

# Whole-genome microsynteny-based phylogeny of angiosperms

Tao Zhao (✉ [tao.zhao@nwfufu.edu.cn](mailto:tao.zhao@nwfufu.edu.cn))

Ghent University

Jiayu Xue

Institute of Botany, Jiangsu Province and Chinese Academy of Sciences

Arthur Zwaenepoel

Ghent University

Shu-min Kao

Ghent University

Zhen Li

Department of Plant Biotechnology and Bioinformatics, Ghent University, 9052 Gent, Belgium.

M. Schranz

Wageningen University <https://orcid.org/0000-0001-6777-6565>

Yves Van de Peer

Ghent University <https://orcid.org/0000-0003-4327-3730>

---

## Article

**Keywords:** Synteny network, gene order, matrix representation, phylogeny, angiosperms, maximum-likelihood

**Posted Date:** August 17th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-51378/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Nature Communications on June 9th, 2021. See the published version at <https://doi.org/10.1038/s41467-021-23665-0>.

1 Whole-genome microsynteny-based phylogeny of angiosperms

2

3 Tao Zhao<sup>1,2,3</sup>, Jia-Yu Xue<sup>4</sup>, Arthur Zwaenepoel<sup>1,2</sup>, Shu-Min Kao<sup>1,2</sup>, Zhen Li<sup>1,2</sup>, M. Eric  
4 Schranz<sup>5</sup>, Yves Van de Peer<sup>1,2,6,7</sup>

5

6 <sup>1</sup>Department of Plant Biotechnology and Bioinformatics, Ghent University, Ghent,  
7 Belgium.

8 <sup>2</sup>Center for Plant Systems Biology, VIB, Ghent, Belgium.

9 <sup>3</sup>State Key Laboratory of Crop Stress Biology for Arid Areas/Shaanxi Key Laboratory  
10 of Apple, College of Horticulture, Northwest A & F University, Yangling, 712100,  
11 China.

12 <sup>4</sup>Institute of Botany, Jiangsu Province and Chinese Academy of Sciences, Nanjing,  
13 China.

14 <sup>5</sup>Biosystematics Group, Wageningen University and Research, Wageningen, The  
15 Netherlands.

16 <sup>6</sup>Center for Microbial Ecology and Genomics, Department of Biochemistry, Genetics  
17 and Microbiology, University of Pretoria, Pretoria, South Africa.

18 <sup>7</sup>College of Horticulture, Nanjing Agricultural University, Nanjing, China.

19

20

21 To whom correspondence should be addressed:

22 [yves.vandeppeer@psb.ugent.be](mailto:yves.vandeppeer@psb.ugent.be)

23 [tao.zhao@nwafu.edu.cn](mailto:tao.zhao@nwafu.edu.cn)

24 **Abstract**

25 Plant genomes are generally very complex and dynamic structures, and vary greatly  
26 in size, organization, and architecture. This is mainly due to the often-excessive  
27 numbers of transposable and repetitive elements, as well as to the fact that many  
28 plants are ancient or recent polyploids. Such (recurrent) whole-genome duplications  
29 are usually followed by genomic rearrangements, gene transpositions and gene loss,  
30 making local gene order-based phylogenetic inference particularly challenging.  
31 Nevertheless, microsynteny, i.e. the conservation of local gene content and order,  
32 has been recognized as a valuable and alternative phylogenetic character to  
33 sequence-based characters (nucleotides or amino acids) for the inference of  
34 phylogenetic trees, but to date its application for reconstructing larger phylogenies  
35 has been, for several reasons, limited. Here, by combining synteny network analysis,  
36 matrix representation, and maximum likelihood, we have reconstructed a  
37 microsynteny-based phylogenetic tree for more than 120 available high-quality plant  
38 genomes, representing more than 50 different plant families and 30 plant orders  
39 within the angiosperms. Comparisons with sequence alignment-based trees and  
40 current phylogenetic classifications show that we reconstruct very accurate and  
41 robust phylogenies, albeit with sometimes important alternative sister-group  
42 relationships. For instance, our synteny-based tree positioned Vitales as early-  
43 diverging eudicots, Saxifragales belongs to superasterids, and magnoliids as sister  
44 to monocots. We discuss how synteny-based phylogeny can be complementary to  
45 traditional methods and could provide additional insights into some long-standing  
46 controversial phylogenetic relationships.

47

48 **Key words**

49 Synteny network, gene order, matrix representation, phylogeny, angiosperms,  
50 maximum-likelihood

## 51 **Introduction**

52 Microsynteny (hereafter also simply referred to as synteny), or the local conservation  
53 of gene order and content, provides a valuable means to infer the shared ancestry of  
54 groups of genes and is commonly used to infer the occurrence of ancient polyploidy  
55 events<sup>1,2</sup>, to identify genomic rearrangements<sup>3</sup>, and to establish gene orthology  
56 relationships<sup>4-6</sup>, particularly for large gene families where sequence-based  
57 phylogenetic methods may be inconclusive<sup>7,8</sup>. In addition, microsynteny has also  
58 been used for the inference of phylogenetic relationships. For instance, recently,  
59 Drillon et al. (2020)<sup>9</sup> developed a bottom-up pairwise (partial splits) distance-based  
60 tree reconstruction approach that starts with the identification of breakpoints between  
61 synteny blocks, followed by the identification of partial splits. This method was  
62 successfully applied to 13 vertebrate genomes and 21 yeast genomes<sup>9</sup>. However, in  
63 general, application of synteny data for phylogenetic inference in complex eukaryotic  
64 genomes has been very limited<sup>9-11</sup>. Compared to sequence-based phylogenetic  
65 approaches, methods based on gene order analyses are usually computationally  
66 costly and ancestral genome reconstruction and inference of gene order  
67 rearrangements is a highly challenging combinatorial and algorithmic problem<sup>12-14</sup>.  
68 Most methods hitherto developed can only handle pairs of genomes, or single-  
69 chromosome organisms and organelles, and such methods have therefore only been  
70 applied to simpler genomes such as plastid genomes or bacterial genomes, or  
71 simplified genome datasets or simulated datasets<sup>15,16</sup>.

72 Plant genomes are highly diverse<sup>17,18</sup> and affected by a complex interplay of small-  
73 and large-scale duplication events<sup>19,20</sup>, hybridization<sup>21</sup>, transposable element (TE)  
74 activity<sup>22</sup>, and gene loss<sup>23,24</sup>. Consequently, employing synteny information for  
75 phylogenetic inference from large plant genome datasets has been notoriously  
76 difficult. Ideally, inferring phylogenetic relationships from synteny data should involve  
77 dealing with (1) multiple chromosomes; (2) genome duplications; and (3) multiple  
78 genomes at the same time, which greatly increases the complexity of the problem  
79 and computational cost. Recently, some of us developed a novel approach in which  
80 microsynteny information is converted into a network data structure, and which has  
81 proven to be well suited for evolutionary synteny comparisons among many  
82 eukaryotic nuclear genomes<sup>7,18</sup>. Using this approach, conservation or divergence of  
83 genome structure can be conveniently summarized and reflected by synteny cluster

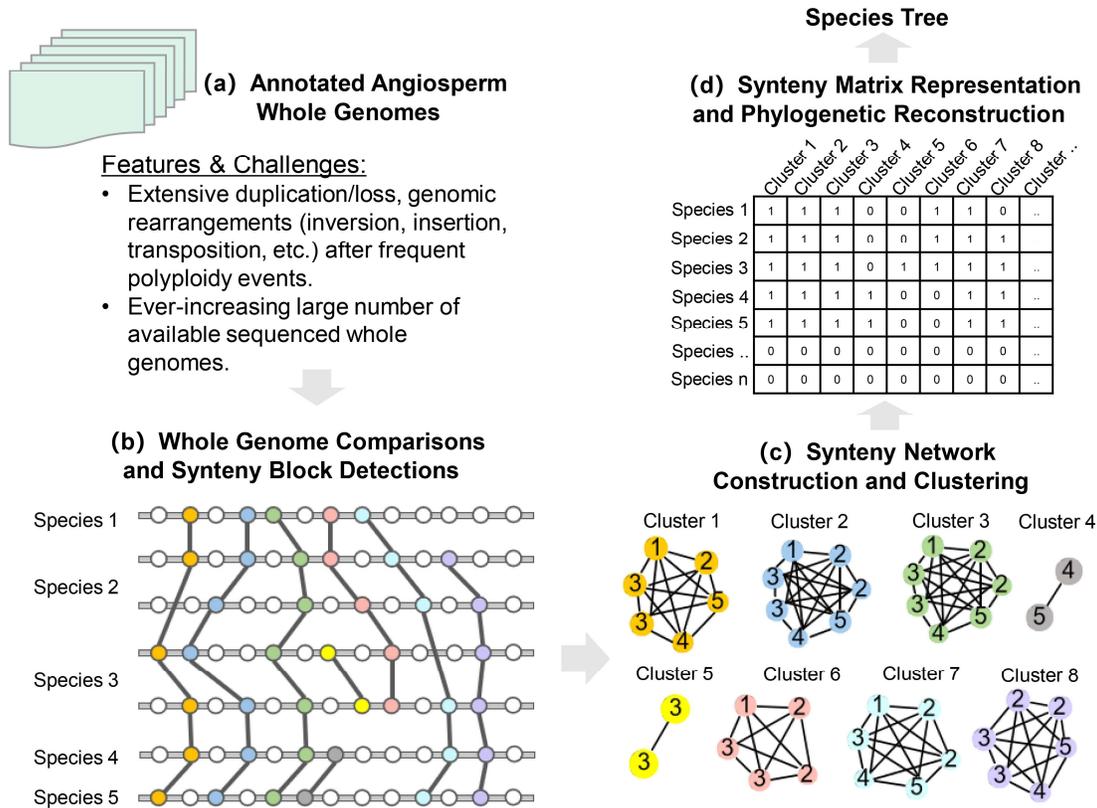
84 sizes and composition. Importantly, the network representation of synteny  
85 relationships provides an abstraction of structural homology across genomes that is  
86 in principle amenable to tree inference using standard approaches from  
87 phylogenetics, however this possibility has hitherto not been explored.

88 Here, we combined synteny networks and their matrix representation with standard  
89 maximum-likelihood based statistical phylogenetics to reconstruct phylogenies for  
90 real-life plant whole genome datasets. We have constructed a well-resolved 'synteny  
91 tree' using currently available representative genomes of flowering plants. The  
92 obtained tree is highly comparable with current phylogenetic classifications, although  
93 some notable differences concerning the phylogenetic positions of magnoliids,  
94 Vitales, Caryophyllales, Saxifragales, and Santalales were observed. We believe  
95 that our approach provides a noteworthy complement to more classical approaches  
96 of tree inference and could help to solve some longstanding problems that may  
97 remain difficult to solve with sequence (alignment) based methods.

## 98 Results

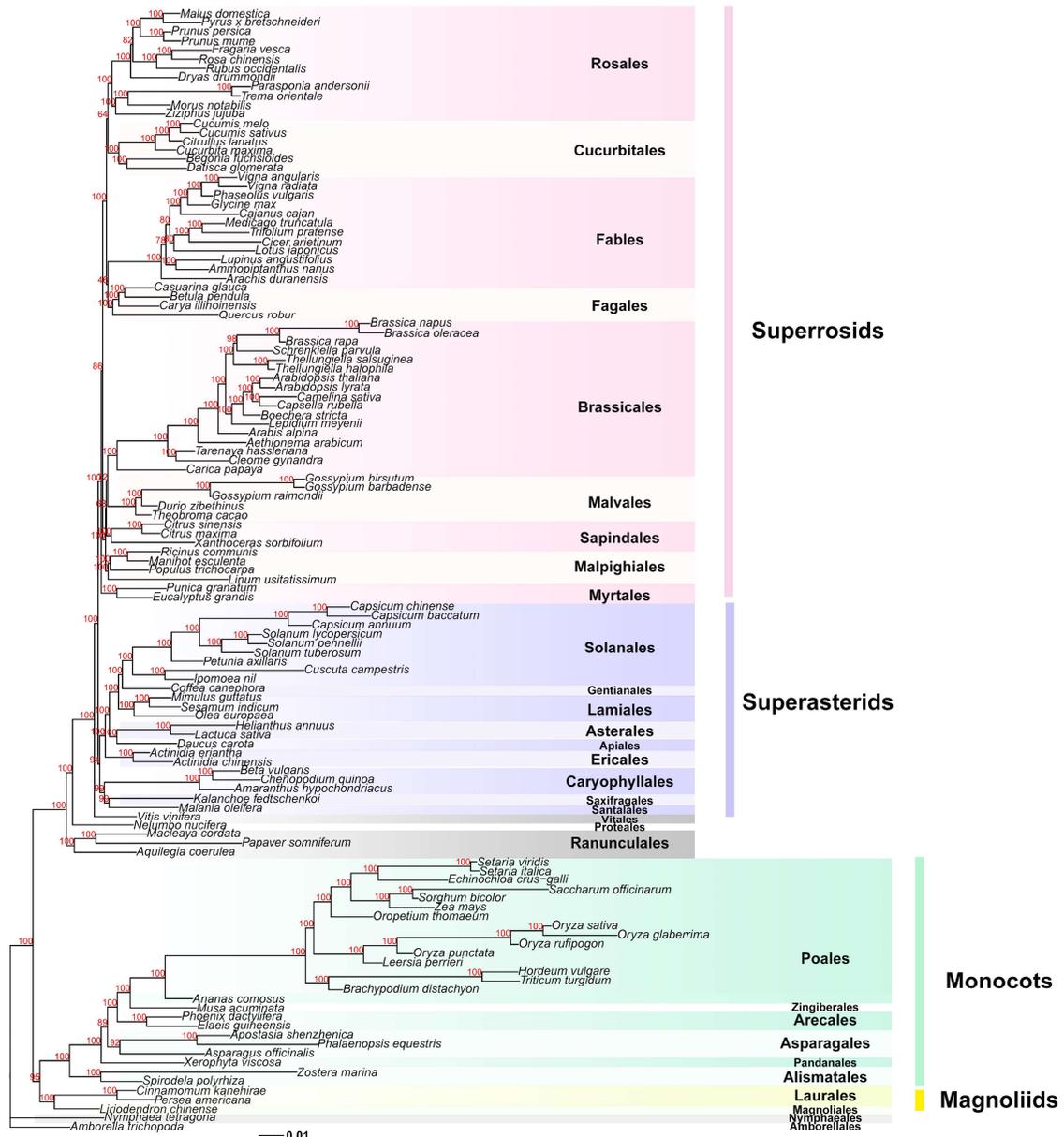
### 99 Highly-resolved microsynteny-based phylogeny for angiosperm genomes

100 After quality control, 123 fully-sequenced plant genomes were used for synteny-  
101 based phylogenetic analysis, which includes synteny network construction and  
102 clustering, matrix representation of synteny followed by maximum likelihood  
103 estimation (Fig. 1, see Methods for details). The size of the matrix used for ML tree  
104 inference is 123 x 137,833, which contains a binary presence/absence coding for  
105 each cluster in the synteny network. The resulting best ML tree demonstrated that  
106 the overall monophyly of most clades was strongly supported, and 110 of the 122  
107 nodes had  $\geq 95\%$  bootstrap support (Fig. 2). *Amborella* was used as the sister group  
108 to other angiosperms in order to root the obtained phylogenies. Nymphaeales was  
109 resolved as a successive lineage right after *Amborella*, again sister to all other  
110 angiosperm genomes (Fig. 2). The monophyly of mesangiosperms was also strongly  
111 supported (Fig. 2). Interestingly, magnoliids (including Laurales and Magnoliales)  
112 form a sister group to monocots (BS = 95%), with the resulting clade sister to the  
113 eudicots (BS = 100%). Notably, our synteny tree strongly favors Vitales as sister to  
114 the rest of the core eudicots, or as another clade of early-diverging eudicots,  
115 branching off after Proteales (*Nelumbo*) (both BS = 100%) (Fig. 2). Moreover, in the  
116 synteny tree, we find Santalales as sister to Saxifragales (BS = 99%), and with both  
117 sister to Caryophyllales (BS = 99%), which in turn was found as sister to all other  
118 asterids (BS = 94%). Synteny trees positioned Myrtales (*Eucalyptus* and *Punica*) as  
119 early-diverging rosids (BS = 100%), and Malpighiales as early-diverging Malvids (BS  
120 = 100%). Almost all the nodes within Brassicales were fully resolved (Fig. 2). The  
121 monophyly of the Nitrogen-Fixing Clade was fully recovered, and supports a  
122 relationship of ((Fagales, Fabales), (Rosales, Cucurbitales)), however with lower  
123 support (46% and 64% BS support, respectively) (Fig. 2).



124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148

**Fig. 1** Whole-genome microsynteny-based species tree inference. **(a)** Whole-genome datasets with all predicted genes are used for phylogeny reconstruction. **(b)** The synteny network approach first conducts all pairwise reciprocal genome comparisons, followed by synteny block detection. All syntenic blocks constitute the synteny network database (see Methods for details). **(c)** We analyzed all synteny clusters after clustering the entire network database. Synteny clusters vary in size and node compositions. Shared genomic rearrangements are reflected by cluster compositions. Specific anchor pairs shared by a lineage/species form specific clusters (e.g. Clusters 4-6). We account for the presence or absence of the same recurring anchors for multiple blocks derived from whole-genome or segmental duplications (e.g. for Species 2 and 3 in Clusters 2, 3, 5, and 8). **(d)** The phylogenomic profiling of all clusters constructs a binary matrix, where rows represent species and columns represent clusters. The synteny matrix comprehensively represents phylogenomic gene order dynamics. It transforms the concept of synteny comparisons from analyzing massive parallel coordinates plots into analyzing profiles of individual clusters/networks. Each cluster stands for a shared homologous 'context'. For example, genes (driven by transposable elements) can be transposed as insertions into new contexts or can be lost from the original context (e.g. genes in Clusters 4-6). As long as such transpositions are shared by different genomes (e.g. genes in Clusters 4 and 6) or within the same genome because of whole genome duplication (e.g. genes in Cluster 5), specific clusters will emerge and corresponding signals will be added to the matrix. This synteny matrix is used as the input for species tree inference by maximum likelihood (referred to as Syn-MRL).



149

150

**Fig. 2** Maximum likelihood (ML) tree for 123 fully-sequenced flowering plant genomes based on the microsynteny approach. The tree is rooted by *Amborella*, and four main clades, i.e. superrosids, superasterids, monocots, and magnoliids are shaded in light-red, light-purple, light-green, and light-yellow, respectively. Ultrafast bootstrapping values are denoted for all the nodes. Names for the different plant orders follow the APG IV classification<sup>25</sup>.

151

152

153

154

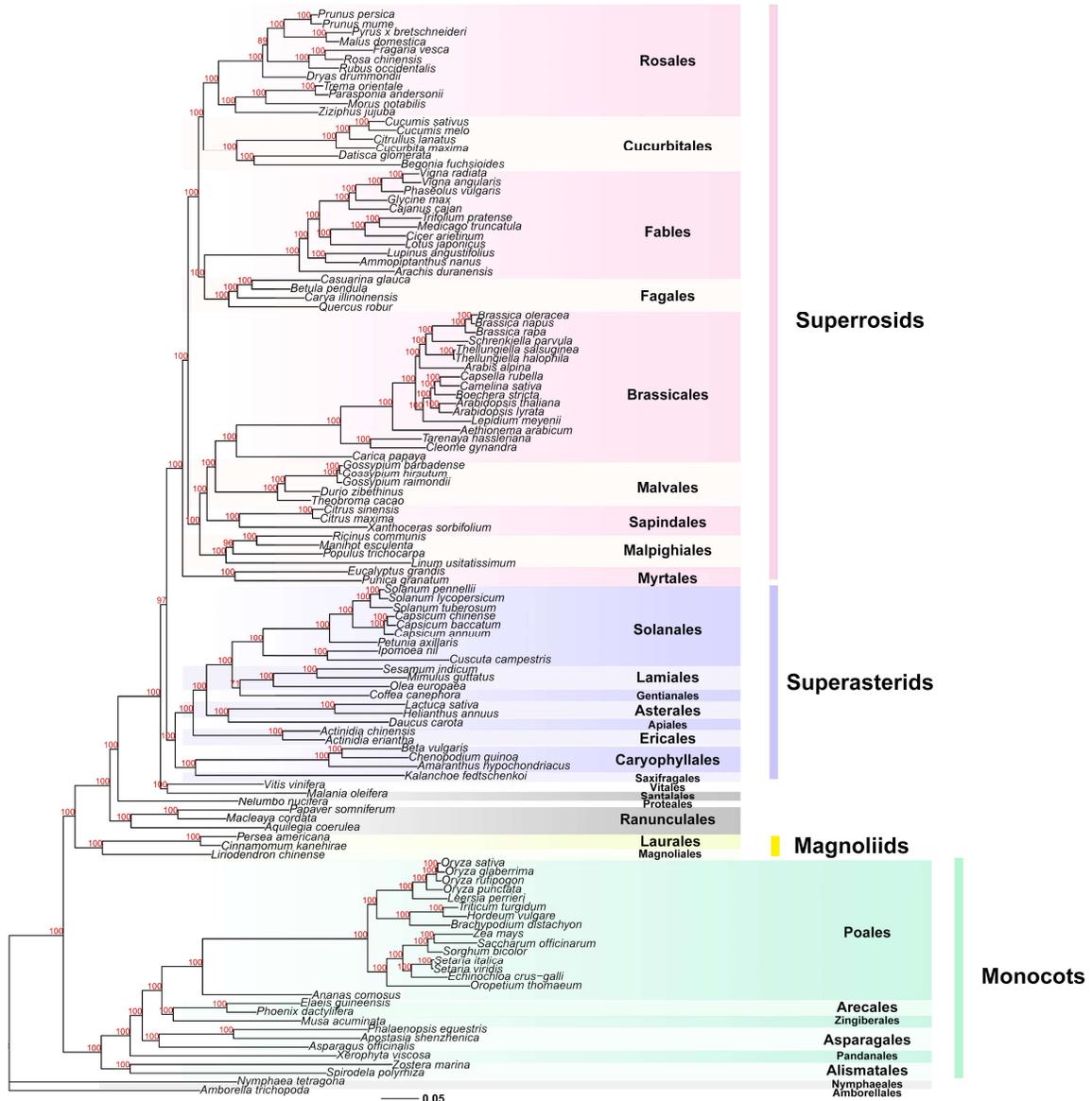
155

156 **Comparison of the synteny tree with sequence alignment-based phylogenies**  
157 **(using whole-genome derived gene markers), and with current phylogenetic**  
158 **classifications**

159 Taking advantage of the large whole-genome plant dataset used in this study, we  
160 also reconstructed phylogenies using widely-adopted sequence alignment-based  
161 approaches for comparison with our microsynteny-based tree. We used BUSCO  
162 gene sets, which are widely adopted in benchmarking genome assembly and  
163 annotation quality. Besides BUSCO, we developed a set of CSSC (Conserved  
164 Single-copy Synteny Clusters) gene markers from a profiling and screening of all  
165 synteny clusters (see Methods). We used supermatrix, supertree, and multispecies  
166 coalescent as representative sequence-based approaches (Supplementary Fig. 1)  
167 for these two sets of gene markers, and generated six phylogenetic trees  
168 (Supplementary Figs. 2-7). Overall, the six sequence alignment-based trees  
169 (hereafter referred to as SA trees) are highly similar, and only differ by a few  
170 bootstrap support values and a few minor sister-group relationships (Supplementary  
171 Figs. 2-7). We used the supermatrix-BUSCO tree as the representative of the SA  
172 trees and compared it with the synteny tree (Fig. 3).

173 Comparing both trees shows that the synteny tree we have obtained is highly  
174 congruent with the SA tree (Figs. 2-3, Supplementary Fig. 8). Interestingly, some  
175 notable differences were found between the two trees, such as the positioning of  
176 magnoliids, Santalales, Zingiberales, and Gentianales (Figs. 2-3, Supplementary Fig.  
177 8). Additional subtle differences were found within the orders, such as in Poales and  
178 Brassicales (Figs. 2-3, Supplementary Fig. 8). We further compared our synteny tree  
179 with the current Angiosperm Phylogeny Group (APG) phylogeny (version IV)<sup>25</sup>, and  
180 the recent 1000 plant transcriptomes (1KP) phylogeny of green plants<sup>20</sup> at the order  
181 level (Supplementary Figs. 9-11). Apart from the orders that do not have  
182 representative genomes yet, our synteny tree generally shows strong congruence to  
183 both the APG and 1KP trees (note however that a certain degree of incongruence  
184 exists between the two latter trees (Supplementary Fig. 9)) (Supplementary Figs. 10-  
185 11). Besides the relationship of magnoliids and monocots and dicots, all  
186 discrepancies are confined to the positions and relationships of Vitales, Santalales,  
187 Saxifragales, Caryophyllales (Supplementary Figs. 9-11), which reflect long-time  
188 controversies in plant systematics and have been well acknowledged in the

189 literature<sup>20,25-33</sup>. Our SA trees agree with the 1KP tree on magnoliids being sister to  
190 eudicots, which is different from the synteny tree (magnoliids sister to monocots) and  
191 APG tree (magnoliids sister to both monocots and eudicots) ([Supplementary Figs. 8-  
192 11](#)); the APG tree and the 1KP tree resolve Vitales and Saxifragales as early-  
193 diverging lineages within Rosids, whereas our SA trees (except the ASTRAL-  
194 BUSCO tree) and synteny tree favor Vitales as sister to the rest of the core eudicots,  
195 and Saxifragales cluster within Superasterids ([Supplementary Figs. 2-11](#)). We tested  
196 whether our synteny matrix (under the current synteny network construction  
197 parameters) would reject representative alternative topologies. To this end, we used  
198 the ‘approximately unbiased’ (AU) test<sup>34</sup> to evaluate the support for alternative  
199 topologies regarding the positions of magnoliids, Vitales, Saxifragales,  
200 Caryophyllales, and Myrtales ([Supplementary Fig. 12](#)). The test shows that our ML  
201 synteny tree was found to be significantly better than the alternative topologies with  
202 all alternative topologies rejected at the  $p = 0.05$  level, except for the scenario where  
203 magnoliids are sister to eudicots ( $p = 0.146$ , [Supplementary Fig. 12a](#)). This means  
204 that, at least based on our synteny matrix and the (admittedly *ad hoc*) Markov model  
205 of character evolution, alternative scenarios such as magnoliids sister to both  
206 monocots and eudicots ([Supplementary Fig. 12b](#)), Vitales ([Supplementary Fig. 12c](#)),  
207 or Saxifragales ([Supplementary Fig. 12d](#)) as early-diverging rosids, and so on, were  
208 significantly less well supported.



209

210 **Fig. 3** Maximum likelihood (ML) tree based on the concatenation of protein  
 211 alignments of BUSCO genes. This tree was used as the representative SA  
 212 (sequence-alignment based) tree.

## 213 **Discussion**

214 Phylogenetic trees serve many purposes and are an indispensable tool for the  
215 interpretation of evolutionary trends and changes. So far, abundant tools and  
216 sophisticated substitution models have been developed for molecular sequence  
217 alignment-based phylogenetic inference. In contrast, although it is well  
218 acknowledged that there is phylogenetic signal in gene order dynamics<sup>4,35-39</sup>,  
219 efficient tools for analyzing phylogenomic synteny properties, reconstructing  
220 ancestral genomes and inference of genome rearrangements, particularly for large  
221 datasets, remain scarce. In this study, we present an approach that bypasses the  
222 challenging combinatorial problem of inferring genome rearrangements and  
223 ancestral genome organization by integrating information on pairwise homologous  
224 genomic organization across many plant genomes in a network representation. As in  
225 classical orthogroup inference, we perform similarity searches across multiple  
226 genomes and cluster a network representation thereof, but our graph representation  
227 now also includes the genomic context (i.e. synteny information), and the resulting  
228 clusters are therefore defined by sequence level homology and structural homology  
229 at a micro-synteny level. The idea developed in the present paper is that each  
230 cluster, no matter how conserved, family- or lineage-specific, contains a phylogenetic  
231 signal that can readily be used for tree inference. By reducing whole-genome  
232 syntenic comparisons to the synteny network representation, information with regard  
233 to the actual syntenic contexts of the genes is ignored while only relations of shared  
234 syntenic contexts between genes are retained. By identifying clusters in the resulting  
235 synteny network, evolutionarily relevant homologous features of genome structure  
236 across species are efficiently captured while we are able to abstract from much of  
237 the complex features of genome evolution.

238

## 239 **Accuracy and reliability of the synteny tree**

240 We exploited the total information residing in all synteny clusters by encoding their  
241 phylogenomic synteny pattern profiles into a large binary matrix (synteny matrix) (Fig.  
242 1). Standard ML-based tree inference was then applied to reconstruct a species tree  
243 (Fig. 1). The process is similar to the MRL (Matrix Representation with Likelihood)  
244 supertree method (which uses the same data matrix as Matrix Representation with

245 Parsimony, but with ML-based inference, yielding higher accuracy<sup>40,41</sup>), except that  
246 MRL is based on a set of input trees, in contrast with our synteny matrix based on  
247 synteny clusters. To avoid misunderstanding, in this study we confine the usage of  
248 ‘MRL’ to the supertree method; and we refer to our ‘supercluster’ approach as Syn-  
249 MRL (Synteny Matrix Representation with Likelihood).

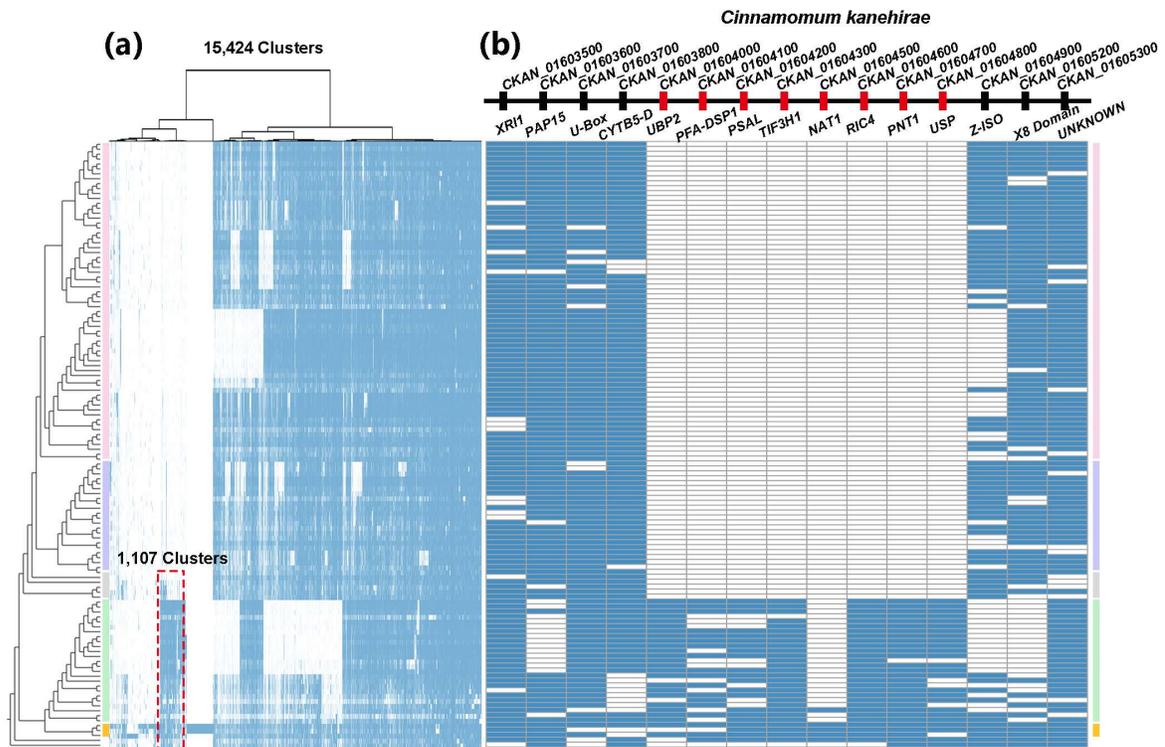
250 For the Syn-MRL approach, we use the Mk model, which is a Jukes-Cantor type  
251 model for discrete morphological data<sup>42</sup>. Each column of the data matrix is regarded  
252 as an independent character. The two-state binary encoding for each column  
253 represents the two groups of related species regarding that specific character. There  
254 is no further ordering or special weighing for the elements in the data matrix. The Mk  
255 model may arouse some concern regarding the symmetric two-state model (i.e.,  
256 there is an equal probability of changing from state 0 to state 1 and from state 1 to  
257 state 0), as generally, time-reversible Markov models of evolution may not be ideal  
258 for the Syn-MRL approach. However, a first and important observation was that  
259 these models do seem to result in very reasonable well-supported phylogenies, in a  
260 similar vein as the MRL approach<sup>40</sup>. We hypothesize this is because, firstly, we may  
261 reasonably suspect that for our data, state 0 or 1 has the same probability of being  
262 the ancestral or derived state, as we can hypothesize that the emergence of a new  
263 lineage-specific synteny cluster and loss of another, are due to the same processes,  
264 e.g. transposition. Secondly, although the Mk model allows numerous state changes,  
265 in practice it yields trees with identical likelihood scores compared to a modified Mk  
266 model where character states can change only once or not at all<sup>42</sup>. This suggests  
267 that more biologically plausible non-reversible models, where for instance the re-  
268 emergence of a synteny cluster (a secondary 0→1 transition after the *initial*  
269 emergence and subsequent loss of the cluster) occurs at a different rate than the  
270 initial 0→1 transition, might not result in a substantially better fit. Nevertheless, we  
271 believe developing other probabilistic models for synteny cluster evolution is a fruitful  
272 avenue for further research.

273 Despite much progress in the last two decades, it is currently well-acknowledged that  
274 some phylogenetic relationships in plants remain especially controversial. For  
275 example, the relationships among monocots, eudicots, magnoliids, Ceratophyllaceae,  
276 and Chloranthales remain unsolved<sup>20,25,43</sup>, while also the relationships among core  
277 rosids, asterids, Saxifragales, Vitales, Santalales, and Caryophyllales have been

278 enigmatic<sup>20,25-29</sup>. We have compared our synteny tree with SA trees and state-of-art  
279 classifications (represented by APG (IV) and 1KP) in a general way ([Supplementary](#)  
280 [Figs. 8-11](#)). Overall, our synteny tree showed great accuracy and congruence on the  
281 classification of the major lineages and clades. On the other hand, the synteny tree  
282 also provides alternative sister-group relationships for some of the recalcitrant clades  
283 mentioned higher. It should be noted that it is not our intention to argue that our tree  
284 is the 'true' tree, finally resolving the contentious relationships between some of the  
285 plant clades, but rather to provide a novel way to reconstruct and consider large-  
286 scale gene order-based phylogenetic trees.

287 The phylogenetic position of magnoliids has long been discussed<sup>43</sup>, and even  
288 recently, based on whole-genome information of new magnoliid genomes, different  
289 sister-group relationships for magnoliids have been proposed<sup>31-33,44</sup>. Also, a sister-  
290 group relationship for magnoliids and monocots, as suggested in the current study,  
291 has been suggested before, based on different approaches and data<sup>45-49</sup>. Based on  
292 microsynteny, we have explored the relationships of magnoliids, monocots, and  
293 eudicots in more detail. Our synteny-based approach provides a means to  
294 investigate phylogenetic signal in the data matrix and trace those back to the  
295 genomic regions where differential gene order arrangements are located. To this end,  
296 we focused on the 'submatrix' of synteny profiles for magnoliids and extracted  
297 15,424 magnoliid-associated synteny clusters ([Fig. 4a](#)). A hierarchical clustering of  
298 the phylogenomic profile of these synteny clusters showed 1,107 synteny clusters  
299 may related to the grouping of magnoliids and monocots ([Fig. 4a](#)). To validate their  
300 contribution to the final ML tree, we first removed these signals from the entire matrix  
301 and reconstructed the phylogenetic tree, after which the species tree obtained favors  
302 magnoliids as sister to eudicots (BS = 100%) ([Supplementary Fig. 13](#)). To further  
303 understand the genomic distribution of specific genes in these clusters joining  
304 magnoliids and monocots, we reorganized the cluster profiles according to the  
305 chromosome gene arrangement of a magnoliid representative (*Cinnamomum*  
306 *kanehirae*)<sup>31</sup>([Supplementary Table S2](#)). In doing so, we observed a number of  
307 'signature' blocks consisting of specific anchor pairs that are shared by monocots  
308 and magnoliids with exclusion of eudicots ([Supplemental Table S2-sheet 1](#), contexts  
309 with highlighted yellow rows). As an example, we highlight a synteny context of 15  
310 genes where 8 genes (highlighted red) are only found in synteny between magnoliid

311 and monocot genomes (Fig. 4b), with flanking genes generally conserved across  
 312 angiosperms (highlighted in blue) (Fig. 4b, Supplemental Table S2-sheet2). However,  
 313 an alternative explanation could be that the synteny context is lost in eudicots, from  
 314 which it might be ‘wrongly’ concluded that monocots and magnoliids share some  
 315 ‘derived’ characters and therefore share a common ancestry. Also, some synteny  
 316 contexts are shared with basal eudicots (Ranunculales) from the 1,107 clusters  
 317 (Supplementary Table S2). Nevertheless, with more representative genomes of  
 318 Chloranthales, Ceratophyllales, and early-diverging angiosperms to be added into  
 319 the analysis, a better resolution of the genomic rearrangements for magnoliids and  
 320 related lineages could be obtained.



321  
 322 **Fig. 4** Magnoliids-associated signals and a representative example of  
 323 phylogenetically informative microsynteny. (a) Hierarchical clustering (ward.D) of  
 324 15,424 magnoliids-associated cluster profiles based on jaccard distance. On the far-  
 325 left, the synteny-based species tree is displayed (same as Fig. 1). Superrosids,  
 326 superasterids, early-diverging eudicots, monocots, and magnoliids are shaded in  
 327 light-red, light-purple, light-grey, light-green, and light-yellow, respectively. 1,107  
 328 clusters supporting a grouping of magnoliids and monocots (Supplemental Table S2).  
 329 (b) One example from all supporting signals. A fifteen-gene context in the genome of  
 330 *Cinnamomum kanehirae* (a magnoliid) shows eight neighboring genes (highlighted in  
 331 orange) only present in magnoliids and monocot genomes, while the flanking genes  
 332 (colored blue) are generally conserved angiosperm-wide.

333 Apart from the contentious phylogenetic relationship of magnoliids with related  
334 lineages, evolutionary relationships of early-diverging core eudicots, such as Vitales,  
335 Saxifragales, Santalales, and Caryophyllales as suggested by microsynteny are  
336 strongly supported compared to competing hypotheses under the Syn-MRL  
337 approach (as reflected by the AU test statistic). For example, our synteny-based  
338 phylogenies strongly support Vitales as early-diverging, right after Proteales  
339 (*Nelumbo*). Trees based on concatenation, MRL, and ASTRAL-CSSC trees in  
340 general also support Vitales as early-diverging but also forming a sister-group  
341 relationship with Santalales (Supplementary Figs. 2-5, and 7). Several studies have  
342 already reported the positioning of Vitales as sister to the core-eudicots or as early-  
343 diverging eudicots, based on mitochondrial genes or genomic data<sup>50,51</sup>, assembly of  
344 multi-nuclear genes<sup>52</sup>, and large-scale transcriptome data<sup>53,54</sup>. The recent 1KP study  
345 also observed substantial gene-tree discordance for Vitales from analyses employing  
346 coalescent and supermatrix approaches, as well as plastid genomes<sup>20</sup>. Moreover, a  
347 recent analysis involving syntenic comparisons of the columbine (*Aquilegia*) and  
348 grapevine genomes revealed that the gamma palaeohexaploidy at the root of the  
349 core-eudicots is possibly a result of hybridization between a tetraploid and a diploid  
350 species (Aköz and Nordborg, 2019). The hybrid origin of core-eudicots could help us  
351 better understand the phylogenetic incongruence of the early-diverging groups within  
352 core-eudicots including Vitales<sup>55</sup>.

353

#### 354 **Insights towards ancestral introgressive hybridization**

355 Within the same plant family or order, our synteny tree also revealed some subtle  
356 differences with the SA trees (Figs 2-3, Supplementary Fig. 8). For example, the  
357 position of *Boechnera* and *Arabis* in Brassicaceae, and the position of the Pooideae  
358 clade (including wheat, barley, and *Brachypodium*) in Poales (Supplementary Fig. 8).

359 Very interestingly, a study has comprehensively investigated the particular species  
360 branching order of three monophyletic groups in the Brassicaceae: *Boechnera* (Clade  
361 B), *Capsella* and *Camelina* (Clade C), and *Arabidopsis* (Clade A)<sup>56</sup>. Forsyth et al.  
362 (2018) explicitly proved the existence of massive nuclear introgression between  
363 Clades B and C, which largely reduced the sequence divergence between them. As  
364 a result, the majority of single-copy gene trees strongly support the branching order

365 of (A,(B,C)), while the true branching order is supposed to be (B,(A,C))<sup>56</sup>. Their study  
366 employed multiple approaches and measurements to show that neither gene  
367 duplication and loss nor ILS, nor phylogenetic noise can adequately explain such  
368 incongruence, except introgressive hybridization<sup>56</sup>. Interestingly, our results are  
369 congruent with the findings of Forsythe et al. (2018). Alignment-based trees  
370 (supermatrix, MRL, and ASTRAL) all support (A,(B,C)) ([Supplementary Fig.s 2-7](#)),  
371 whereas only the synteny tree supports the 'true' species branching order (B,(A,C)).

372 Regarding *Arabidopsis* in Brassicaceae, another recent study identified shared genomic  
373 block associations between *Aethionema* (which is sister to all other Brassicaceae)  
374 and *Arabidopsis* by extending the analysis of Brassicaceae genomic blocks to  
375 *Aethionema*<sup>57</sup>. Their study of macrosynteny supports Arabideae as the next  
376 diverging clade following *Aethionema* in Brassicaceae<sup>57</sup>, which is also consistent  
377 with our synteny tree.

378 For Pooideae in Poales, in order to test whether potential bias exists caused by  
379 specific genome(s), we further compared the clustering of Pooideae (wheat, barley,  
380 and Brachypodium) with the rest of the Poaceae family by removing one or two  
381 genomes from Pooideae from the synteny matrix and rebuilding phylogenetic trees.  
382 We obtained consistent results and topologies regarding the positioning of the clade  
383 of Pooideae ([Supplementary Figs. 14-17](#)). To our knowledge, similar relationships  
384 regarding the PACMAD and the BEP clades have rarely been reported, however, the  
385 comparative phylogenomic large-scale gene family expansion and contraction  
386 analysis from the wheat genome project showed a similar clustering pattern of  
387 species based on gene family profiles, thus provided some evidence in support of  
388 our result (Figure 4A of International Wheat Genome Sequencing Consortium,  
389 2018<sup>58</sup>). If the synteny tree is indeed true, it would reshape our understanding of the  
390 evolution of some of the most important crops including wheat, rice, and maize, as  
391 well as the origin of the C<sub>4</sub> lineages<sup>59</sup>.

392 The notion that species boundaries can be obscured by introgressive hybridization is  
393 increasingly accepted<sup>56,60-62</sup>. Thus, it is plausible that the incongruences discussed  
394 above may indicate an effect from introgression caused by (recent or more ancestral)  
395 hybridization. Since the focal point of synteny information is to reflect genomic  
396 structure variation instead of gene sequence changes, we believe it is plausible that  
397 genome-level synteny-based phylogeny inference may avoid the bias caused by

398 (extensive) introgression for sequence-based species tree reconstruction and can  
399 reflect the ‘true’ (or perhaps better, ‘main’) species branching order. Also, TE  
400 mobilization often follows hybrid speciation or introgression<sup>63-65</sup> which leads to gene  
401 transpositions. However, our approach can capture potential gene transpositions and  
402 rearrangements as novel anchor pairs within synteny blocks and convert those into  
403 phylogenetic signals (explained in Fig. 1). Nevertheless, in the case of true hybrid  
404 speciation, where a bifurcating species tree does not exist, our method can be  
405 similarly unsuitable as other phylogenetic inference methods that do not allow  
406 phylogenetic network inference. In any case, the whole-genome synteny approach  
407 could provide unique opportunities towards developing phylogenetic inference  
408 methods that are robust to introgressive hybridization<sup>66</sup>.

409

#### 410 **Caveats and conclusion**

411 Here, we have presented a methodological roadmap to reconstruct species trees  
412 based on synteny information from large volumes of available whole-genome data,  
413 which can be applied to any set of genomes. However, it should be noted that our  
414 approach depends on the quality of genome assemblies and their gene annotations,  
415 which is the basis for synteny detection. Also, parameter settings for synteny  
416 detection should be tested and compared beforehand (here we adopt the  
417 parameterization found to be most appropriate for the analysis of angiosperm  
418 genomes in a previous study<sup>18</sup>). One should consider the evolutionary distance and  
419 genomic properties of the added genomes, as synteny conservation can be scarce  
420 for distantly related species. For example, caution should be taken when comparing  
421 a gymnosperm or fern genome to an angiosperm genome because in such  
422 comparisons, microsynteny might be inadequate due to extensive genome  
423 rearrangements. On the contrary, for highly similar genomes, our approach might  
424 simply not provide enough resolution due to the lack of informative rearrangements.  
425 However, in such cases, stricter synteny detection parameter settings, as well as  
426 consideration of gene orientation might help to increase the number of informative  
427 signals for resolving the tree topology.

428 To conclude, using a character matrix derived from a network representation of  
429 pairwise microsynteny relations, we here explored a maximum likelihood approach to

430 reconstruct phylogenies using genome structure data across a large set of  
431 angiosperm genomes. Our resulting synteny-based species tree showed high  
432 resolution and strong consistency with phylogenies of angiosperms using more  
433 classical methods to infer tree topologies and current classifications, although some  
434 notable differences were identified. We hope that our approach might offer a  
435 complementary way to consider and evaluate ambiguous phylogenetic relationships.  
436 Furthermore, as more and more high-quality genomes from underrepresented plant  
437 phyla are becoming available at increasing rates, we expect our approach to become  
438 more sensitive and informative in future applications.

439 **Methods**

440 **Genome resources**

441 Reference genomes were obtained from public repositories, including Phytozome,  
442 CoGe, GigaDB, and NCBI. For each genome, we downloaded FASTA format files  
443 containing protein sequences of all predicted gene models and the genome  
444 annotation files (GFF/BED) containing the positions of all the genes. We modified all  
445 peptide sequence files and genome annotation GFF/BED files with corresponding  
446 species abbreviation identifiers. After constructing our synteny network database and  
447 clustering (see further), poor quality genomes could be relatively easily identified  
448 (see further), and were removed from the database for further analysis. After quality  
449 control, the final list of genomes used in the current study and related information for  
450 each genome can be found in [Supplemental Table S1-Sheet1](#), genomes that were  
451 filtered out due to low contiguity were listed in [Supplemental Table S1-Sheet2](#).

452 We acknowledge that our taxon sampling is limited to currently available high-quality  
453 genomes and thus many important lineages could not be included. In our current  
454 study, we have included genome sequences of as many representative angiosperm  
455 families and orders as possible  
456 (<https://www.plabipd.de/portal/web/guest/angiosperm-phylogenetic-view>). However,  
457 some important plant lineages still lack well-assembled genomes (such as for  
458 Dilleniales and Chloranthales), while others only have few representatives (such as  
459 Santalales and Saxifragales), and still others are relatively 'over-represented' (such  
460 as Brassicales, Fabales, and Poales). For 'over-represented' orders, we kept all  
461 qualified genomes because this provides an opportunity to test the resolution and  
462 robustness of our approach at different levels (e.g. order-level, family-level, and  
463 species-level).

464 We manually downloaded each genome and checked the completeness of the  
465 annotation files. Selecting genomes to be used for the analysis initially not based on  
466 certain cutoffs, such as N50 or BUSCO completeness. First, candidate genomes  
467 were all used for the synteny network construction. Next, 'genome quality' is checked  
468 manually in the phylogenomic profiling plot/matrix (of all synteny clusters, where  
469 rows represent species/genomes and columns are clusters). Genomes of poor  
470 completeness and contiguity - indicated by lighter rows (for an example, see Figure 5

471 of Zhao and Schranz, 2019<sup>18</sup>, rows indicated by black arrows) - are removed from  
472 the microsynteny network. After this step, 123 fully sequenced plant genomes were  
473 used for further analysis (Supplemental Table 1). The overall sampling covered 31  
474 orders and 52 different families of angiosperms. We are aware that for some orders  
475 there is an overrepresentation (Brassicales, Fabales, and Poales), while for orders  
476 such as Santalales, Saxifragales, and Gentianales, there is only one representative.  
477 Moreover, for some genera, several genomes are included, for example for *Oryza*,  
478 *Solanum*, *Gossypium*, and *Brassica*. Among those, *B. napus*, *G. herbaceum*, and *G.*  
479 *arboreum* are polyploid (allopolyploid) genomes.

480

### 481 **Pipeline for whole-genome microsynteny-based phylogenetic inference**

482 Our synteny-based phylogenetic reconstruction approach includes four main steps,  
483 in turn namely phylogenomic synteny network construction, network clustering,  
484 matrix representation, and maximum-likelihood estimation. Together we call our  
485 approach 'Syn-MRL' for short.

486 The synteny network construction consists of two main steps: first, all-vs-all  
487 reciprocal annotated-protein comparisons of the whole genome using DIAMOND  
488 was performed<sup>67</sup>, followed by MCScanX<sup>68</sup>, which was used for pairwise synteny  
489 block detection. Parameter settings for MCScanX have been tested and compared  
490 before<sup>18</sup>; here we adopt 'b5s5m25' (b: number of top homologous pairs, s: number of  
491 minimum matched syntenic anchors, m: number of max gene gaps), which has  
492 proven to be appropriate by various studies for the evolutionary distances among  
493 angiosperm genomes. To avoid large numbers of local collinear gene pairs due to  
494 tandem arrays, if consecutive homologs (up to five genes apart) share a common  
495 gene, homologs are collapsed to one representative pair (with the smallest E-value).  
496 Further details regarding phylogenomic synteny network construction can be found  
497 in a tutorial available in the associated GitHub repository  
498 (<https://github.com/zhaotao1987/SynNet-Pipeline>). Each pairwise synteny block  
499 represents pairs of connected nodes (syntenic genes), all pairwise identified synteny  
500 blocks together form a comprehensive synteny network with millions of nodes and  
501 edges. In this synteny network, nodes are genes (from the synteny blocks), while  
502 edges connect syntenic genes. For our work, the entire synteny network summarizes

503 information from 7,435,502 pairwise syntenic blocks, and contains 3,098,333 nodes  
504 (genes) and 94,980,088 edges (syntenic connections).

505 The entire synteny network (database) is clustered for further analysis. We used the  
506 Infomap algorithm for detecting synteny clusters within the map equation  
507 framework<sup>69</sup> (<https://github.com/mapequation/infomap>). We have discussed before  
508 why Infomap is more appropriate for clustering phylogenomic synteny networks<sup>18</sup>.  
509 We used the two-level partitioning mode with ten trials (--clu -N 10 --map -2). The  
510 network was treated as undirected and unweighted. Resulting synteny clusters vary  
511 in size and composition, which is associated with synteny either being well-  
512 conserved or rather lineage-/species-specific. A typical synteny cluster comprises of  
513 syntenic genes shared by groups of species, which precisely represent phylogenetic  
514 relatedness of genomic architecture among species (Fig. 1). Here, we classified the  
515 entire synteny network into 137,833 synteny clusters.

516 A cluster phylogenomic profile shows its composition by the number of nodes in  
517 each species. We summarize the total information residing in all synteny clusters as  
518 a data matrix for tree inference. Phylogenomic profiles of all clusters construct a  
519 large data matrix, where rows represent species, and columns as clusters (Fig. 1).  
520 The matrix was then reduced to a binary presence-absence matrix to obtain the final  
521 synteny matrix (Fig. 1). Here, the dimension of our input synteny matrix is 123 ×  
522 137,833.

523 Tree estimation was based on maximum-likelihood as implemented in IQ-TREE  
524 (version 1.7-beta7) (Nguyen et al., 2014), using the MK+R+FO model. (where “M”  
525 stands for “Markov” and “k” refers to the number of states observed, in our case, k =  
526 2). The +R (FreeRate) model was used to account for site-heterogeneity, and  
527 typically fits data better than the Gamma model for large datasets<sup>70,71</sup>. State  
528 frequencies were optimized by maximum-likelihood (by using ‘+FO’). We generated  
529 1000 bootstrap replicates for the SH-like approximate likelihood ratio test (SH-aLRT),  
530 and 1000 ultrafast bootstrap (UFBoot) replicates (-alrt 1000 -bb 1000)<sup>72</sup>.

531

### 532 **Sequence alignment-based phylogenetic reconstruction**

533 For sequence-based phylogenetic inference, we employed three commonly used  
534 approaches, namely a supermatrix (also called superalignment or concatenation)

535 approach, a reconciliation approach based on the multispecies coalescent (MSC),  
536 and a supertree approach using matrix representation with likelihood (MRL). For  
537 each of these, we used two sets of whole-genome derived gene markers  
538 independently, namely BUSCO genes (Benchmarking Universal Single-Copy  
539 Orthologs)<sup>73</sup> and CSSC genes (Conserved Single-copy Synteny Clusters).

540 The criterion for characterizing of CSSC genes was: median number of nodes across  
541 species < 2, present in ≥ 90% genomes, and presence within Poaceae, monocots  
542 (except Poaceae), Asterids, Rosales, Brassicaceae, and Fabaceae must ≥ 50%.  
543 BUSCO analysis (v3.0, embryophyta\_odb9, with 1440 profiles) identified a total of  
544 1438 conserved single-copy genes from the 123 angiosperm genome sequences,  
545 compared to 883 identified as CSSC. Multiple sequence alignments were performed  
546 using MAFFT (version 7.187)<sup>74</sup>. Two rounds of alignment trimming and filtering were  
547 conducted by trimAl<sup>75</sup>. First, the alignments were trimmed through heuristic selection  
548 of the automatic method (-automated1). Second, sequences with less than 50%  
549 residues that pass the minimum residue overlap threshold (overlap score: 0.5) were  
550 removed (-resoverlap 0.5 -seqoverlap 50). Alignment concatenation was conducted  
551 by catfasta2phymI (<https://github.com/nylander/catfasta2phymI>). The length of our  
552 concatenated gene sequence alignments of BUSCO and CSSC genes were 591,196  
553 and 341,431 amino acids, respectively.

554 Maximum-likelihood analyses were conducted using IQ-TREE<sup>76</sup>. For sequence-  
555 based tree construction, we used the JTT+R model for protein alignments, both for  
556 the construction of trees based on single alignments, as well as for the concatenated  
557 sequence alignments. For all trees, we performed bootstrap analysis with 1000  
558 bootstrap replicates (SH-aLRT and UFBoot (-alrt 1000 -bb 1000)).

559 ASTRAL-Pro<sup>77</sup> was used for the tree summary approach based on multi-species  
560 coalescence to estimate the species trees from 1438 BUSCO gene trees and 883  
561 CSSC gene trees, respectively. ASTRAL-Pro is the latest update of ASTRAL, which  
562 can now account for multi-copy trees.

563 For the MRL supertree analysis, 1438 BUSCO gene trees and 883 CSSC gene trees  
564 were used as two independent data sets. For each tree of the two sets, branch  
565 support from bootstrap analysis has been included as a standard output of IQ-TREE  
566 (suffixed as \*.phy.splits.nex'). We then encoded all splits (bipartitions) with ≥ 85%

567 UFBoot support of all the trees into the data matrix. Thus, each column (matrix  
568 element) in the matrix represents a well-supported bipartition. Coding is similar to the  
569 "Baum-Ragan" coding method (0, 1, ?)<sup>78</sup>, but without question marks because '?'  
570 was originally designed for missing taxa as trees were multi-sourced. The  
571 dimensions of the matrices are 123 × 139,538 for the BUSCO gene trees, and 123 ×  
572 102,617 for the CSSC gene trees. We used the same binary model (MK+R+FO) in  
573 IQ-TREE, and parameter settings were as the one described earlier in Syn-MRL  
574 supercluster analysis.

575 To assess whether alternative sister-group relationships of certain plant clades could  
576 be statistically rejected given the synteny matrix, we performed approximate  
577 unbiased (AU) tests<sup>34</sup>, as implemented in IQ-TREE<sup>76</sup>, under the 'MK+R+FO' model,  
578 with 10,000 replicates.

579

## 580 **Acknowledgements**

581 Z. L. is supported by a postdoctoral fellowship from the Special Research Fund of  
582 Ghent University (BOFPDO2018001701). A. Z. is supported by a PhD fellowship of  
583 the Research Foundation Flanders (FWO). Yves Van de Peer acknowledges funding  
584 from the European Research Council (ERC) under the European Union's Horizon  
585 2020 research and innovation programme (grant agreement No 833522).

586 **Supplemental documents**

587 **Supplementary Fig. 1** Representative sequence-based methods for reconstructing  
588 species trees used in this study, we used superalignment concatenation/supermatrix,  
589 multispecies coalescent, and MRL-based supertree methods for two sets of whole-  
590 genome-derived markers (BUSCO and CSSC) (see text for details). First, multiple  
591 sequence alignments (MSA) were first build for each of the whole-genome-derived  
592 markers (BUSCO and CSSC) and used as input for the Supermatrix method inferring  
593 species trees based on the concatenation of gene alignments. Second, independent  
594 gene trees can be inferred for each alignment, after which a species tree can be  
595 inferred from the set of obtained gene trees under the multispecies coalescent model,  
596 or by using a supertree method. In the latter case, we used a MRL-based method  
597 (see text for details). For example, Clade 1+2 is well-supported ( $BS \geq 85\%$ ) in Tree  
598 1, then this branching order (phylogenetic grouping) is coded as the first column of  
599 the matrix, similarly, Clade 1+2+3 of Tree 1 is coded as the second column. Leaves  
600 from well-supported nodes of all trees construct a binary branch matrix, which is then  
601 used for phylogenetic analysis by maximum likelihood.

602 **Supplementary Fig. 2** Concatenation-BUSCO tree.

603 **Supplementary Fig. 3** Concatenation-CSSC tree.

604 **Supplementary Fig. 4** MRL-BUSCO tree.

605 **Supplementary Fig. 5** MRL-CSSC tree.

606 **Supplementary Fig. 6** ASTRAL-BUSCO tree.

607 **Supplementary Fig. 7** ASTRAL-CSSC tree.

608 **Supplementary Fig. 8** Comparison of synteny tree and the SA tree. Both trees are  
609 rooted by *Amborella*, and three main clades, i.e. superrosids, superasterids, and  
610 monocots are shaded in light-red, light-purple, and light-green, respectively. Eight  
611 differences are indicated by indexed black dots. Branches are not drawn to scale.  
612 Ultrafast bootstrapping values (see text for details) were marked for nodes with less  
613 than 100% support.

614 **Supplementary Fig. 9** Comparison of phylogenetic relationships between the  
615 phylogeny of Angiosperm Phylogeny Group IV and 1KP.

616 **Supplementary Fig. 10** Comparison of phylogenetic relationships between APG IV  
617 and the synteny tree.

618 **Supplementary Fig. 11** Comparison of phylogenetic relationships between the  
619 phylogeny of 1KP study (angiosperms part) and the synteny tree.

620 **Supplementary Fig. 12** Approximate unbiased (AU) test for alternative topologies,  
621 with resulting  $p$  values indicated under the trees. Alternative (tested) topologies  
622 include (a) magnoliids as sister to eudicots, which is the only scenario tested that  
623 cannot be significantly rejected (see text for details). (b) magnoliids as sister to both  
624 monocots and eudicots, (c) Vitales as early-diverging rosids, (d) Saxifragales as  
625 early-diverging rosids, (e) Vitales and Saxifragales as early-diverging rosids, (f)  
626 Malpighiales as early-diverging Fabids, and (g) Malpighiales as early-diverging  
627 Fabids plus Myrtales as early-diverging Malvids.

628 **Supplementary Fig. 13** Syn-MRL tree based on the synteny matrix without 1107  
629 (Figure 4a) specific signals.

630 **Supplementary Fig. 14** Syn-MRL tree based on the synteny matrix without *Triticum*  
631 *turgidum* (wheat).

632 **Supplementary Fig. 15** Syn-MRL tree based on the synteny matrix without  
633 *Hordeum vulgare* (barley).

634 **Supplementary Fig. 16** Syn-MRL tree based on the synteny matrix without  
635 *Brachypodium distachyon* (Brachypodium).

636 **Supplementary Fig. 17** Syn-MRL tree based on the synteny matrix without *Triticum*  
637 *turgidum* (wheat) and *Hordeum vulgare* (barley).

638 **Supplementary Table S1** List of genomes used in this study.

639 **Supplementary Table S2** Phylogenomic profiling rearranged by *Cinnamomum*  
640 *kanehirae* chromosomes and gene orders.

641

## 642 **Data availability**

643 Datasets used in this study are available at DataVerse  
644 (<https://doi.org/10.7910/DVN/7ZZWIH>). This includes all annotated protein  
645 sequences in FASTA format of each genome, entire synteny network database

646 (edgelist), network clustering result, trimmed alignments and corresponding  
647 phylogenetic trees of BUSCO and CSSC genes, bipartitions with support values for  
648 each tree, and binary data matrices. Related codes and software parameters are  
649 available at Github (<https://github.com/zhaotao1987/Syn-MRL>).

650

651 **Author contributions**

652 Y.V.d.P., M.E.S, and J.X. conceived the idea. T.Z. and Y.V.d.P. designed the study.  
653 T.Z. and S.K. performed the analysis. T.Z., J.X., Y.V.d.P, Z.L., M.E.S., and A.Z.  
654 analyzed data. T.Z. and Y.V.d.P. wrote the manuscript. All authors discussed the  
655 results and commented on the manuscript.

656 **References**

657

- 658 1. Van de Peer, Y. Computational approaches to unveiling ancient genome duplications.  
659 *Nat Rev Genet* **5**, 752-763 (2004).
- 660 2. Bowers, J.E., Chapman, B.A., Rong, J. & Paterson, A.H. Unravelling angiosperm  
661 genome evolution by phylogenetic analysis of chromosomal duplication events.  
662 *Nature* **422**, 433-8 (2003).
- 663 3. Pevzner, P. & Tesler, G. Genome rearrangements in mammalian evolution: lessons  
664 from human and mouse genomes. *Genome Res* **13**, 37-45 (2003).
- 665 4. Dewey, C.N. Positional orthology: putting genomic evolutionary relationships into  
666 context. *Brief Bioinform* **12**, 401-412 (2011).
- 667 5. Koonin, E.V. Orthologs, paralogs, and evolutionary genomics. *Annu Rev Genet* **39**,  
668 309-338 (2005).
- 669 6. Van Bel, M. *et al.* Dissecting plant genomes with the PLAZA comparative genomics  
670 platform. *Plant Physiol* **158**, 590-600 (2012).
- 671 7. Zhao, T. *et al.* Phylogenomic synteny network analysis of MADS-box transcription  
672 factor genes reveals lineage-specific transpositions, ancient tandem duplications,  
673 and deep positional conservation. *Plant Cell* **29**, 1278-1292 (2017).
- 674 8. Sampedro, J., Lee, Y., Carey, R.E., Depamphilis, C.W. & Cosgrove, D.J. Use of  
675 genomic history to improve phylogeny and understanding of births and deaths in a  
676 gene family. *Plant J* **44**, 409-419 (2005).
- 677 9. Drillon, G., Champeimont, R., Oteri, F., Fischer, G. & Carbone, A. Phylogenetic  
678 reconstruction based on synteny block and gene adjacencies. *Mol Biol Evol* (2020).
- 679 10. Luo, H. *et al.* Phylogenetic analysis of genome rearrangements among five  
680 mammalian orders. *Mol Phylogen Evol* **65**, 871-882 (2012).
- 681 11. Feng, B. *et al.* Reconstructing yeasts phylogenies and ancestors from whole genome  
682 data. *Sci Rep* **7**, 1-12 (2017).
- 683 12. Watterson, G.A., Ewens, W.J., Hall, T.E. & Morgan, A. The chromosome inversion  
684 problem. *J Theor Biol* **99**, 1-7 (1982).
- 685 13. Blanchette, M., Bourque, G. & Sankoff, D. Breakpoint phylogenies. *Genome*  
686 *informatics* **8**, 25-34 (1997).
- 687 14. Moret, B.M., Bader, D.A., Wyman, S., Warnow, T. & Yan, M. A new implementation  
688 and detailed study of breakpoint analysis. in *Biocomputing 2001* 583-594 (World  
689 Scientific, 2000).
- 690 15. Bourque, G. & Pevzner, P.A. Genome-scale evolution: Reconstructing gene orders in  
691 the ancestral species. *Genome Res* **12**, 26-36 (2002).

- 692 16. Luo, H., Shi, J., Arndt, W., Tang, J. & Friedman, R. Gene order phylogeny of the  
693 genus *Prochlorococcus*. *PloS one* **3**, e3837-e3837 (2008).
- 694 17. Murat, F., Peer, Y.V.d. & Salse, J. Decoding plant and animal genome plasticity from  
695 differential paleo-evolutionary patterns and processes. *Genome Biol Evol* **4**, 917-928  
696 (2012).
- 697 18. Zhao, T. & Schranz, M.E. Network-based microsynteny analysis identifies major  
698 differences and genomic outliers in mammalian and angiosperm genomes. *Proc Natl  
699 Acad Sci U S A* **116**, 2165-2174 (2019).
- 700 19. Van de Peer, Y., Mizrachi, E. & Marchal, K. The evolutionary significance of  
701 polyploidy. *Nat Rev Genet* **18**, 411 (2017).
- 702 20. Leebens-Mack, J.H. *et al.* One thousand plant transcriptomes and the phylogenomics  
703 of green plants. *Nature* **574**, 679-685 (2019).
- 704 21. Payseur, B.A. & Rieseberg, L.H. A genomic perspective on hybridization and  
705 speciation. *Mol Ecol* **25**, 2337-60 (2016).
- 706 22. Fedoroff, N. Transposons and genome evolution in plants. *Proc Natl Acad Sci USA*  
707 **97**, 7002-7007 (2000).
- 708 23. Freeling, M. Bias in plant gene content following different sorts of duplication: tandem,  
709 whole-genome, segmental, or by transposition. *Annu Rev Plant Biol* **60**, 433-453  
710 (2009).
- 711 24. Bowles, A.M., Bechtold, U. & Paps, J. The origin of land plants is rooted in two bursts  
712 of genomic novelty. *Curr Biol* (2020).
- 713 25. Chase, M.W. *et al.* An update of the Angiosperm Phylogeny Group classification for  
714 the orders and families of flowering plants: APG IV. *Bot J Linn Soc* **181**, 1-20 (2016).
- 715 26. Zeng, L. *et al.* Resolution of deep eudicot phylogeny and their temporal diversification  
716 using nuclear genes from transcriptomic and genomic datasets. *New Phytol* **214**,  
717 1338-1354 (2017).
- 718 27. Moore, M.J., Soltis, P.S., Bell, C.D., Burleigh, J.G. & Soltis, D.E. Phylogenetic  
719 analysis of 83 plastid genes further resolves the early diversification of eudicots. *Proc  
720 Natl Acad Sci USA* **107**, 4623 (2010).
- 721 28. Worberg, A. *et al.* Phylogeny of basal eudicots: Insights from non-coding and rapidly  
722 evolving DNA. *Org Divers Evol* **7**, 55-77 (2007).
- 723 29. Soltis, D.E. *et al.* Angiosperm phylogeny inferred from 18S rDNA, *rbcL*, and *atpB*  
724 sequences. *Bot J Linn Soc* **133**, 381-461 (2000).
- 725 30. Soltis, D.E. & Soltis, P.S. Nuclear genomes of two magnoliids. *Nat Plants* **5**, 6 (2019).
- 726 31. Chaw, S.M. *et al.* Stout camphor tree genome fills gaps in understanding of flowering  
727 plant genome evolution. *Nat Plants* **5**, 63-73 (2019).

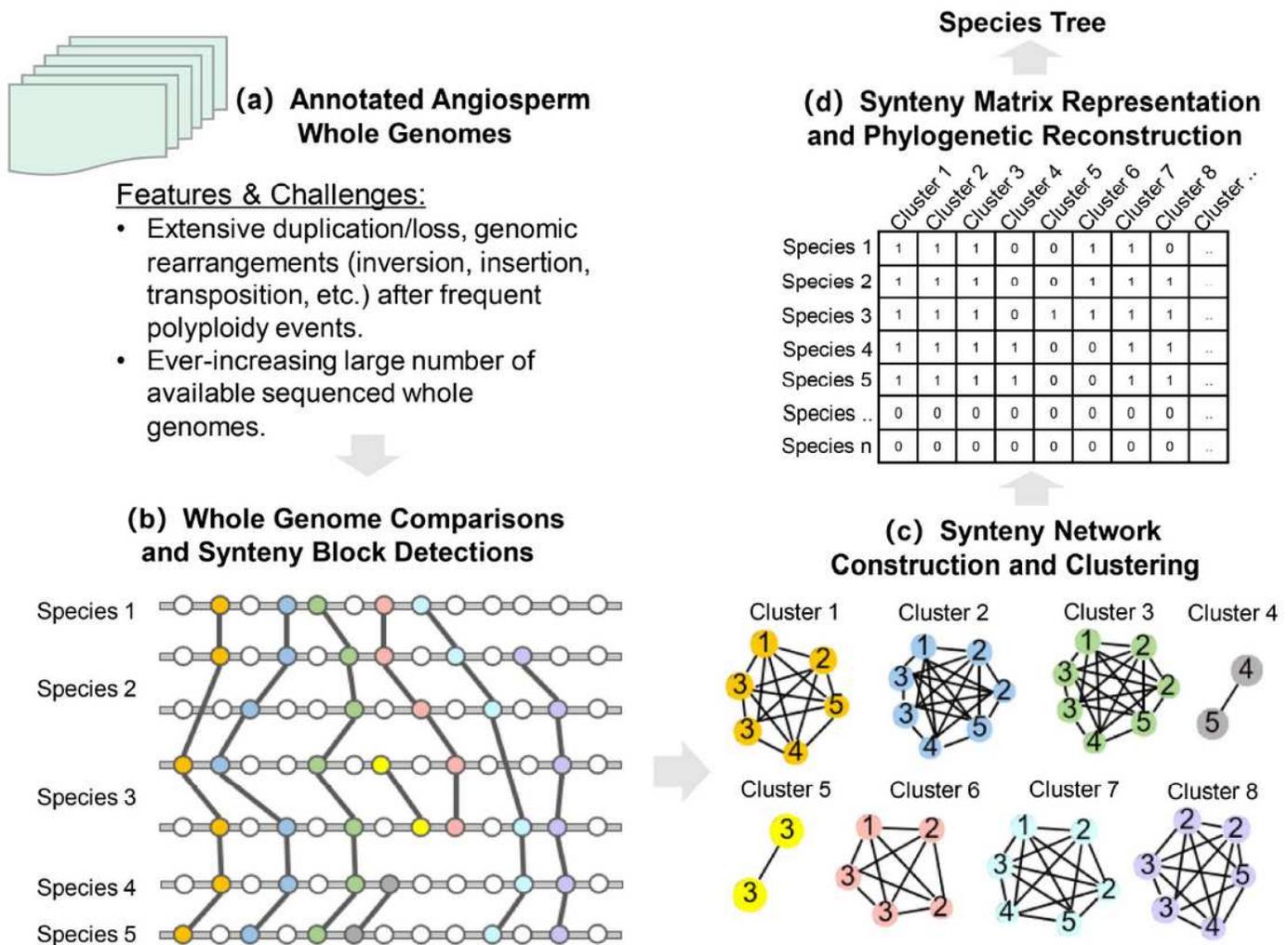
- 728 32. Chen, J. *et al.* *Liriodendron* genome sheds light on angiosperm phylogeny and  
729 species–pair differentiation. *Nat Plants* **5**, 18 (2019).
- 730 33. Rendón-Anaya, M. *et al.* The avocado genome informs deep angiosperm phylogeny,  
731 highlights introgressive hybridization, and reveals pathogen-influenced gene space  
732 adaptation. *Proc Natl Acad Sci USA* **116**, 17081-17089 (2019).
- 733 34. Shimodaira, H. An approximately unbiased test of phylogenetic tree selection. *Syst*  
734 *Biol* **51**, 492-508 (2002).
- 735 35. Boore, J.L. The use of genome-level characters for phylogenetic reconstruction.  
736 *Trends Ecol Evol* **21**, 439-46 (2006).
- 737 36. Rokas, A. & Holland, P.W. Rare genomic changes as a tool for phylogenetics.  
738 *Trends Ecol Evol* **15**, 454-459 (2000).
- 739 37. Nadeau, J.H. & Taylor, B.A. Lengths of chromosomal segments conserved since  
740 divergence of man and mouse. *Proc Natl Acad Sci USA* **81**, 814-818 (1984).
- 741 38. Sankoff, D. *et al.* Gene order comparisons for phylogenetic inference: Evolution of  
742 the mitochondrial genome. *Proc Natl Acad Sci USA* **89**, 6575-6579 (1992).
- 743 39. Sankoff, D. & Nadeau, J.H. Chromosome rearrangements in evolution: From gene  
744 order to genome sequence and back. *Proc Natl Acad Sci USA* **100**, 11188-11189  
745 (2003).
- 746 40. Nguyen, N., Mirarab, S. & Warnow, T. MRL and SuperFine+MRL: new supertree  
747 methods. *Algorithms Mol Biol* **7**, 3 (2012).
- 748 41. Mirarab, S., Bayzid, M.S. & Warnow, T. Evaluating summary methods for multilocus  
749 species tree estimation in the presence of incomplete lineage sorting. *Syst Biol* **65**,  
750 366-80 (2016).
- 751 42. Lewis, P.O. A likelihood approach to estimating phylogeny from discrete  
752 morphological character data. *Syst Biol* **50**, 913-25 (2001).
- 753 43. Soltis, D.E. & Soltis, P.S. Nuclear genomes of two magnoliids. *Nat Plants* **5**, 6-7  
754 (2019).
- 755 44. Hu, L. *et al.* The chromosome-scale reference genome of black pepper provides  
756 insight into piperine biosynthesis. *Nat Commun* **10**, 1-11 (2019).
- 757 45. Soltis, P.S., Soltis, D.E. & Chase, M.W. Angiosperm phylogeny inferred from multiple  
758 genes as a tool for comparative biology. *Nature* **402**, 402-4 (1999).
- 759 46. Sun, G. *et al.* Archaeofractaceae, a new basal angiosperm family. *Science* **296**, 899-  
760 904 (2002).
- 761 47. Endress, P.K. & Doyle, J.A. Reconstructing the ancestral angiosperm flower and its  
762 initial specializations. *Am J Bot* **96**, 22-66 (2009).

- 763 48. Zhang, N., Zeng, L., Shan, H. & Ma, H. Highly conserved low-copy nuclear genes as  
764 effective markers for phylogenetic analyses in angiosperms. *New Phytol* **195**, 923-  
765 937 (2012).
- 766 49. Sun, M. *et al.* Deep phylogenetic incongruence in the angiosperm clade Rosidae. *Mol*  
767 *Phylogen Evol* **83**, 156-166 (2015).
- 768 50. Qiu, Y.L. *et al.* Angiosperm phylogeny inferred from sequences of four mitochondrial  
769 genes. *J Syst Evol* **48**, 391-425 (2010).
- 770 51. Liu, Y., Cox, C.J., Wang, W. & Goffinet, B. Mitochondrial phylogenomics of early land  
771 plants: mitigating the effects of saturation, compositional heterogeneity, and codon-  
772 usage bias. *Syst Biol* **63**, 862-878 (2014).
- 773 52. Finet, C., Timme, R.E., Delwiche, C.F. & Marlétaz, F. Multigene phylogeny of the  
774 green lineage reveals the origin and diversification of land plants. *Curr Biol* **20**, 2217-  
775 2222 (2010).
- 776 53. Wickett, N.J. *et al.* Phylotranscriptomic analysis of the origin and early diversification  
777 of land plants. *Proc Natl Acad Sci U S A* **111**, E4859-68 (2014).
- 778 54. Zeng, L. *et al.* Resolution of deep angiosperm phylogeny using conserved nuclear  
779 genes and estimates of early divergence times. *Nat Commun* **5**, 4956 (2014).
- 780 55. Aköz, G. & Nordborg, M. The *Aquilegia* genome reveals a hybrid origin of core  
781 eudicots. *Genome Biol* **20**, 256 (2019).
- 782 56. Forsythe, E.S., Nelson, A.D. & Beilstein, M.A. Biased gene retention in the face of  
783 massive nuclear introgression obscures species relationships. *bioRxiv*, 197087  
784 (2018).
- 785 57. Walden, N., Nguyen, T.-P., Mandáková, T., Lysak, M.A. & Schranz, M.E. Genomic  
786 blocks in *Aethionema arabicum* support *Arabideae* as next diverging clade in  
787 Brassicaceae. *Front Plant Sci* **11**, 719 (2020).
- 788 58. International Wheat Genome Sequencing, C. *et al.* Shifting the limits in wheat  
789 research and breeding using a fully annotated reference genome. *Science* **361**(2018).
- 790 59. Sage, R.F., Sage, T.L. & Kocacinar, F. Photorespiration and the evolution of C<sub>4</sub>  
791 photosynthesis. *Annu Rev Plant Biol* **63**, 19-47 (2012).
- 792 60. Dasmahapatra, K.K. *et al.* Butterfly genome reveals promiscuous exchange of  
793 mimicry adaptations among species. *Nature* **487**, 94 (2012).
- 794 61. Fontaine, M.C. *et al.* Mosquito genomics. Extensive introgression in a malaria vector  
795 species complex revealed by phylogenomics. *Science* **347**, 1258524 (2015).
- 796 62. Edelman, N.B. *et al.* Genomic architecture and introgression shape a butterfly  
797 radiation. *Science* **366**, 594-599 (2019).

- 798 63. Shan, X. *et al.* Mobilization of the active MITE transposons *mPing* and *Pong* in rice  
799 by introgression from wild rice (*Zizania latifolia* Griseb.). *Mol Biol Evol* **22**, 976-990  
800 (2005).
- 801 64. Wang, N. *et al.* Transpositional reactivation of the *Dart* transposon family in rice lines  
802 derived from introgressive hybridization with *Zizania latifolia*. *BMC Plant Biol* **10**, 190  
803 (2010).
- 804 65. Serrato-Capuchina, A. & Matute, D.R. The role of transposable elements in  
805 speciation. *Genes* **9**, 254 (2018).
- 806 66. Folk, R.A., Soltis, P.S., Soltis, D.E. & Guralnick, R. New prospects in the detection  
807 and comparative analysis of hybridization in the tree of life. *Am J Bot* **105**, 364-375  
808 (2018).
- 809 67. Buchfink, B., Xie, C. & Huson, D.H. Fast and sensitive protein alignment using  
810 DIAMOND. *Nat Methods* **12**, 59-60 (2015).
- 811 68. Wang, Y. *et al.* MCScanX: a toolkit for detection and evolutionary analysis of gene  
812 synteny and collinearity. *Nucleic Acids Res* **40**, e49-e49 (2012).
- 813 69. Rosvall, M. & Bergstrom, C.T. Maps of random walks on complex networks reveal  
814 community structure. *Proc Natl Acad Sci USA* **105**, 1118-1123 (2008).
- 815 70. Yang, Z. A space-time process model for the evolution of DNA sequences. *Genetics*  
816 **139**, 993-1005 (1995).
- 817 71. Soubrier, J. *et al.* The influence of rate heterogeneity among sites on the time  
818 dependence of molecular rates. *Mol Biol Evol* **29**, 3345-3358 (2012).
- 819 72. Hoang, D.T., Chernomor, O., Von Haeseler, A., Minh, B.Q. & Vinh, L.S. UFBoot2:  
820 improving the ultrafast bootstrap approximation. *Mol Biol Evol* **35**, 518-522 (2017).
- 821 73. Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. & Zdobnov, E.M.  
822 BUSCO: assessing genome assembly and annotation completeness with single-copy  
823 orthologs. *Bioinformatics* **31**, 3210-3212 (2015).
- 824 74. Katoh, K. & Standley, D.M. MAFFT multiple sequence alignment software version 7:  
825 improvements in performance and usability. *Mol Biol Evol* **30**, 772-780 (2013).
- 826 75. Capella-Gutiérrez, S., Silla-Martínez, J.M. & Gabaldón, T. trimAl: a tool for automated  
827 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972-  
828 1973 (2009).
- 829 76. Nguyen, L.-T., Schmidt, H.A., von Haeseler, A. & Minh, B.Q. IQ-TREE: a fast and  
830 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol*  
831 *Biol Evol* **32**, 268-274 (2014).
- 832 77. Zhang, C., Scornavacca, C., Molloy, E. & Mirarab, S. ASTRAL-Pro: quartet-based  
833 species tree inference despite paralogy. *bioRxiv* (2019).

834 78. Baum, B.R. Combining trees as a way of combining data sets for phylogenetic  
835 inference, and the desirability of combining gene trees. *Taxon* **41**, 3-10 (1992).  
836  
837

# Figures



**Figure 1**

Whole-genome microsynteny-based species tree inference. (a) Whole genome datasets with all predicted genes are used for phylogeny reconstruction. (b) The synteny network approach first conducts all pairwise reciprocal genome comparisons, followed by synteny block detection. All syntenic blocks constitute the synteny network database (see Methods for details). (c) We analyzed all synteny clusters after clustering the entire network database. Synteny clusters vary in size and node compositions. Shared genomic rearrangements are reflected by cluster compositions. Specific anchor pairs shared by a lineage/species form specific clusters (e.g. Clusters 4-6). We account for the presence or absence of the same recurring anchors for multiple blocks derived from whole-genome or segmental duplications (e.g. for Species 2 and 3 in Clusters 2, 3, 5, and 8). (d) The phylogenomic profiling of all clusters constructs a binary matrix, where rows represent species and columns represent clusters. The synteny matrix comprehensively represents phylogenomic gene order dynamics. It transforms the concept of synteny comparisons from analyzing massive parallel coordinates plots into analyzing profiles of individual clusters/networks. Each cluster stands for a shared homologous 'context'. For example, genes (driven by

transposable elements) can be transposed as insertions into new contexts or can be lost from the original context (e.g. genes in Clusters 4-6). As long as such transpositions are shared by different genomes (e.g. genes in Clusters 4 and 6) or within the same genome because of whole genome duplication (e.g. genes in Cluster 5), specific clusters will emerge and corresponding signals will be added to the matrix. This synteny matrix is used as the input for species tree inference by maximum likelihood (referred to as Syn-MRL).

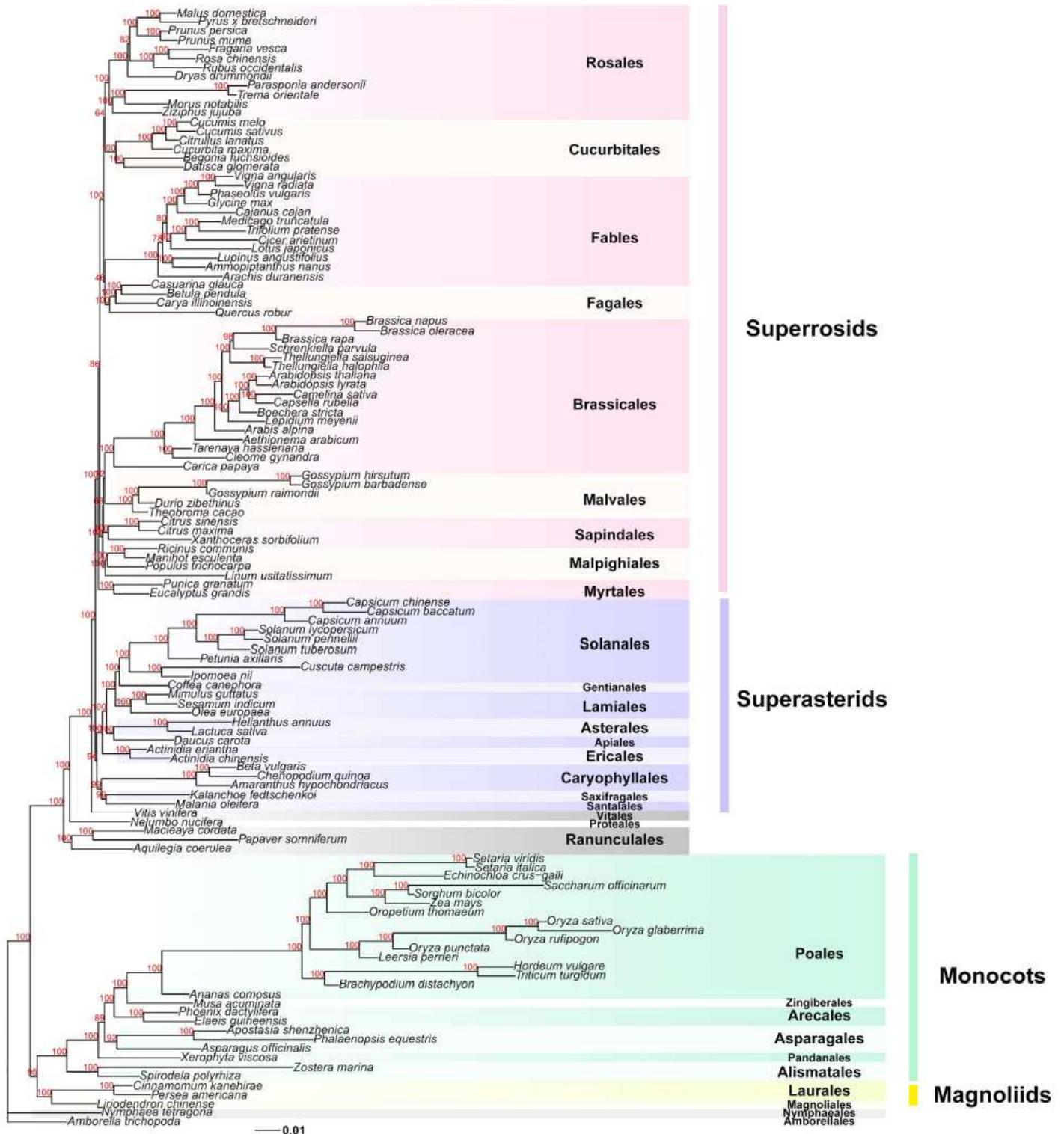


Figure 2

Maximum likelihood (ML) tree for 123 fully-sequenced flowering plant genomes based on the microsynteny approach. The tree is rooted by Amborella, and four main clades, i.e. superrosids, superasterids, monocots, and magnoliids are shaded in light-red, light-purple, light-green, and light-yellow, respectively. Ultrafast bootstrapping values are denoted for all the nodes. Names for the different plant orders follow the APG IV classification<sup>25</sup>.

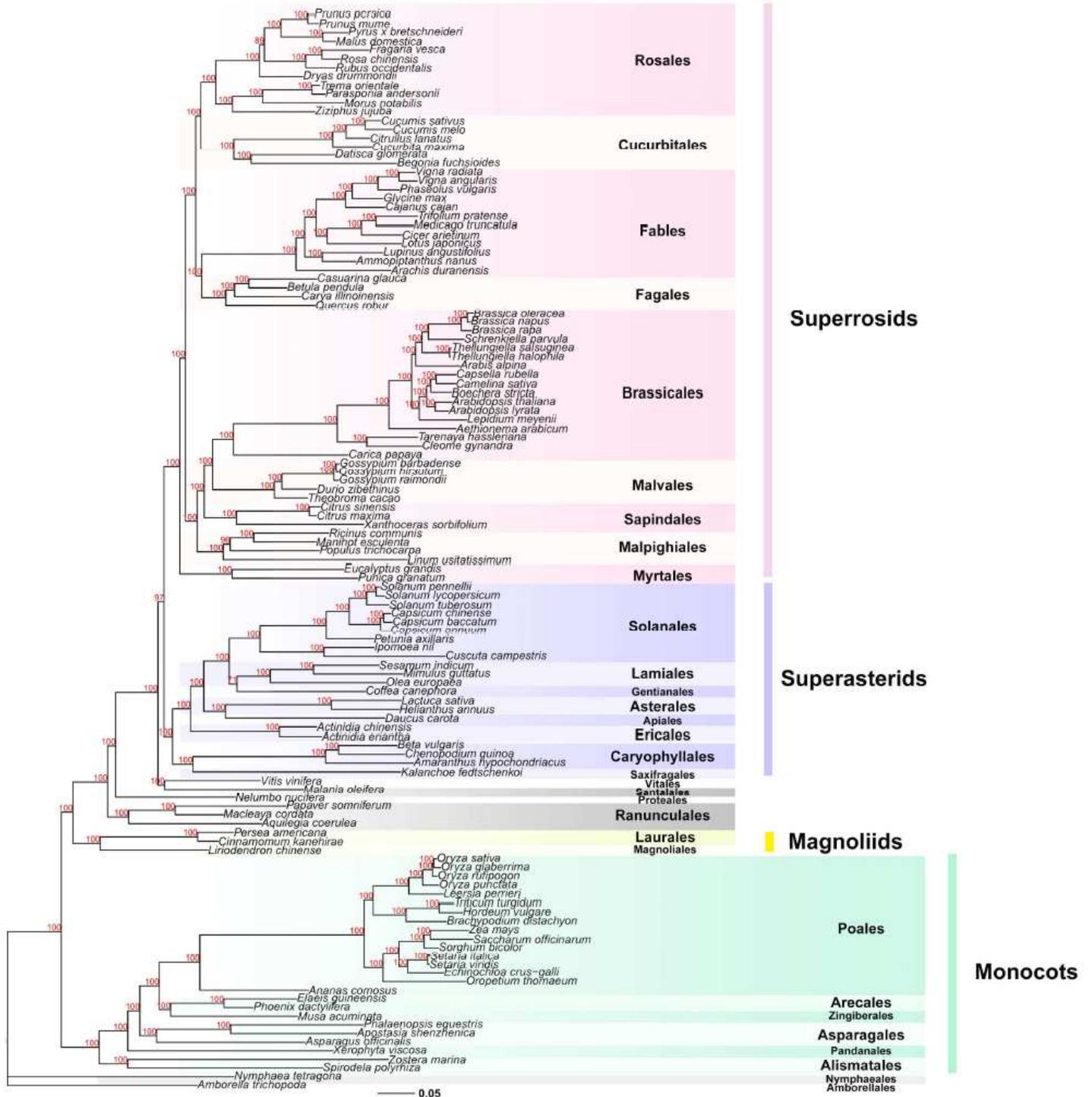
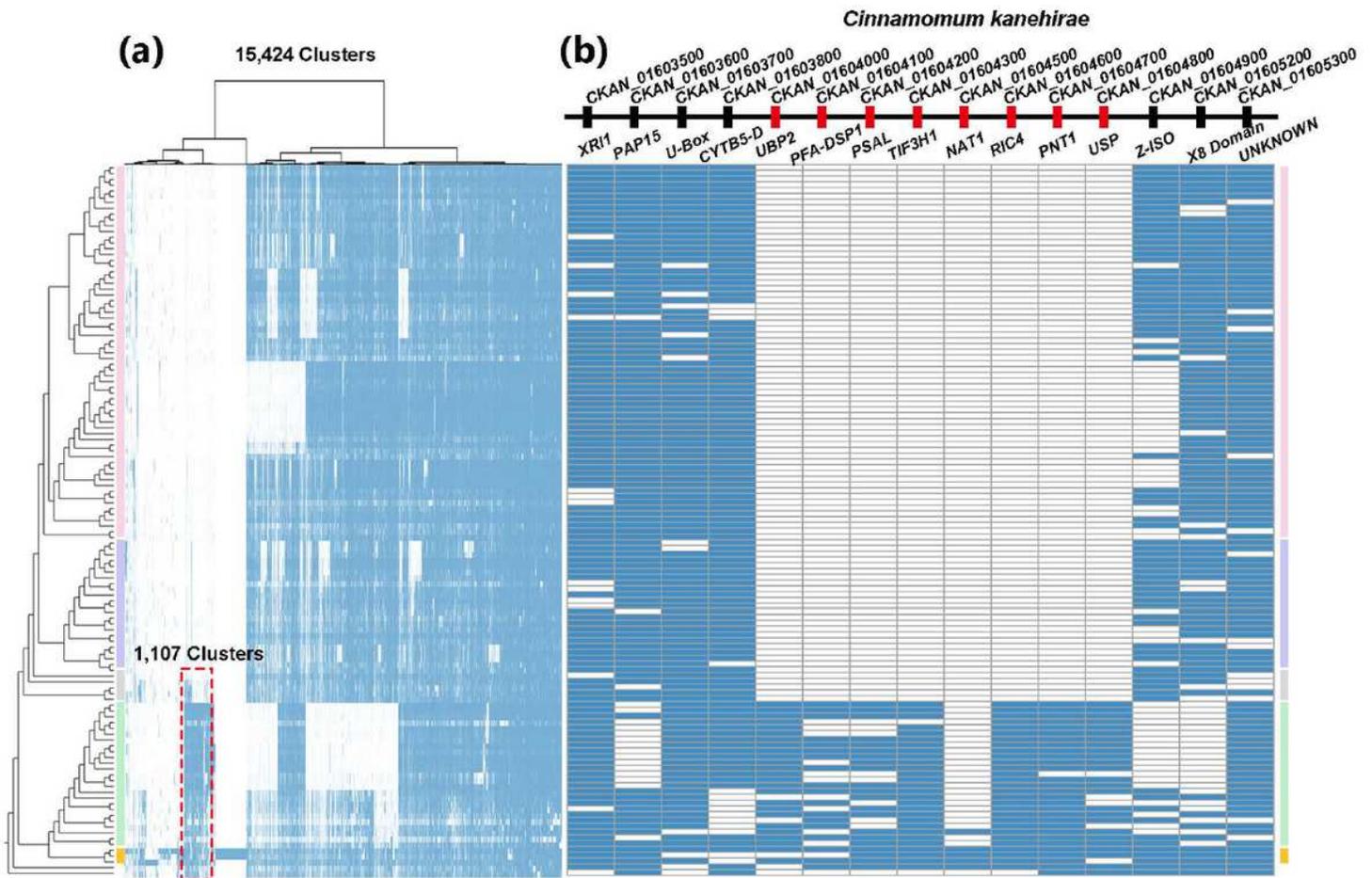


Figure 3

Maximum likelihood (ML) tree based on the concatenation of protein alignments of BUSCO genes. This tree was used as the representative SA (sequence-alignment based) tree.



**Figure 4**

Magnoliids-associated signals and a representative example of phylogenetically informative microsynteny. (a) Hierarchical clustering (ward.D) of 15,424 magnoliids-associate cluster profiles based on jaccard distance. On the far left, the synteny-based species tree is displayed (same as Fig. 1). Superrosids, superasterids, early-diverging eudicots, monocots, and magnoliids are shaded in light-red, light-purple, light-grey, light-green, and light-yellow, respectively. 1,107 clusters supporting a grouping of magnoliids and monocots (Supplemental Table S2). (b) One example from all supporting signals. A fifteen-gene context in the genome of *Cinnamomum kanehirae* (a magnoliid) shows eight neighboring genes (highlighted in orange) only present in magnoliids and monocot genomes, while the flanking genes (colored blue) are generally conserved angiosperm-wide.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryInformation.pdf](#)

- [TableS1.xlsx](#)
- [TableS2.xlsx](#)