

# 1 **The *Aquilegia* genome reveals a hybrid origin of core eudicots**

2

3 Gökçe Aköz<sup>1,2</sup> and Magnus Nordborg<sup>1,\*</sup>

4 <sup>1</sup>Gregor Mendel Institute, Austrian Academy of Sciences, Vienna Biocenter, Vienna,  
5 Austria

6 <sup>2</sup>Vienna Graduate School of Population Genetics, Vienna, Austria

7 \*magnus.nordborg@gmi.oeaw.ac.at

8

## 9 **Abstract**

10 **Background:** Whole-genome duplications (WGD) have dominated the evolutionary  
11 history of plants. One consequence of WGD is a dramatic restructuring of the genome as  
12 it undergoes diploidization, a process under which deletions and rearrangements of  
13 various sizes scramble the genetic material, leading to a repacking of the genome and  
14 eventual return to diploidy. Here, we investigate the history of WGD in the columbine  
15 genus *Aquilegia*, a basal eudicot, and use it to illuminate the origins of the core eudicots.

16 **Results:** Within-genome synteny confirms that columbines are ancient tetraploids, and  
17 comparison with the grape genome reveals that this tetraploidy is shared with the core  
18 eudicots. Thus, the ancient *gamma* hexaploidy found in all core eudicots must have  
19 involved a two-step process: first tetraploidy in the ancestry of all eudicots, then  
20 hexaploidy in the ancestry of core eudicots. Furthermore, the precise pattern of synteny  
21 sharing suggests that the latter involved allopolyploidization, and that core eudicots  
22 thus have a hybrid origin.

23 **Conclusions:** Novel analyses of synteny sharing together with the well-preserved  
24 structure of the columbine genome reveal that the *gamma* hexaploidy at the root of core  
25 eudicots is likely a result of hybridization between a tetraploid and a diploid species.

26

## 27 **Background**

28 Whole-genome duplication (WGD) is common in the evolutionary history of plants  
29 [reviewed in 1,2]. All flowering plants are descended from a polyploid ancestor, which in  
30 turn shows evidence of an even older WGD shared by all seed plants [3]. These repeated  
31 cycles of polyploidy dramatically restructure plant genomes. Presumably driven by the  
32 “diploidization” process, whereby genomes are returned to an effectively diploid state,  
33 chromosomes are scrambled via fusions and fissions, lose both repetitive and genic  
34 sequences, or are lost entirely [4–11]. Intriguingly, gene loss after WGD is non-random:  
35 not only is there a bias against the retention of certain genes [12,13], but also against the  
36 retention of one of the WGD-derived paralog chromosomes [6,9,14–16].

37

38 We investigated the history of WGDs in the columbine genus *Aquilegia* for two reasons.  
39 The first is related to its phylogenetic position: as a basal eudicot, columbines are  
40 members of the very earliest diverging branch of the eudicots [17,18]. This matters

41 because our understanding of eudicot karyotype evolution is limited to the heavily  
42 sampled core eudicots. Using the recently published *Aquilegia coerulea* genome [19],  
43 we are able to address key questions about the history of polyploidization in all eudicots.  
44 Second, we traced the origins of the columbine chromosomes with a particular focus on  
45 chromosome 4, which, compared to the rest of the genome: harbors more genetic  
46 polymorphism, has a higher transposable element density, has a lower gene density and  
47 reduced gene expression, shows less population structure worldwide, appears more  
48 permeable to gene flow, and carries the rDNA clusters [19].

49

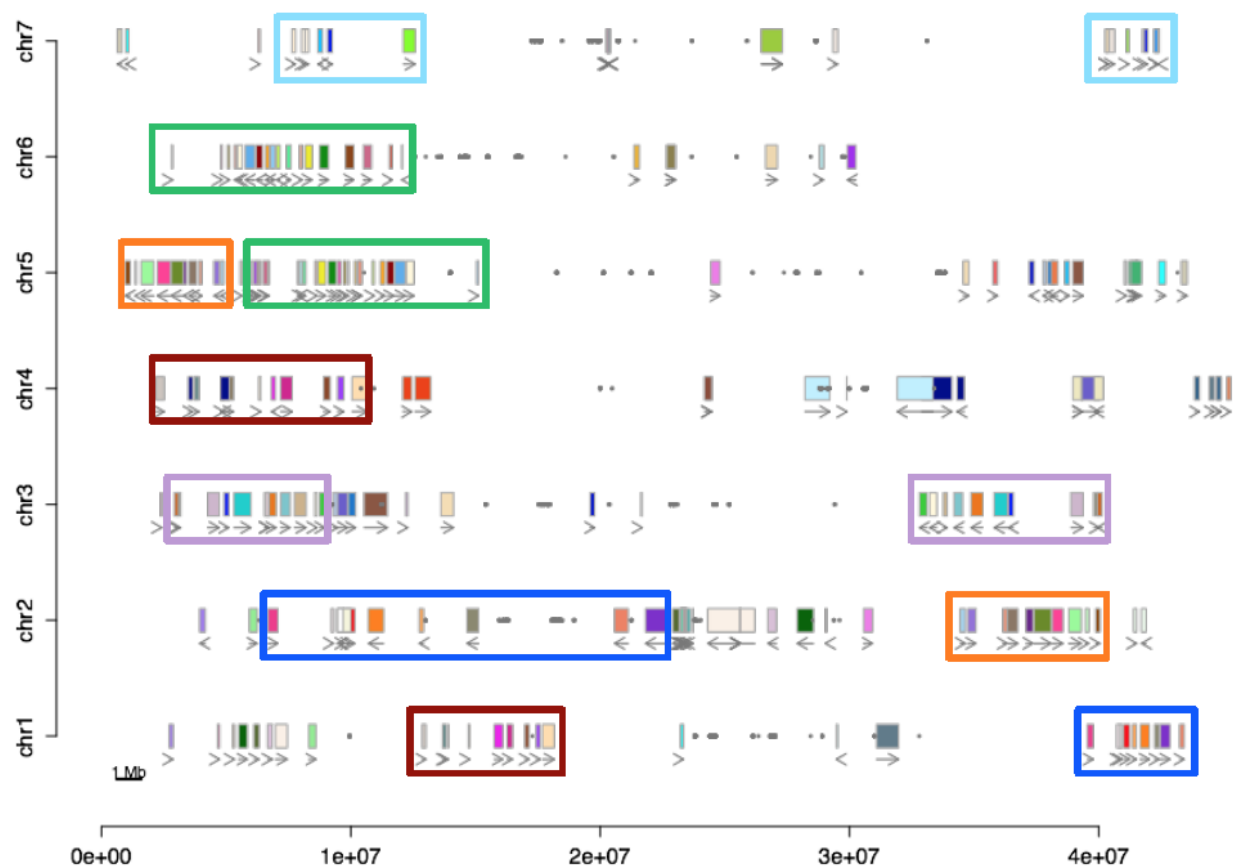
## 50 **Results**

### 51 **Within genome synteny confirms columbine paleotetraploidy**

52 Ancient WGDs have been commonly inferred from the distribution of divergences  
53 between gene duplicates. The simultaneous generation of gene duplicates via WGD is  
54 expected to produce a peak in the age distribution relative to the background age  
55 distribution of single gene duplicates [20–22]. Such a spike of ancient gene birth was  
56 the first evidence of paleotetraploidy in columbines [23], and was later supported by  
57 gene count-based modelling [24].

58

59 Given an assembled genome, a more direct method to infer ancient polyploidy is to look  
60 for regions with conserved gene order [25,26]. Such conservation (a.k.a., synteny)  
61 decreases over time due to gene loss and rearrangements, but will still be detectable  
62 unless the extent of change is extreme. We detected a total of 121 synteny block pairs  
63 harboring at least five paralogous gene pairs within the columbine genome. The  
64 distribution of these blocks across the seven columbine chromosomes indicates pairwise  
65 synteny between large genomic regions (Fig. 1). This 1:1 relationship suggests a single  
66 round of WGD in columbine, and is further supported by similar levels of divergence  
67 between synteny pairs (Figs. S1 and S2).



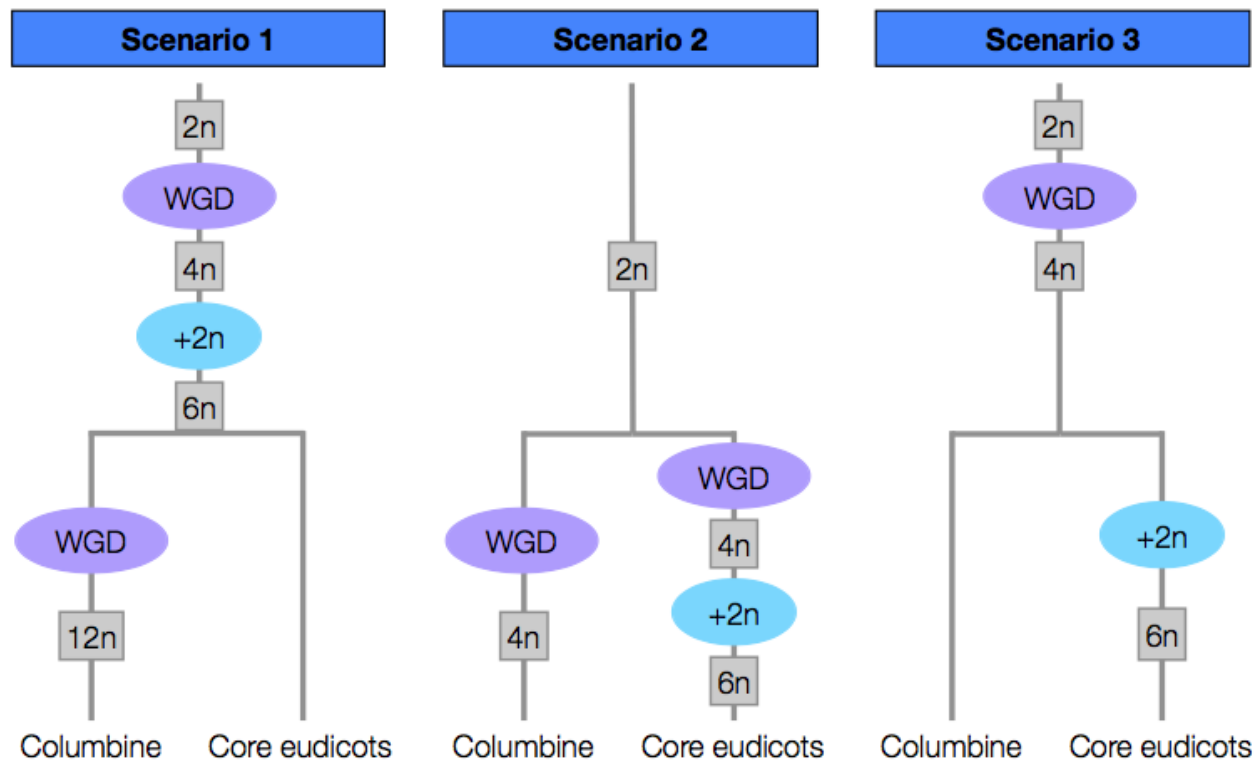
68  
69

70 **Fig. 1: Intragenomic synteny blocks in the columbine genome.** Pairs of synteny  
71 blocks are denoted as uniquely colored small rectangles. Larger rectangles of the same  
72 color outline large regions of synteny. Arrows under the synteny blocks show the  
73 orientation of the alignment between collinear genes. Grey dots highlight BLAST hits of  
74 a 329 bp centromeric repeat monomer [19,27].

## 75 **Columbines share ancient tetraploidy with core eudicots**

76 All sequenced core eudicots appear to share a triplicate genome structure due to  
77 paleohexaploidy postdating the separation of monocots and eudicots [9,28–32, and  
78 Supplementary Note 5 in 33]. The tetraploidy in columbines, a basal eudicot, might be  
79 independent of this ancient “*gamma*” hexaploidy (Scenarios 1 and 2 in Fig. 2) or might  
80 be a remnant of a WGD at the base of all eudicots, which formed the first step of the  
81 *gamma* hexaploidy in core eudicots (Scenario 3 in Fig.2).

82



83 Columbine Core eudicots Columbine Core eudicots Columbine Core eudicots

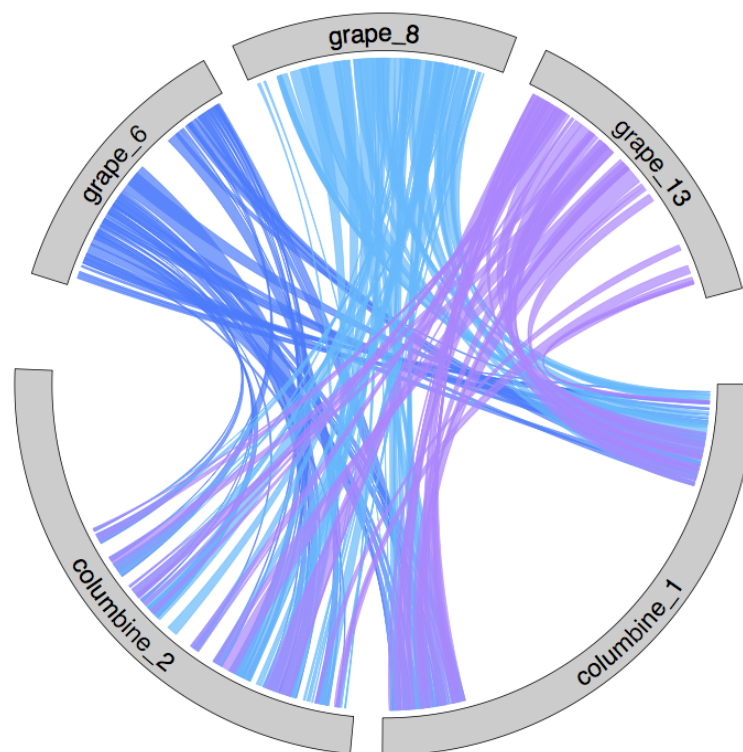
84

85 **Fig. 2: Three scenarios for the relationship between columbine tetraploidy**  
 86 **and core eudicot “gamma” hexaploidy.** The *gamma* hexaploidy is a two-step  
 87 process: a single round of WGD creates tetraploids ( $4n$ ) whose unreduced gametes then  
 88 fuse with diploid gametes ( $+2n$ ). **Scenario 1:** *gamma* hexaploidy precedes the split  
 89 between columbine and core eudicots, with the former undergoing an additional  
 90 tetraploidy. **Scenario 2:** Both *gamma* hexaploidy and columbine tetraploidy occur  
 91 after the split between columbines and core eudicots. **Scenario 3:** Columbine  
 92 tetraploidy is derived from the ancient tetraploidy that was the first step of the process  
 93 leading to *gamma* hexaploidy.

94

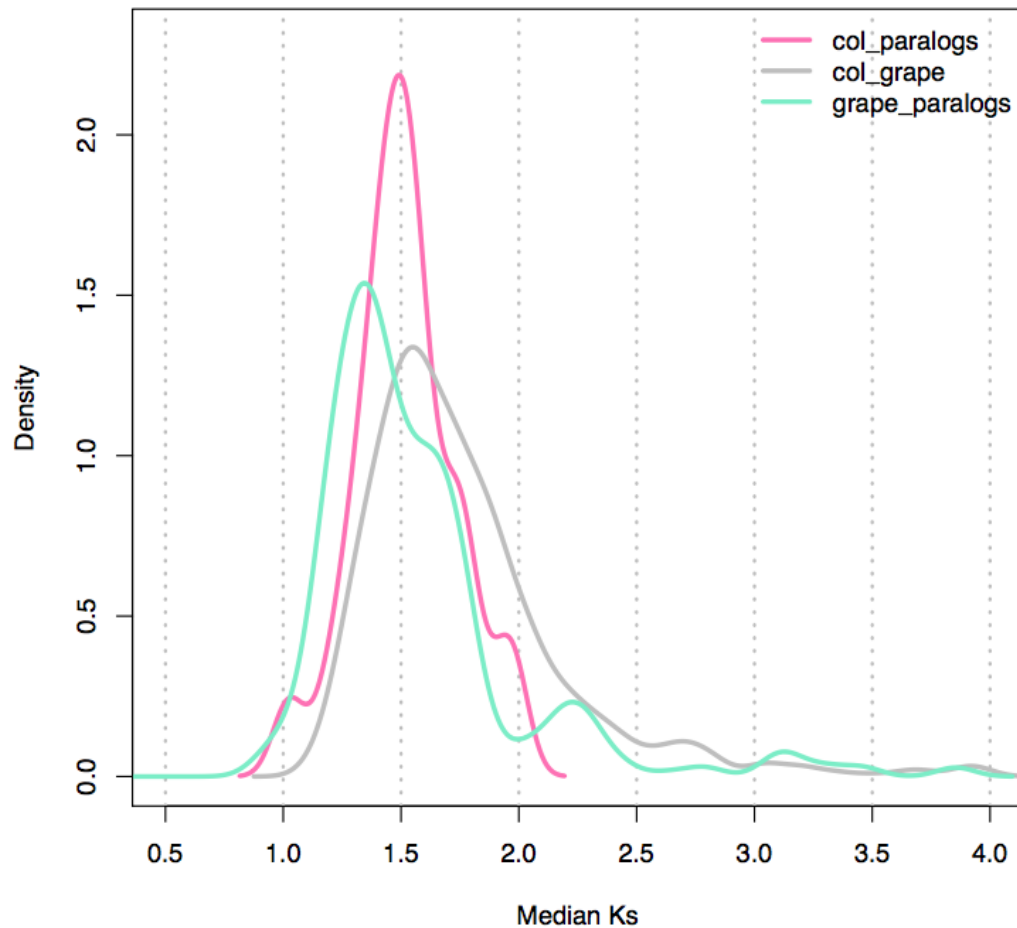
95 We used the grape (*Vitis vinifera*) genome as a representative of the core eudicots to  
 96 distinguish between the three scenarios in Fig. 2. Grape has experienced the least  
 97 number of chromosomal rearrangements post-*gamma* and thus strongly resembles the  
 98 ancestral pre-hexaploid genome [34]. Given the ploidy level of columbine under each  
 99 scenario, we can predict the synteny relationship between the homologous  
 100 chromosomes of grape and columbine, which is simply the ratio of haploid chromosome  
 101 set in grape to that of in columbine. If tetraploidy in columbines is lineage-specific and  
 102 superimposed on the *gamma* hexaploidy (Scenario 1), we would expect to find a 3:6  
 103 ratio of grape and columbine synteny blocks. Instead, we observe a 3:2 relationship  
 104 (Figs. 3 and S3) as expected under Scenarios 2 or 3. A similar 3:2 pattern is found in  
 105 comparisons between grape and sacred lotus [35]. This strongly suggests that basal

106 eudicots do not share the triplicate genome structure of core eudicots, ruling out  
107 Scenario 1.  
108



109  
110  
111 **Fig. 3: Synteny between the homologous regions of columbine and grape.**  
112 The results are shown here only for the columbine chromosomes 1,2 and the grape  
113 chromosomes 6,8 and 13 but reflect the overall synteny relationship of 3:2 between  
114 grape:columbine chromosomes (see Fig. S3 for the genome-wide synteny). This pattern  
115 argues against Scenario 1, but is consistent with either Scenario 2 or Scenario 3.

116  
117 To distinguish between the two remaining scenarios, we compared the divergence at  
118 synonymous sites (Ks) between columbine paralogs, grape paralogs and columbine-  
119 grape homologs. In agreement with the analysis of Jiao et al. [36], the Ks distribution  
120 for grape paralogs shows two major peaks, as expected under the two-step model for  
121 gamma hexaploidy (Fig. 4). However, columbine paralogs and columbine-grape  
122 homologs each show a single peak of divergence — and the peaks overlap each other and  
123 the “older” divergence peak of grape paralogs. This suggests that columbine tetraploidy  
124 is derived from the tetraploidy that eventually led to *gamma* hexaploidy in core eudicots  
125 (Scenario 3).  
126



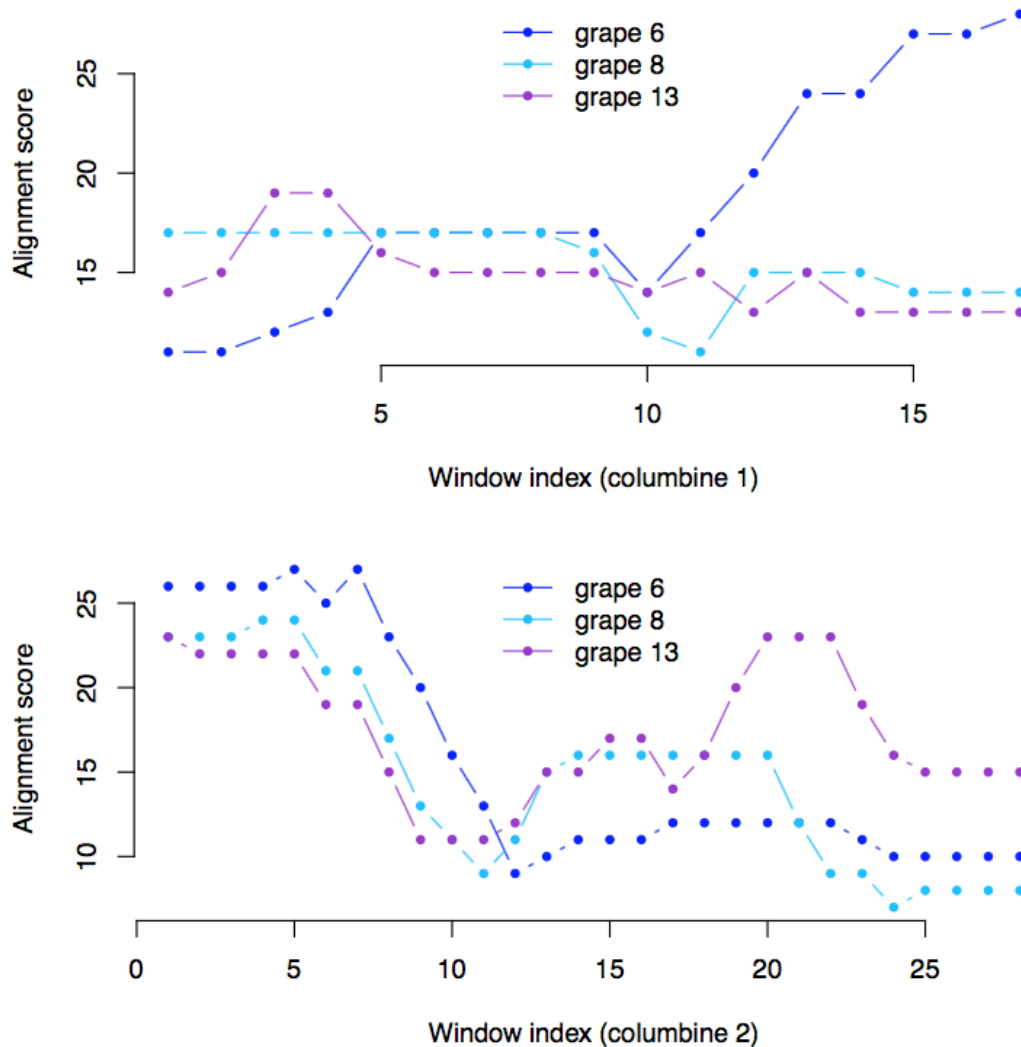
127  
128

129 **Fig. 4: The distribution of the median Ks across syntenic regions.** Synteny  
130 blocks are identified within columbine (col\_paralogs), between columbine and grape  
131 (col\_grape) and within grape (grape\_paralogs). Note that only the putative WGD-  
132 derived blocks (median Ks=1-2) are kept in columbine (Fig. S2).

133

134 To further explore the hypothesis of a shared WGD by all eudicots, we focused on the  
135 gene order similarity between the homologous regions of columbine and grape. If  
136 columbine and grape have descended from a common tetraploid ancestor, they should  
137 have inherited diploidization-driven differential gene order on the paralogous  
138 chromosomes of the ancestor (Fig. S5). As a result, we expect to see the alternative  
139 paralogous gene orders to be uniquely shared between two different pairs of columbine  
140 and grape chromosomes. To detect this, we first searched for at least three consecutive  
141 genes aligning between a pair of columbine and grape chromosomes and then looked at  
142 the distribution of these genes on all the columbine and grape chromosomes. This way  
143 of reconstructing chromosomes clearly shows that each of the paralogous chromosome  
144 pairs in columbine has a match to a single grape chromosome, with respect to its gene  
145 order (Fig. S6). This result was corroborated by a second approach where we quantify

146 the similarity between columbine and grape chromosomes. When we performed a  
147 pairwise alignment between each sliding window of genes on a columbine chromosome  
148 and all the genes on a grape homolog, we again see that each member of columbine  
149 paralogs gets the best hit to a single grape chromosome (Figs. 5 and S7-9). Reshuffling  
150 genes on grape chromosomes further indicates that this pattern of clustering is highly  
151 unlikely to be produced by chance alone ( $p=0-0.05$ ).  
152



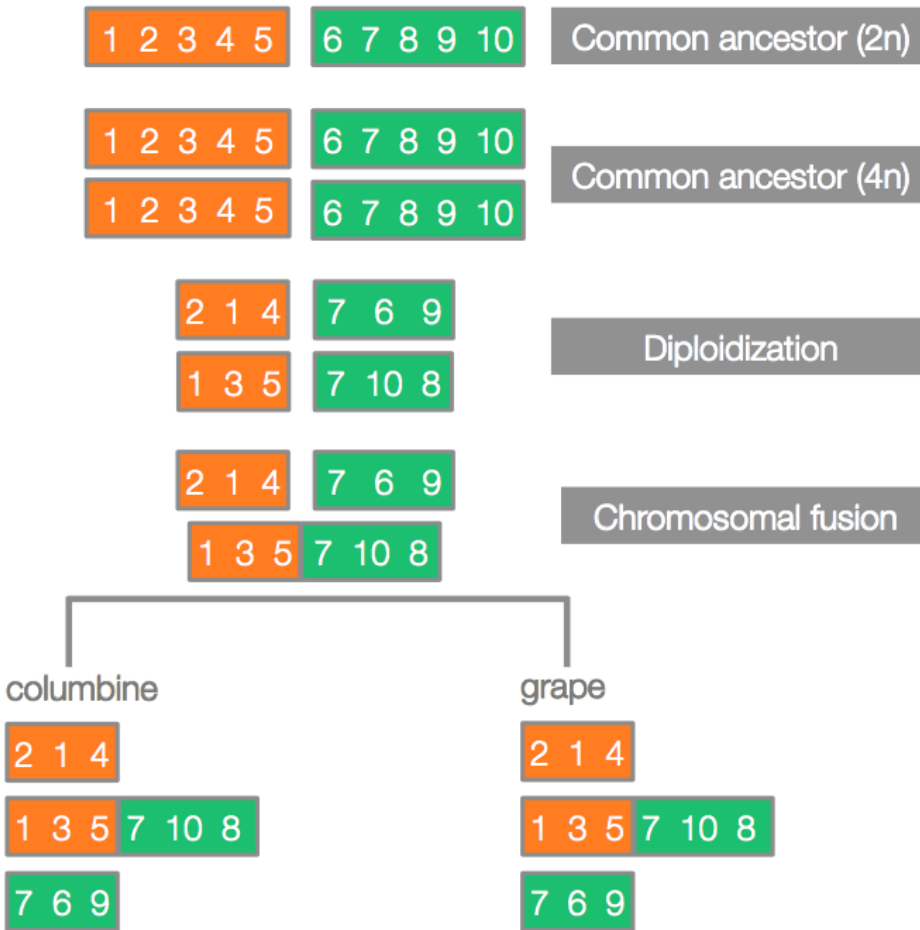
153  
154

155 **Fig. 5: Examples of gene order similarity between the homologous regions**  
156 **of columbine and grape.** For successive windows of genes within a given columbine  
157 chromosomal region, the best alignment score with respect to each of the three grape  
158 chromosomes harboring homologous regions, is given. For example, columbine  
159 chromosomes 1 and 2 share a paralogous region homologous to grape chromosomes 6,  
160 8, and 13 (Figs. 7 and S3). The chromosome 1 region (**top panel**) appears to be most  
161 closely related to grape chromosome 6, whereas its paralogous counterpart on  
162 chromosome 2 (**bottom panel**) appears to be most closely related to grape

163 chromosome 13. See Fig. S6 for the correspondence between gene orders and scores,  
164 which peak towards the end of each columbine region. Note that the results presented  
165 here are shown for a window size of 12 genes but remain significant for all the window  
166 sizes tested ( $p=0-0.05$ ).

167  
168 An eudicot-wide WGD is further supported by the observation that a chromosomal  
169 fusion, presumably experienced by the common tetraploid ancestor, is still detectable in  
170 the genomes of columbine and grape despite their separation of around 125 million  
171 years [37]. The first hint comes from the composition of the chromosomes: columbine  
172 chromosome 5 and grape chromosome 7 share the two chromosomal origins (Fig. 7). If  
173 these fused chromosomes were created by a single fusion event in the common  
174 tetraploid ancestor of eudicots, they should match each other with respect to gene order  
175 on each of the two homologous portions (“orange” and “green” portions in Fig. 6). This  
176 is what we see: columbine chromosome 5 and grape chromosome 7 cluster together with  
177 respect to their gene order on the “orange” portion (Fig. S7). For the “green” portion,  
178 columbine chromosome 5 matches grape chromosome 4 (Fig. S8), which used to be  
179 fused to grape chromosome 7 [38]. Additional support for shared ancestral fusion comes  
180 from the cacao (*Theobroma cacao*) genome [39]. The first chromosome of cacao does  
181 not only show a similar pattern of chromosomal ancestry [38,39], but also shares the  
182 gene order exclusively with the grape chromosomes 4 and 7 on the corresponding  
183 homologous portions (Fig. S10). In summary, the columbine fusion clusters with that of  
184 grape, which, in turn, clusters with that of cacao, strongly favoring a common origin of  
185 the fusion between “orange” and “green” ancestral chromosomes (Fig. 7).

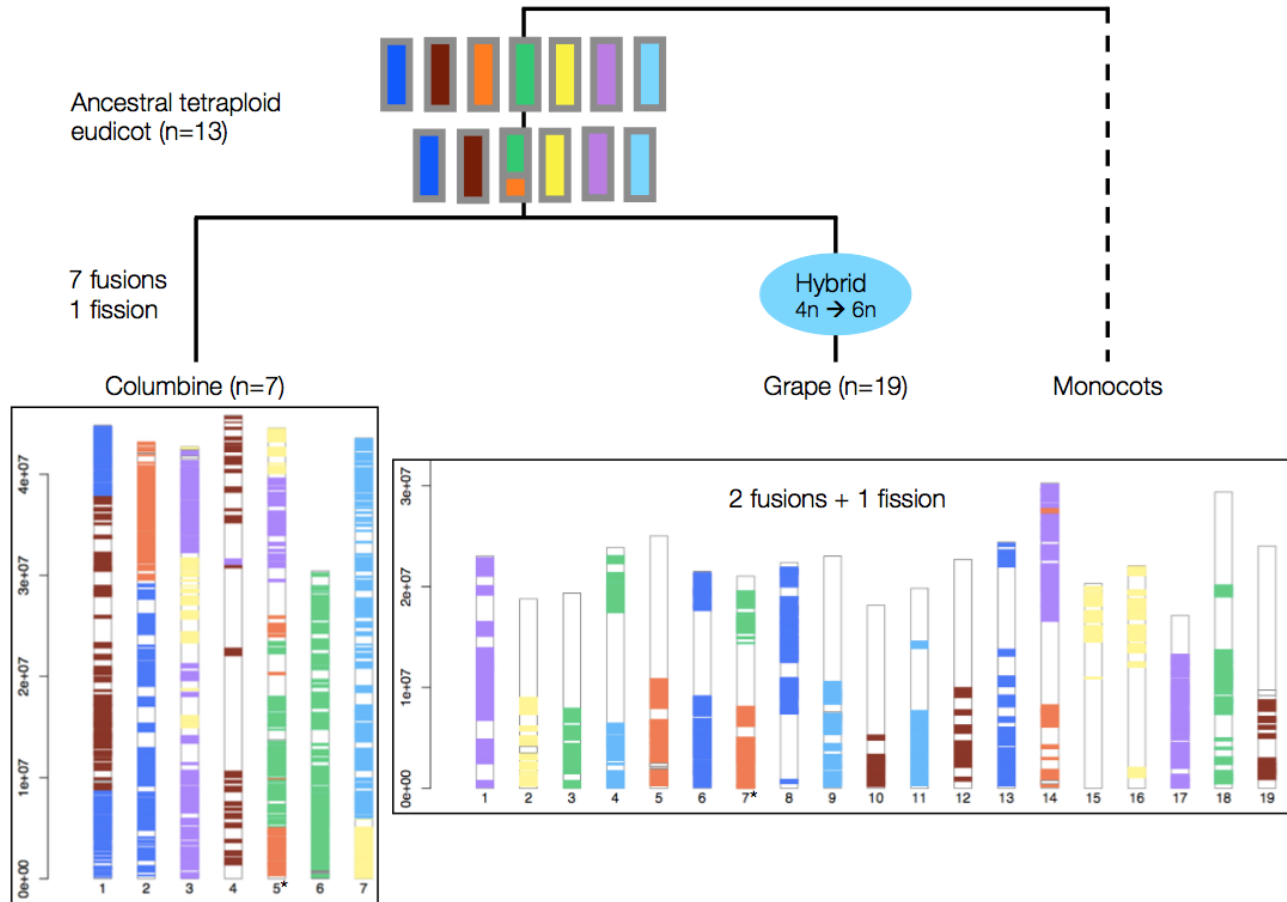




186  
187

188 **Fig. 6: Schematic of predicted synteny patterns in the case of shared**  
189 **ancestral fusion.** Two ancestral chromosomes (orange and green rectangles, with  
190 genes depicted as numbers) undergo WGD. Paralogous chromosome pairs diverge as a  
191 part of the diploidization process. A fusion joins one version of the “orange”  
192 chromosome (‘1, 3, 5’) with one version of the “green” chromosome (‘7, 10, 8’). If this  
193 took place in the common tetraploid ancestor of eudicots, the fused chromosomes in  
194 columbine and grape should also carry these versions on their “orange” and “green”  
195 portions. In the hypothetical example here, diploidization precedes the fusion event but  
196 may well happen afterwards with no effect on the predicted synteny patterns.

197



198  
199

200 **Fig. 7: Tracing the genome reshuffling in columbine following tetraploidy.**

201 Grape chromosomes (**bottom right**) are colored by within-genome synteny. Seven  
202 distinct colors represent the haploid set of seven ancestral chromosomes before the  
203 eudicot-wide WGD. Each color is shared by three grape chromosomes reflecting the  
204 triplicate genome structure of core eudicots. The only exception is the “green”  
205 chromosome which is shared by four grape chromosomes due to a fission event [38].  
206 Columbine chromosomes (**bottom left**) are colored by their synteny to grape  
207 chromosomes. Each color is generally shared by two chromosomes, reflecting columbine  
208 paleotetraploidy. As few as 7 fusions and a single fission are enough to explain the  
209 current structure of the columbine genome. Of these 7 fusions, 5 are between different  
210 chromosomes while 2 are between WGD-derived paralogous chromosomes. Columbine  
211 chromosomes 3 and 7 are examples of the latter (Figs. 1 and S4). Note that chromosome  
212 5 of columbine and chromosome 7 of grape (\*) both have the colors “orange” and  
213 “green” (cf. Fig. 6).

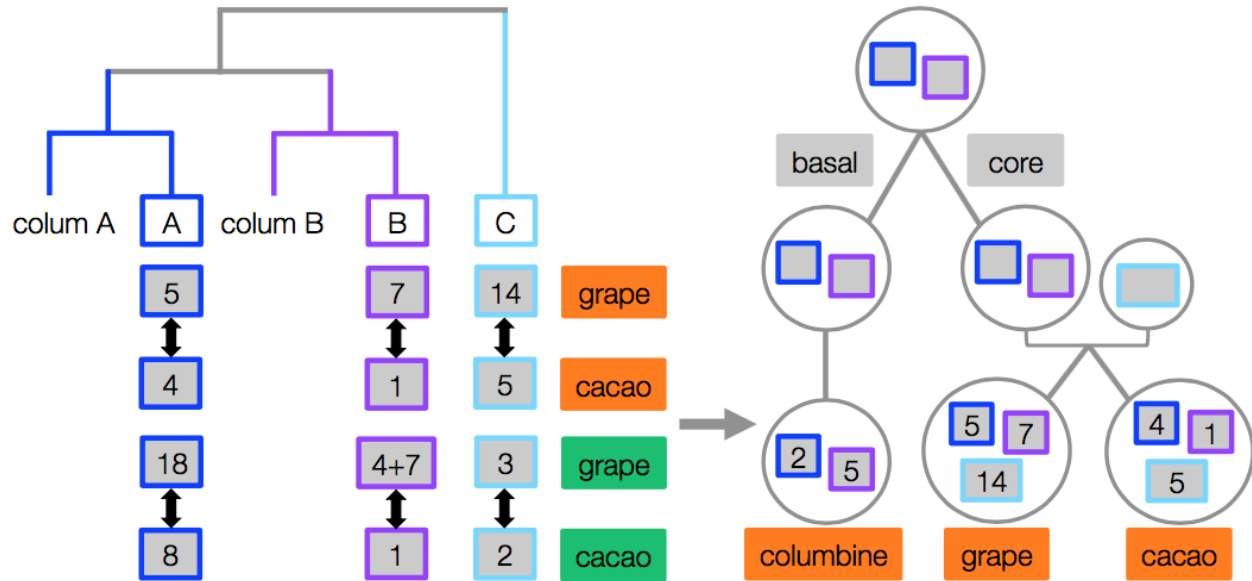
214

## 215 **The core eudicots have a hybrid origin**

216 Our inference of shared tetraploidy between basal and core eudicots makes use of the  
217 signals presumably generated by diplodization (Figs. 6 and S5). However, hybridization  
218 of unreduced gametes from two divergent diploid genomes, “allotetraploidy”, would also  
219 lead to gene order-based clustering between two different pairs of grape and columbine  
220 chromosomes (Figs. S11-12). In this case, the alternative paralogous gene orders of the  
221 tetraploid ancestor reflect the gene orders on progenitor chromosomes. Thus, the  
222 clustering pattern does not depend on whether the eudicot tetraploid genome evolved  
223 via “auto-” or “allopolyploidy”. The same is not true for the second of the process leading  
224 to hexaploidy. Only allohexaploidy would lead to one of the three paralogous grape  
225 chromosomes being an “outlier” to the two grape-columbine pairing (Figs. S11-12) –  
226 which is what we see in our data (light blue lines in Figs. 5 and S7-9).

227  
228 If our interpretation is correct and all core eudicots have a hybrid origin, the pattern of  
229 gene order-based clustering should be conserved. That is, we should be able to identify  
230 the hexaploidy-derived “outlier” chromosomes in other core eudicot genomes as well. To  
231 check this expectation, we again used the cacao genome, one of the most conserved  
232 genomes after grape [9,39]. Pairwise alignment between the homologous regions of  
233 columbine and cacao confirms our expectation: each member of columbine paralogs  
234 pairs up with a single cacao chromosome, leaving one of the cacao paralogs as an outlier  
235 (Figs. S13-14). Furthermore, as shown in Fig. 8 (see also Fig. S10), the cacao regions  
236 putatively derived from tetraploidy and hexaploidy, respectively, show a very clear one-  
237 to-one match to those in grape (detected in the grape-columbine comparison). As  
238 expected, the putatively orthologous pairs of cacao and grape regions show similar levels  
239 of synteny conservation with their paralogous counterparts, with the “outlier” regions  
240 being the most divergent [38]. Thus, the cacao genome provides an independent line of  
241 evidence for a hybrid origin, and highlights the key role of the columbine genome in  
242 unravelling the history of the eudicot genome.

243



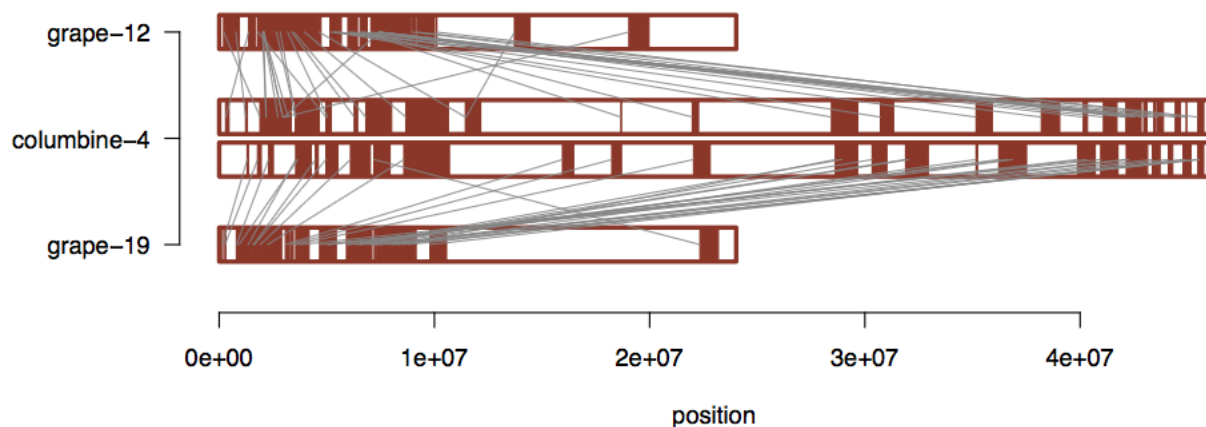
244  
245

246 **Fig. 8: The shared history of chromosomes in columbine, grape and cacao.**  
 247 Gene order-based clustering results (**left panel**) are summarized here for the  
 248 chromosomes harboring the “orange” and “green” homologous portions. The former  
 249 corresponds to 5, 7, 14 in grape and 1, 4, 5 in cacao. The latter corresponds to 3, 4+7  
 250 (products of a fission), 18 in grape and 1, 2, 8 in cacao. In columbine, the “orange”  
 251 portions are on chromosomes 2 and 5 while the “green” portions are on chromosomes 6  
 252 and 5, each pair of which being denoted as *colum A* and *colum B*, respectively. Both  
 253 grape- and cacao-columbine pairing distinguish tetraploidy-derived regions (blue and  
 254 purple rectangles) from hybridization-derived ones (light blue rectangles), defining the  
 255 orthologous sets of regions across the three eudicot genomes (**right panel**). The  
 256 conservation of gene order exclusively between the putatively orthologous regions of  
 257 grape and cacao (black arrows, Fig. S10) further strengthens our columbine-based  
 258 inference of orthology.

### 259 **Current columbine chromosomes have mostly been generated via fusions**

260 It is widely accepted that genome shuffling post-WGD has shaped the present-day  
 261 karyotypes of all plant genomes [34]. Nevertheless, the extent of genome shuffling as a  
 262 part of the “re-diploidization” process seems to vary widely: only 3 chromosomal  
 263 rearrangements post-*gamma* are enough to explain the current structure of the grape  
 264 genome (Fig. 7) while almost 150 chromosomal rearrangements were necessary for the  
 265 sunflower genome to reach its current karyotype after several rounds of WGD [11]. To  
 266 check where columbine falls in this spectrum, we identified chromosomal  
 267 rearrangements likely to have happened after the tetraploidy shared by all eudicots: if  
 268 the pre-WGD ancestral eudicot karyotype had a haploid number of 7 chromosomes [28],  
 269 only seven columbine-specific fusions and a single fission are enough to explain the

270 reduction in columbine chromosome number from  $n=13$  to  $n=7$  after the ancestral  
271 fusion event (Fig. 7). These rearrangements involve all the chromosomes in columbine  
272 apart from chromosomes 4 and 6, the former of which paradoxically shows the greatest  
273 erosion of synteny with grape chromosomes (Figs. 7 and S3). Given all the evidence  
274 suggesting a “decaying” nature of columbine chromosome 4 [19], we repeated the  
275 analysis of grape-columbine synteny detection with relaxed parameter settings. We did  
276 this by decreasing the minimum number of aligned gene pairs within a block (from 5 to  
277 3) and increasing the maximum genic distance between matches (from 20 to 30). This  
278 allowed us to extend the synteny blocks towards more proximal regions (Fig. S15).  
279 Further zooming into the synteny relationship between grape chromosomes that are  
280 homologous to columbine chromosome 4 confirmed that there is no evidence of a fusion  
281 event (Fig. 9).  
282



283  
284 **Fig. 9: Synteny between columbine chromosome 4 and grape chromosomes**  
285 **12 and 19.** Much smaller grape chromosomes look like the compact versions of  
286 columbine chromosome 4. Note that this result is generated with the most relaxed  
287 parameter combination in Fig. S15, but holds true for a less relaxed combination of  
288 parameters as well (Fig. S16).  
289

290 The lack of a fusion event on columbine chromosome 6 might explain the fact that it is  
291 the smallest chromosome of columbine (Fig. 7). However, chromosome 4 is comparable  
292 in size to the remaining chromosomes, all of which are products of ancient fusion  
293 events. The observations that chromosome 4 has a higher proportion of genes in tandem  
294 duplicates (0.37 versus genome-wide mean of 0.22) and a greater extent of intra-  
295 chromosomal synteny (indicative of segmental duplications) (Fig. S17) suggest that  
296 chromosome 4 has reached a comparable size partly due to numerous tandem and  
297 segmental duplications and partly due to an expansion of repetitive DNA [19]. These  
298 results reinforce the idea that chromosome 4 has followed a distinct evolutionary path  
299 from the rest of genome.

300 Fusion-dominated genome shuffling [34] is not the only facet of diploidization [40].  
301 Following WGDs, gene duplicates get lost and this happens in a non-random manner.  
302 Genes involved in connected molecular functions like kinases, transcription factors and  
303 ribosomal proteins are retained in pairs [41–45] potentially due to dosage-related  
304 constraints [46]: losing or duplicating some, but not all of these dosage-sensitive genes  
305 might upset the stoichiometric relationship between their protein products [47–49].  
306 Consistent with this dosage balance hypothesis, columbine genes potentially retained  
307 post-WGD (1302 genes across 76 syntenic regions; Supplementary Data 1) are enriched  
308 for the GO categories “structural constituent of ribosome”, “transcription factor  
309 activity”, “translation” ( $p < 0.001$ ) and “protein tyrosine kinase activity” ( $p < 0.01$ ).  
310 Tandemly duplicated genes ( $n = 6972$ ), on the other hand, are depleted for the GO  
311 categories “structural constituent of ribosome”, and “translation” ( $p = 10^{-17}$ ), reflecting  
312 the role of dosage-related purifying selection.

313

## 314 **Discussion**

315 The evolutionary history of plants is characterized by multiple WGDs and columbine is  
316 no exception. The alignment between chromosomal regions in a 1:1 ratio (Fig. 1)  
317 confirms a single round of WGD in columbine [23,24,50]. Furthermore, we demonstrate  
318 that this tetraploidy is shared with all eudicots, refuting a lineage-specific  
319 polyploidization in columbine [23,50], in favor of an eudicot-wide WGD [26,35,36,51].  
320 Unlike previous attempts based on genetic distance, our approach simply relies on gene  
321 order conservation. It also takes advantage of the well-preserved genome structure of  
322 columbines: free from recent WGDs, the columbine genome carries only the traces of  
323 the ancient tetraploidy.

324

325 This approach also helps us shed light on the nature of the *gamma* hexaploidy found in  
326 all core eudicots [9,28–32, and Supplementary Note 5 in 33]. WGDs have often been  
327 discussed as if they were “events”, ignoring the process by which they originated. We  
328 show here that core eudicot hexaploidy is the result of two processes: an ancient  
329 tetraploidization shared by all eudicots, followed by allopolyploidization leading to core  
330 eudicots. In other words, all core eudicots have a hybrid origin. An allohexaploid origin  
331 has indeed been previously suggested by Murat et al. [9], who identified the three  
332 subgenomes of grape using differential patterns of gene loss on “dominant” versus  
333 “sensitive” subgenomes. Their classification assumes that the most recently added set of  
334 paralogous chromosomes will be “dominant”, because they have spent a shorter amount  
335 of time in the polyploid genome and thus experienced fewer gene losses. Contrary to  
336 this, our results suggest that the most recently added grape chromosomes  
337 (chromosomes 3, 8, 9 and 14) largely corresponds to the “sensitive” grape chromosomes  
338 identified by Murat et al. [9]. Instead, we argue that the extensive gene loss in the  
339 “youngest” subgenome reflects its divergence from the other two subgenomes at the  
340 time of hexaploid formation, perhaps similar to the situation in the allotetraploid

341 *Arabidopsis suecica*, which is a hybrid between the more ancestral (n=8) genome of *A.*  
342 *arenosa*, and the heavily reduced (n=5) genome of *A. thaliana* [52]. Another example is  
343 hexaploid wheat, which is a hybrid between tetraploid emmer wheat and a wild diploid  
344 grass, *Aegilops tauschii* [53 and references therein].

345

## 346 **Conclusions**

347 Our findings help us understand the hybrid structure of core eudicot genomes and will  
348 hopefully encourage larger scale analyses to understand what hybridization has meant  
349 for core eudicots — a group which comprises more than 70% of all living flowering  
350 plants [54]. What are the hybridization-coupled changes that has led to the current  
351 patterns of gene expression, methylation, transposable element density/distribution?  
352 All these questions call for additional genomes from basal eudicots which — as this  
353 study illustrates — have great values as outgroup to the core eudicots.

354

## 355 **Materials and Methods**

### 356 **Synteny detection**

357 We performed all genes (CDS)-against-all genes (CDS) BLAST for the latest version of  
358 *Aquilegia coerulea* reference genome (v3.1) using SynMap tool [29] in the online CoGe  
359 portal [55]. We also looked at the synteny within *Vitis vinifera* (v12) and between *V.*  
360 *vinifera* and *A. coerulea* using both default and more relaxed parameter combinations  
361 in DAGChainer. We filtered the raw output files for both within grape and grape-to-  
362 columbine synteny. For the former, we only kept the blocks that are syntenic between  
363 the polyploidy-derived paralogous chromosomes of grape as identified by Jaillon et al.  
364 [28]. For the latter, we required that a given columbine chromosome is overall syntenic  
365 to all the three paralogous chromosomes of grape. So, for a given pair of columbine and  
366 grape chromosomes, we only kept the blocks if the columbine chromosome also matches  
367 to the other members of paralogous grape chromosomes.

368

369 The raw output files can be regenerated at the CoGe portal [55] using the id numbers  
370 provided below for each species (Availability of data and material) and changing the  
371 default parameter combination in DAGChainer (D:A=20:5) when needed. D and A  
372 specify the maximum genic distance between two matches and the minimum number of  
373 aligned gene pairs, respectively, to form a collinear syntenic block.

### 374 **Estimating the divergence between synteny block pairs**

375 We used Ks (the number of synonymous substitutions per synonymous site) values  
376 provided for each duplicate gene pair by the CoGe portal [55]. We estimated the median

377 Ks of all gene duplicates in a syntenic block after filtering duplicates with  $K_s > 10$  due to  
378 saturation effect [56].

### 379 **Quantifying gene order similarity**

380 We first detected three consecutive genes aligning between a pair of columbine and  
381 grape chromosomes harboring homologous regions (D:A=0:3). We particularly chose  
382 three genes since it is the most stringent value we could use to detect homologous  
383 synteny blocks; we detected almost nothing when we required 4 consecutive genes  
384 (D:A=0:4). We then looked at the distribution of these genes on a given pair of  
385 columbine and grape chromosomes and also on their paralogous counterparts  
386 (D:A=0:1). Once we had the gene order for each chromosome, we assigned a unique  
387 word to each synteny block and the genes forming the block to be able to use the text  
388 alignment provided by the R package *align\_local* [57]. Having each chromosome  
389 represented by a sentence, we quantified the gene (“word”) similarity as such: for an  
390 initial N number of words on a columbine chromosome (N=window size), we did a  
391 pairwise alignment between these N words and a grape chromosome (match=4, gap=-  
392 1). We repeated the same analysis with the inverted order of N words and picked the  
393 maximum alignment score. We repeated these steps by sliding the window by one word  
394 and keeping the N constant to get a distribution of scores as in Fig. 5. We used different  
395 N values ranging from 4 to 15. Note that we excluded columbine chromosomes 3 and 4  
396 from this analysis since both have a complex history of lineage-specific chromosomal  
397 reshuffling events: fusions and a fission gave rise to columbine chromosome 3 (Figs. 1  
398 and 7) while duplications have shaped the current structure of columbine chromosome  
399 4 (Fig. S17).

400  
401 We applied the same stringent criteria (D:A=0:3) to detect the homologous regions  
402 between grape and cacao (*Theobroma cacao*, v1). The same criteria led to very few  
403 homologous regions between columbine and cacao. So, we relaxed the parameters for  
404 the synteny detection between these two genomes (D:A=0:2) and quantified the gene  
405 order similarity with greater window sizes (N=20, 30, 35, 40 and 50). Note that we  
406 focused on the triplicated regions distributed across 3 different cacao chromosomes  
407 (Figs. 8, S13-14), which are rather unaffected by lineage-specific shuffling [38].

### 408 **Statistical testing of gene order similarity**

409 Given the gene order similarity between the two different pairs of columbine and grape  
410 chromosomes harboring homologous regions, we performed permutation tests to  
411 estimate the probability of observing such a clustering just by chance. To do so, we first  
412 combined all the grape genes and sampled the same number of genes (“words”) as we  
413 observe to reconstruct each of the paralogous grape chromosome. We repeated the  
414 quantification step as above to get a permuted distribution of alignment score between a



415 pair of columbine and grape chromosomes. We used Wilcoxon rank sum test (W-  
416 statistic) to quantify the shift in the distribution of alignment scores between one of the  
417 members of columbine paralogous chromosomes and its best grape hit when combined  
418 with the alignment scores between the same columbine chromosome and other grape  
419 chromosomes. We repeated the same analysis for the other member of columbine  
420 paralogous chromosomes as well. Having these *observed* W-statistics, we counted the  
421 number of cases (out of 100) where the permuted distributions generate W-statistics as  
422 high as or higher than the observed ones. Note that for columbine chromosome 7, whose  
423 structure has been greatly shaped by the fusion of WGD-derived paralog chromosomes  
424 (Fig. S4), we created two paralogous chromosomes using the observed distribution of  
425 alignment scores (Fig. S9). Columbine chromosome 7 matches best to grape  
426 chromosome 11 for the first 14 “words” and to grape chromosome 4 for most of the  
427 remaining “words”, which define the putative boundaries of columbine paralogous  
428 chromosomes before the fusion event. We ran permutation tests for the columbine-  
429 cacao pairing as well (Figs. S13-14).

### 430 **GO enrichment analysis**

431 We used gene annotations provided by JGI [19] to test the null hypothesis that the  
432 property for a gene to be retained post-WGD and to belong to a given GO category are  
433 independent. We created a 2x2 contingency table as shown below and applied Fisher’s  
434 exact test for each GO category independently. We repeated the same analysis for  
435 tandem gene duplicates as identified by SynMap [29,55]; this time testing the null  
436 hypothesis that the property for a gene to be tandemly duplicated and to belong to a  
437 given GO category are independent. We excluded genes on scaffolds and reported  
438 enriched/depleted categories if they remain significant ( $p < 0.05$ ) after multiple test  
439 correction (fdr).

440

441 **Table 1:** 2x2 contingency table obtained by classifying genes into 2 categorical  
442 variables. The letters denote the number of genes for a given category (e.g. “a” denotes  
443 the number of retained genes annotated with the tested GO category).

444 -----

	<b>GO</b>	<b>not-GO</b>	<b>SUM</b>
<b>retained</b>	a	b	a+b*
<b>not-retained</b>	c	d	c+d
<b>SUM</b>	a+c	b+d	N=total number of genes =29550 (across 7 chromosomes)

449 -----

451 \*equal to 1302 and 6972 for candidate WGD-derived paralogs and tandem gene  
452 duplicates, respectively.

453

## 454 **List of abbreviations**

455 WGD: whole genome duplication; Ks: the number of synonymous substitutions per  
456 synonymous site; GO: Gene Ontology.

457

## 458 **Declarations**

### 459 **Availability of data and material**

460 The columbine, grape and cacao genomes are available at the CoGE portal for the  
461 synteny analyses with the id numbers 28620, 19990 and 25287, respectively [55].

### 462 **Competing interests**

463 The authors declare no competing interests.

### 464 **Funding**

465 G.A. was supported by the Vienna Graduate School of Population Genetics (Austrian  
466 Science Fund, FWF: DK W1225-B20).

### 467 **Authors' contributions**

468 G.A. performed all analyses. G.A. and M.N. wrote the manuscript.

### 469 **Acknowledgements**

470 We thank Robin Burns and Claus Vogl for their comments on the manuscript; Daniel  
471 Gómez Sánchez and Benjamin Jaegle for the fruitful discussions.

472

## 473 **References**

474 1. Wendel JF. The wondrous cycles of polyploidy in plants. *Am J Bot.* 2015;102:1753–6.

475 2. Van de Peer Y, Mizrachi E, Marchal K. The evolutionary significance of polyploidy.  
476 *Nat Rev Genet.* 2017;18:411–24.

477 3. Jiao Y, Wickett NJ, Ayyampalayam S, Chanderbali AS, Landherr L, Ralph PE, et al.  
478 Ancestral polyploidy in seed plants and angiosperms. *Nature.* 2011;473:97–100.

479 4. Leitch IJ, Bennett MD. Genome downsizing in polyploid plants. *Biol J Linn Soc Lond.*  
480 Oxford University Press; 2004;82:651–63.

481 5. Comai L. The advantages and disadvantages of being polyploid. *Nat Rev Genet.*  
482 2005;6:836–46.

- 483 6. Thomas BC, Pedersen B, Freeling M. Following tetraploidy in an Arabidopsis  
484 ancestor, genes were removed preferentially from one homeolog leaving clusters  
485 enriched in dose-sensitive genes. *Genome Res.* 2006;16:934–46.
- 486 7. Otto SP. The evolutionary consequences of polyploidy. *Cell.* 2007;131:452–62.
- 487 8. Renny-Byfield S, Chester M, Kovařík A, Le Comber SC, Grandbastien M-A, Deloger  
488 M, et al. Next generation sequencing reveals genome downsizing in allotetraploid  
489 *Nicotiana tabacum*, predominantly through the elimination of paternally derived  
490 repetitive DNAs. *Mol Biol Evol.* 2011;28:2843–54.
- 491 9. Murat F, Zhang R, Guizard S, Gavranović H, Flores R, Steinbach D, et al. Karyotype  
492 and gene order evolution from reconstructed extinct ancestors highlight contrasts in  
493 genome plasticity of modern rosid crops. *Genome Biol Evol.* 2015;7:735–49.
- 494 10. Murat F, Armero A, Pont C, Klopp C, Salse J. Reconstructing the genome of the most  
495 recent common ancestor of flowering plants. *Nat Genet.* 2017;49:490–6.
- 496 11. Badouin H, Gouzy J, Grassa CJ, Murat F, Staton SE, Cottret L, et al. The sunflower  
497 genome provides insights into oil metabolism, flowering and Asterid evolution. *Nature.*  
498 2017;546:148–52.
- 499 12. Blanc G, Wolfe KH. Functional divergence of duplicated genes formed by polyploidy  
500 during Arabidopsis evolution. *Plant Cell.* 2004;16:1679–91.
- 501 13. De Smet R, Adams KL, Vandepoele K, Van Montagu MCE, Maere S, Van de Peer Y.  
502 Convergent gene loss following gene and genome duplications creates single-copy  
503 families in flowering plants. *Proc Natl Acad Sci U S A.* 2013;110:2898–903.
- 504 14. Woodhouse MR, Schnable JC, Pedersen BS, Lyons E, Lisch D, Subramaniam S, et al.  
505 Following Tetraploidy in Maize, a Short Deletion Mechanism Removed Genes  
506 Preferentially from One of the Two Homeologs. *PLoS Biol. Public Library of Science;*  
507 2010;8:e1000409.
- 508 15. Schnable JC, Springer NM, Freeling M. Differentiation of the maize subgenomes by  
509 genome dominance and both ancient and ongoing gene loss. *Proc Natl Acad Sci U S A.*  
510 2011;108:4069–74.
- 511 16. Tang H, Woodhouse MR, Cheng F, Schnable JC, Pedersen BS, Conant G, et al.  
512 Altered patterns of fractionation and exon deletions in *Brassica rapa* support a two-step  
513 model of paleohexaploidy. *Genetics.* 2012;190:1563–74.
- 514 17. Hodges SA, Kramer EM. Columbines. *Curr Biol.* 2007;17:R992–4.
- 515 18. Soltis DE, Smith SA, Cellinese N, Wurdack KJ, Tank DC, Brockington SF, et al.  
516 Angiosperm phylogeny: 17 genes, 640 taxa. *Am J Bot.* 2011;98:704–30.
- 517 19. Filiault D, Ballerini E, Mandakova T, Akoz G, Derieg N. The *Aquilegia* genome:  
518 adaptive radiation and an extraordinarily polymorphic chromosome with a unique

- 519 history. bioRxiv [Internet]. biorxiv.org; 2018; Available from:  
520 <https://www.biorxiv.org/content/early/2018/02/12/264101.abstract>
- 521 20. Lynch M, Conery JS. The evolutionary fate and consequences of duplicate genes.  
522 *Science*. 2000;290:1151–5.
- 523 21. Blanc G, Wolfe KH. Widespread paleopolyploidy in model plant species inferred  
524 from age distributions of duplicate genes. *Plant Cell*. 2004;16:1667–78.
- 525 22. Cui L, Wall PK, Leebens-Mack JH, Lindsay BG, Soltis DE, Doyle JJ, et al.  
526 Widespread genome duplications throughout the history of flowering plants. *Genome*  
527 *Res*. 2006;16:738–49.
- 528 23. Vanneste K, Baele G, Maere S, Van de Peer Y. Analysis of 41 plant genomes supports  
529 a wave of successful genome duplications in association with the Cretaceous–Paleogene  
530 boundary. *Genome Res*. 2014;24:1334–47.
- 531 24. Tiley GP, Ané C, Burleigh JG. Evaluating and Characterizing Ancient Whole-  
532 Genome Duplications in Plants with Gene Count Data. *Genome Biol Evol*. 2016;8:1023–  
533 37.
- 534 25. Doyle JJ, Egan AN. Dating the origins of polyploidy events. *New Phytol*.  
535 2010;186:73–85.
- 536 26. Jiao Y, Paterson AH. Polyploidy-associated genome modifications during land plant  
537 evolution. *Philos Trans R Soc Lond B Biol Sci* [Internet]. 2014;369. Available from:  
538 <http://dx.doi.org/10.1098/rstb.2013.0355>
- 539 27. Melters DP, Bradnam KR, Young HA, Telis N, May MR, Ruby JG, et al. Comparative  
540 analysis of tandem repeats from hundreds of species reveals unique insights into  
541 centromere evolution. *Genome Biol*. 2013;14:R10.
- 542 28. Jaillon O, Aury J-M, Noel B, Policriti A, Clepet C, Casagrande A, et al. The grapevine  
543 genome sequence suggests ancestral hexaploidization in major angiosperm phyla.  
544 *Nature*. 2007;449:463–7.
- 545 29. Lyons E, Pedersen B, Kane J, Freeling M. The Value of Nonmodel Genomes and an  
546 Example Using SynMap Within CoGe to Dissect the Hexaploidy that Predates the  
547 Rosids. *Trop Plant Biol*. 2008;1:181–90.
- 548 30. Potato Genome Sequencing Consortium, Xu X, Pan S, Cheng S, Zhang B, Mu D, et  
549 al. Genome sequence and analysis of the tuber crop potato. *Nature*. 2011;475:189–95.
- 550 31. Truco MJ, Ashrafi H, Kozik A, van Leeuwen H, Bowers J, Wo SRC, et al. An Ultra-  
551 High-Density, Transcript-Based, Genetic Map of Lettuce. *G3* . 2013;3:617–31.
- 552 32. Denoeud F, Carretero-Paulet L, Dereeper A, Droc G, Guyot R, Pietrella M, et al. The  
553 coffee genome provides insight into the convergent evolution of caffeine biosynthesis.  
554 *Science*. 2014;345:1181–4.

- 555 33. Bombarely A, Moser M, Amrad A, Bao M, Bapaume L, Barry CS, et al. Insight into  
556 the evolution of the Solanaceae from the parental genomes of *Petunia hybrida*. *Nat*  
557 *Plants*. 2016;2:16074.
- 558 34. Salse J. Ancestors of modern plant crops. *Curr Opin Plant Biol*. 2016;30:134–42.
- 559 35. Ming R, VanBuren R, Liu Y, Yang M, Han Y, Li L-T, et al. Genome of the long-living  
560 sacred lotus (*Nelumbo nucifera* Gaertn.). *Genome Biol*. 2013;14:R41.
- 561 36. Jiao Y, Leebens-Mack J, Ayyampalayam S, Bowers JE, McKain MR, McNeal J, et al.  
562 A genome triplication associated with early diversification of the core eudicots. *Genome*  
563 *Biol*. 2012;13:R3.
- 564 37. Zeng L, Zhang Q, Sun R, Kong H, Zhang N, Ma H. Resolution of deep angiosperm  
565 phylogeny using conserved nuclear genes and estimates of early divergence times. *Nat*  
566 *Commun*. 2014;5:4956.
- 567 38. Zheng C, Chen E, Albert VA, Lyons E, Sankoff D. Ancient eudicot hexaploidy meets  
568 ancestral eurosid gene order. *BMC Genomics*. 2013;14 Suppl 7:S3.
- 569 39. Argout X, Salse J, Aury J-M, Gaultier M, Droc G, Gouzy J, et al. The genome of  
570 *Theobroma cacao*. *Nat Genet*. 2011;43:101–8.
- 571 40. Soltis DE, Visger CJ, Marchant DB, Soltis PS. Polyploidy: Pitfalls and paths to a  
572 paradigm. *Am J Bot*. 2016;103:1146–66.
- 573 41. Seoighe C, Wolfe KH. Yeast genome evolution in the post-genome era. *Curr Opin*  
574 *Microbiol*. 1999;2:548–54.
- 575 42. Tian C-G, Xiong Y-Q, Liu T-Y, Sun S-H, Chen L-B, Chen M-S. Evidence for an  
576 ancient whole-genome duplication event in rice and other cereals. *Yi Chuan Xue Bao*.  
577 2005;32:519–27.
- 578 43. Aury J-M, Jaillon O, Duret L, Noel B, Jubin C, Porcel BM, et al. Global trends of  
579 whole-genome duplications revealed by the ciliate *Paramecium tetraurelia*. *Nature*.  
580 2006;444:171–8.
- 581 44. Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, et al. The B73 maize  
582 genome: complexity, diversity, and dynamics. *Science*. 2009;326:1112–5.
- 583 45. Rodgers-Melnick E, Mane SP, Dharmawardhana P, Slavov GT, Crasta OR, Strauss  
584 SH, et al. Contrasting patterns of evolution following whole genome versus tandem  
585 duplication events in *Populus*. *Genome Res*. 2012;22:95–105.
- 586 46. Birchler JA, Bhadra U, Bhadra MP, Auger DL. Dosage-dependent gene regulation in  
587 multicellular eukaryotes: implications for dosage compensation, aneuploid syndromes,  
588 and quantitative traits. *Dev Biol*. 2001;234:275–88.
- 589 47. Papp B, Pál C, Hurst LD. Dosage sensitivity and the evolution of gene families in

- 590 yeast. *Nature*. 2003;424:194–7.
- 591 48. Freeling M, Thomas BC. Gene-balanced duplications, like tetraploidy, provide  
592 predictable drive to increase morphological complexity. *Genome Res*. 2006;16:805–14.
- 593 49. Birchler JA, Veitia RA. Gene balance hypothesis: connecting issues of dosage  
594 sensitivity across biological disciplines. *Proc Natl Acad Sci USA*. 2012;109:14746–53.
- 595 50. Guo L, Winzer T, Yang X, Li Y, Ning Z, He Z, et al. The opium poppy genome and  
596 morphinan production. *Science* [Internet]. 2018; Available from:  
597 <http://dx.doi.org/10.1126/science.aat4096>
- 598 51. Malacarne G, Perazzolli M, Cestaro A, Sterck L, Fontana P, Van de Peer Y, et al.  
599 Deconstruction of the (paleo)polyploid grapevine genome based on the analysis of  
600 transposition events involving NBS resistance genes. *PLoS One*. 2012;7:e29762.
- 601 52. O’Kane SL, Schaal BA, Al-Shehbaz IA. The Origins of *Arabidopsis suecica*  
602 (Brassicaceae) as Indicated by Nuclear rDNA Sequences. *Syst Bot. American Society of*  
603 *Plant Taxonomists*; 1996;21:559–66.
- 604 53. Matsuoka Y. Evolution of polyploid triticum wheats under cultivation: the role of  
605 domestication, natural hybridization and allopolyploid speciation in their  
606 diversification. *Plant Cell Physiol*. 2011;52:750–64.
- 607 54. Friis EM, Pedersen KR, Crane PR. The emergence of core eudicots: new floral  
608 evidence from the earliest Late Cretaceous. *Proc Biol Sci* [Internet]. 2016;283. Available  
609 from: <http://dx.doi.org/10.1098/rspb.2016.1325>
- 610 55. Lyons E, Freeling M. How to usefully compare homologous plant genes and  
611 chromosomes as DNA sequences. *Plant J*. 2008;53:661–73.
- 612 56. Blanc G, Hokamp K, Wolfe KH. A recent polyploidy superimposed on older large-  
613 scale duplications in the *Arabidopsis* genome. *Genome Res*. 2003;13:137–44.
- 614 57. Smith DA, Cordell R, Dillon EM. Infectious texts: Modeling text reuse in nineteenth-  
615 century newspapers. 2013 IEEE International Conference on Big Data. 2013. p. 86–94.
- 616