; on the other hand, alignments of 19g2 to CUR2G, CUR5G and
777 CUR7G revealed mis-joins of sequences from independent linkage groups.

778 **Figure 5. Structural genome variations and meiotic recombination.** Top: recombination
779 landscape created with sliding windows of 500 kb at a step of 50 kb with COs detected
780 in all single pollen nuclei (with coverage over 0.1x), coupled with SNP density and
781 gene density. For x-axis, coordinates were based on the haploid assembly of ‘Currot’-
782 genotype. For y-axis, all features were scaled to 1.0, which stands for a maximum of
783 18 for recombination frequency (*cM/Mb*), 7,410 for SNP density and 480 for gene
784 density. Bottom: structural variations (>50 kb) identified between the two haploid
785 assemblies. In general, crossovers are almost completely absent in SVs, for example,
786 at LG2:11.0–14.5 Mb (inversion case) and LG5:16.0–18.2 Mb (translocation case).
787 Variants spanning over 500 kb are labelled as v_{xy} , where *x* denotes the chromosome
788 number and *y* the identifier of the variant in the chromosome. All these large variants
789 were confirmed within Hi-C contact maps (Fig. 3e, Supplementary Figures 9-16).

790 **Figure 6. Non-allelic crossovers and its consequences.** **a.** Illustration of a non-allelic
791 crossover which results in a chromosomal anomaly. **b.** Analysis of a single-pollen
792 nuclei, which revealed a non-allelic CO resulting in the duplication of a large
793 chromosomal segment. The short-read alignments of a haploid nuclei revealed a
794 pseudo-heterozygous region with increased read coverage, which is the hallmark of a
795 long duplication specific to this genome. All other chromosomes were haploid (not
796 shown). (Top row: ‘Currot’ allele frequency, SNP density (in sliding windows of 500 kb
797 at a step of 50 kb), and read coverage scaled by SNP density. Middle row: count of
798 ‘Currot’ or ‘Orange Red’ alleles at SNP markers. Bottom row: diagram illustrating how a
799 non-allelic CO in transposed regions (as indicated by yellow rectangles) resulted in a
800 large duplication, i.e., the original homologous chromosomal regions labelled with “4”
801 and “5” are now part of the same newly formed chromosome.

802 **Tables**

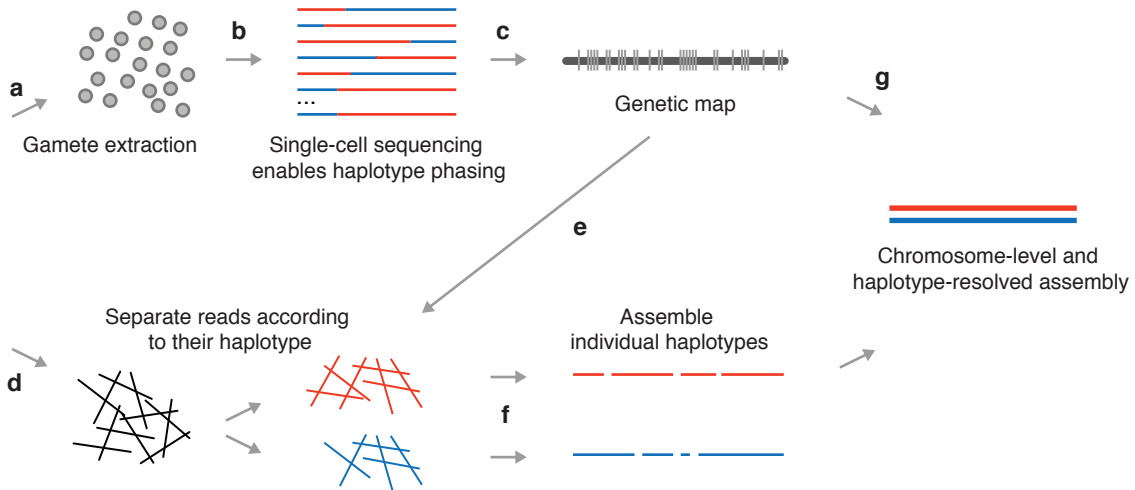
803 **Table 1 Assembly and validation statistics of two haplotype-resolved genome assemblies.** Note, the eight main chromosome-level scaffolds of
 804 each haplotype made up ~99% of the respective assembly.

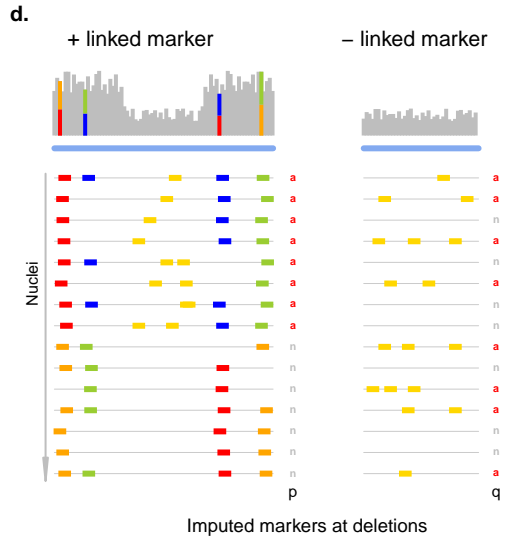
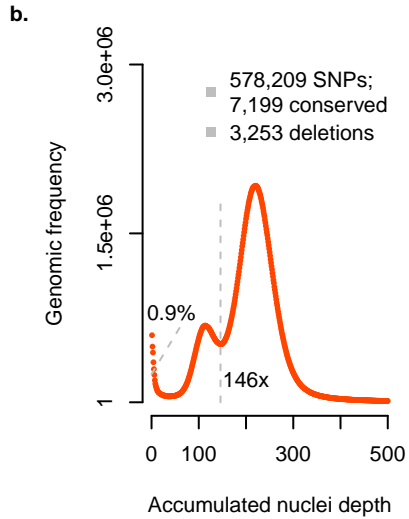
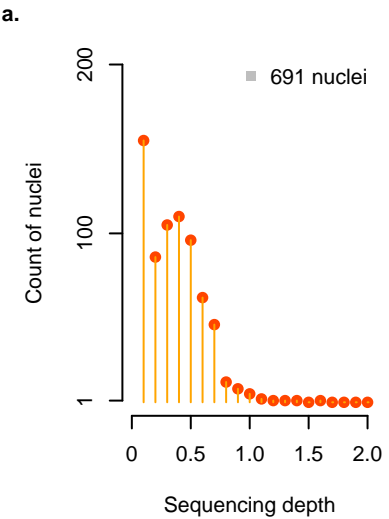
Haploid assemblies of 'Rojo Pasion'	Genome-specific <i>k</i> -mers common with parental WGS		Precision in haplotyping	Size [Mb]	Chromosome scaffolds	Contig N50 [Mb]	N50 [Mb]	Protein-coding genes (Total genes)
	'Currot'	'Orange Red'						
'Currot'-haplotype	12,983,934	129,874	99.1%	216.0	8	0.662	25.8	30,661 (52,472)
'Orange Red'-haplotype	81,422	16,807,958	99.5%	215.2	8	0.664	25.5	30,378 (51,701)

805
 806
 807
 808
 809
 810
 811
 812



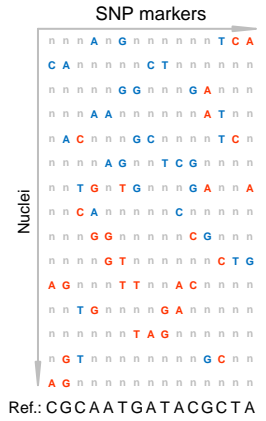
Diploid individual





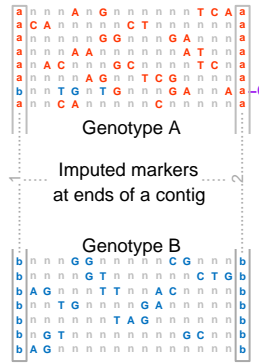
c.

SNP-marker genotyping

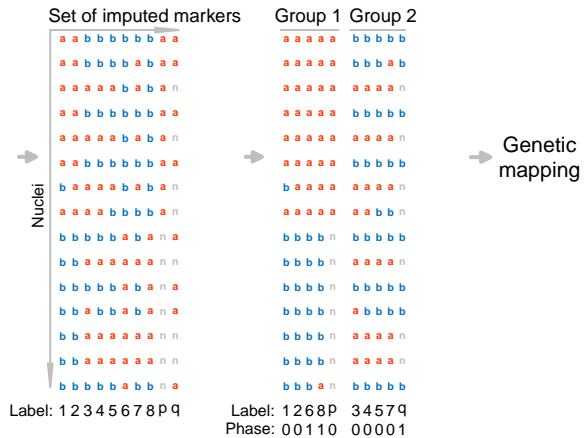


e.

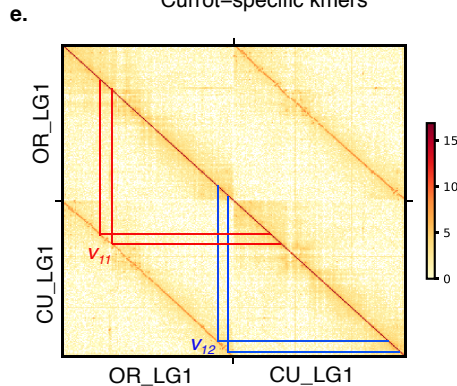
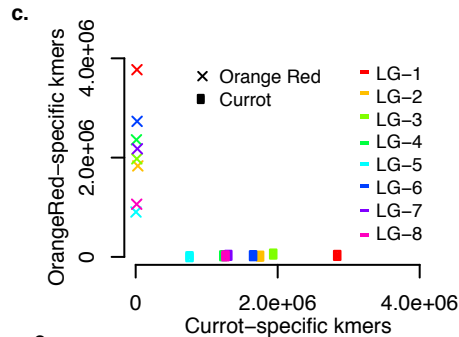
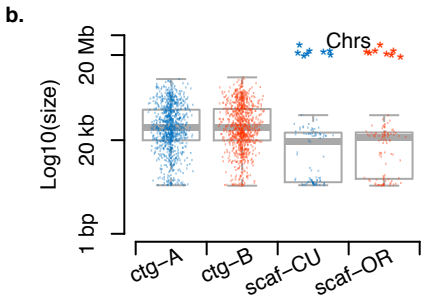
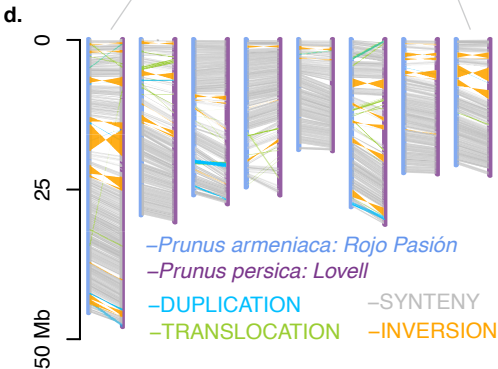
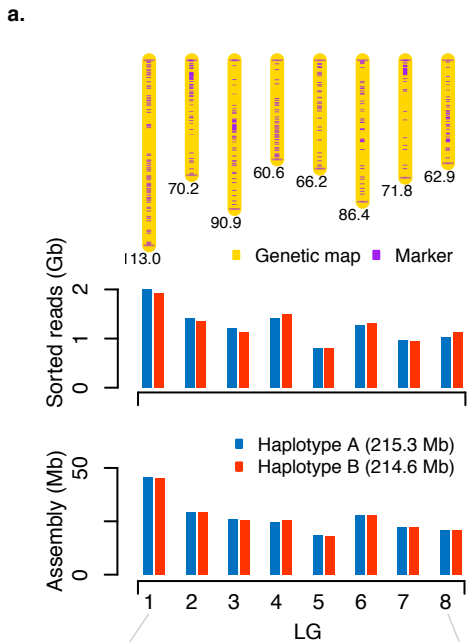
SNP phasing

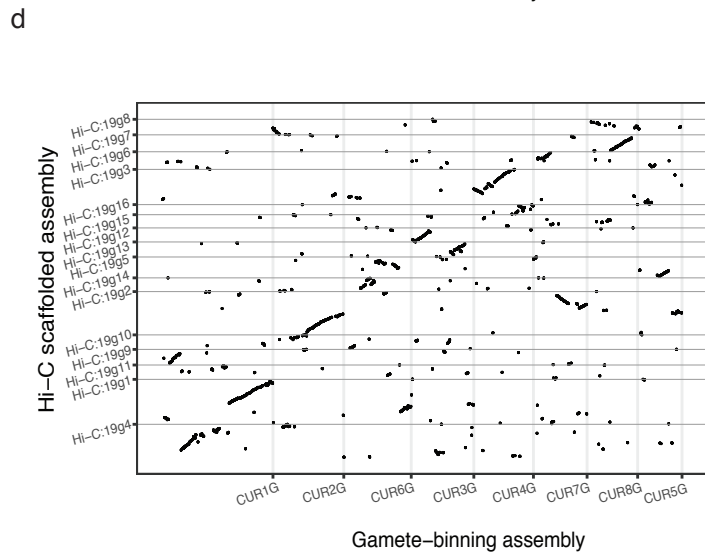
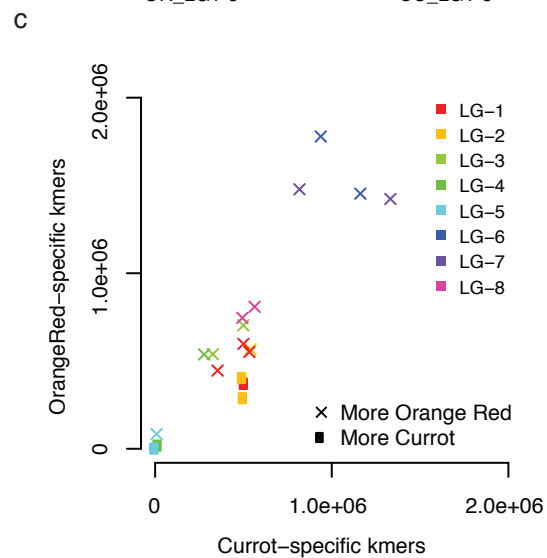
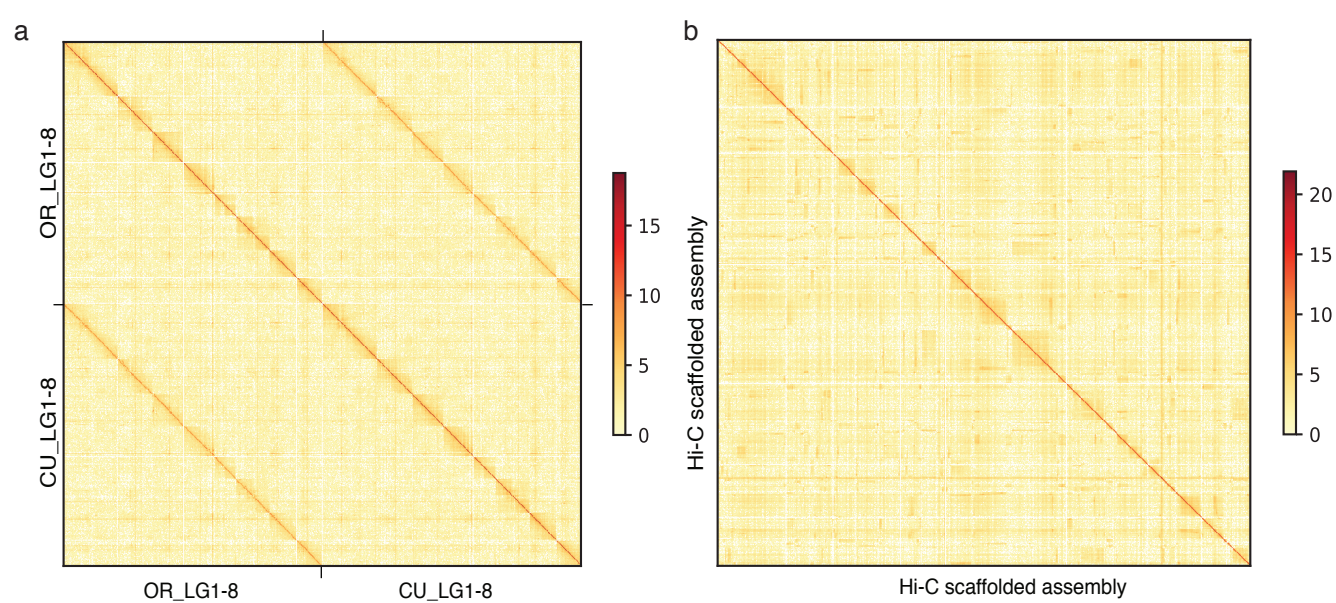


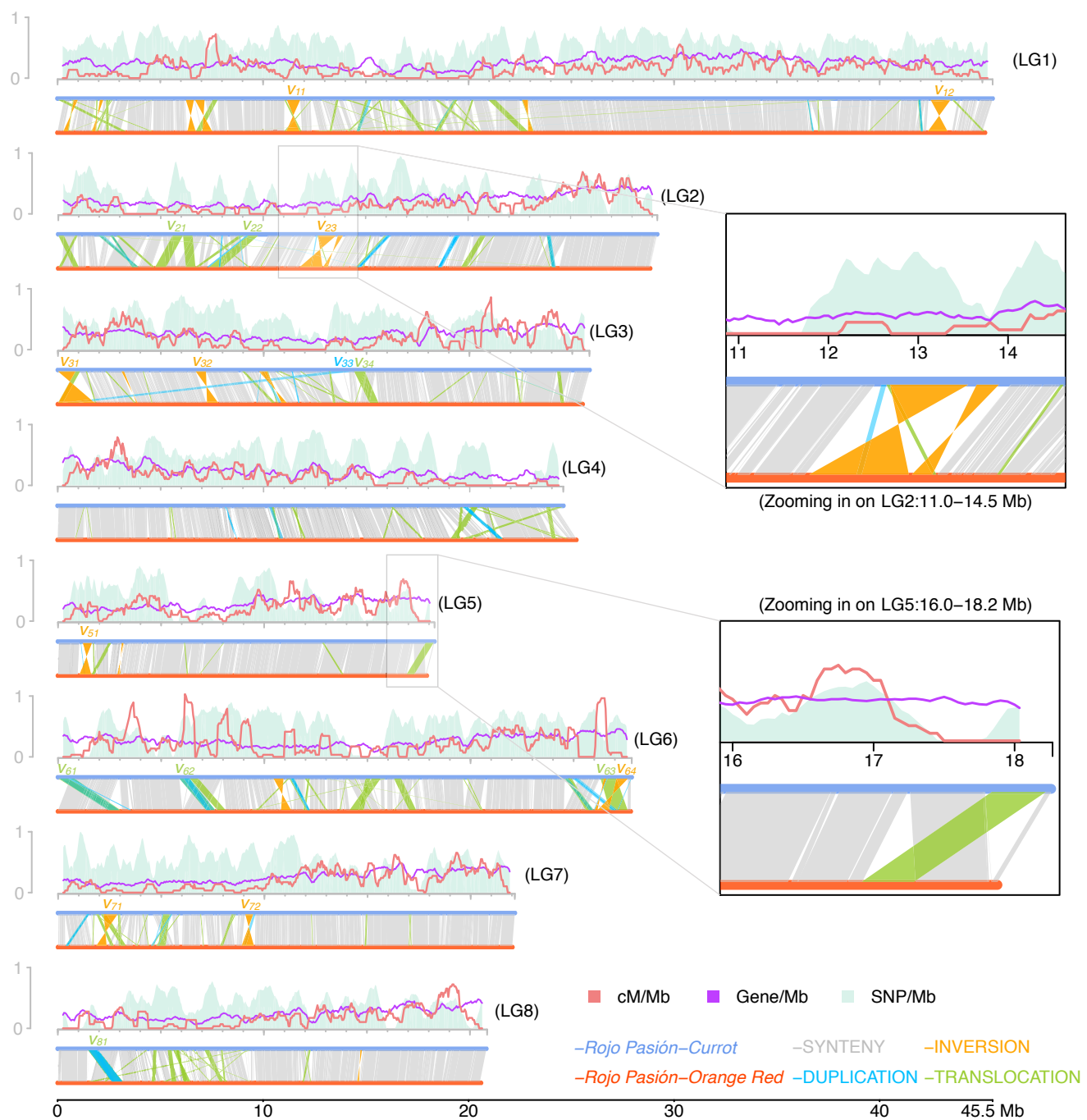
Contig phasing/grouping

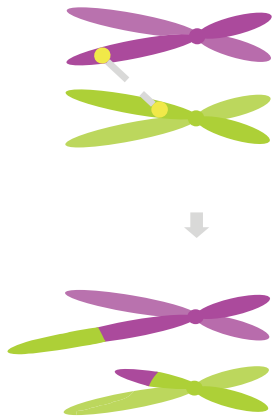


Genetic mapping







a.**b.**