1 **Evolution of Auxin Response Factors in plants characterized by phylogenomic**

2 **synteny network analyses**

3

4 Bei Gao[1,2], Liuqiang Wang[3], Melvin Oliver[4], Moxian Chen[5*], Jianhua Zhang[2,6*]

5

6 [1]School of Life Sciences, The Chinese University of Hong Kong, Hong Kong, China.

7 [2]State Key Laboratory of Agrobiotechnology, The Chinese University of Hong Kong, Hong Kong, China.

8 [3]State Key Laboratory of Tree Genetics and Breeding, Research Institute of Forestry, Chinese Academy of Forestry,

9 Beijing 100091, China.

10 [4]USDA-ARS, Plant Genetics Research Unit, University of Missouri, Columbia, MO 65211, USA.

11 [5]Shenzhen Research Institute, The Chinese University of Hong Kong, Shenzhen, China.

12 [6]Department of Biology, Faculty of Science, Hong Kong Baptist University, Hong Kong, China.

13

14 **\*Correspondence:** Jianhua Zhang: jzhang@hkbu.edu.hk and Moxian Chen: zhanglab06@gmail.com

15

16 **SUMMARY**

17 Auxin response factors (ARFs) have long been a research focus and represent a class of key regulators of plant growth

18 and development. Previous studies focusing genes from limited number of species were unable to uncover the

19 evolutionary trajectory of this family. Here, more than 3,500 ARFs collected from plant genomes and transcriptomes

20 covering major streptophyte lineages were used to reconstruct the broad-scale family phylogeny, where the early origin

21 and diversification of *ARF* in charophytes was delineated. Based on the family phylogeny, we proposed a unified six-

22 group classification system for angiosperm ARFs. Phylogenomic synteny network analyses revealed the deeply

23 conserved genomic syntenies within each of the six ARF groups and the interlocking syntenic relationships connecting

24 distinct groups. Recurrent duplication events, such as those that occurred in seed plant, angiosperms, core eudicots and

25 grasses contributed the expansion of ARF genes which facilitated functional diversification. Ancestral transposition

26 activities in important plant families, including crucifers, legumes and grasses, were unveiled by synteny network

27 analyses. Ancestral gene duplications along with transpositions have profound evolutionary significance which may

28 have accelerated the functional diversification process of paralogues. Our study provides insights into the evolution of

29 ARFs which will enhance our current understandings for this important transcription factor family.

30

31 **Key Words:** Auxin, ARF, transcription factor, gene duplication, genomic synteny

32

33

34

## INTRODUCTION

The plant hormone auxin, the chemically simple molecule indole-3-acetic acid, controls many physiological and developmental processes in land plants including but not limited to organogenesis, tissue differentiation, apical dominance, gravitropism as reviewed previously (Kieffer *et al.*, 2010). Completion of the genome of the moss *Physcomitrella patens* makes revealed that many core functional proteins required for auxin biosynthesis, perception, and signaling are present in this early-diverging land plant (embryophyte) lineage (Rensing *et al.*, 2008, Lang *et al.*, 2018), and suggested the auxin molecular regulatory network evolved at the latest in the common ancestor of mosses and angiosperms. A recent review suggested that with the exception of the PGP/ABCB auxin transporters, homologues of all the other core components for hormonal control of physiology and development by auxin could be identified in *P. patens* (Thelander *et al.*, 2018). Changes in auxin perception and signaling that occurred through evolution could have generated the diversification of plant forms that occurred during the past ~474-515 million-year history of the land plants (Morris *et al.*, 2018), which eventually led to the complex vegetative innovations that shape the modern terrestrial and freshwater ecosystems.

Auxin response factors (ARF), as core components in auxin signaling, have long been a focus of plant signaling research (Chapman and Estelle, 2009). The twenty-three *ARFs* identified in the *Arabidopsis thaliana* genome were phylogenetically clustered into three subfamilies (Clades A, B and C) which were subsequently divided into seven groups (ARF9, ARF1, ARF2, ARF3/4, ARF6/8, ARF5/7 and ARF10/16/17);  a classification that was well supported by *ARF* genes from other angiosperms and representative non-flowering lineages (Finet *et al.*, 2013). Generally, ARF proteins can be functionally divided into transcriptional activators (ARF5-8 and 19 in *A. thaliana*) and repressors (remaining ARFs in *A. thaliana*) with well-characterized functional domain architectures (Guilfoyle and Hagen, 2007, Finet *et al.*, 2013). ARFs bind to the auxin response elements (AuxRE: TGTCTC) in the promoter region of downstream auxin-inducible genes (Ulmasov *et al.*, 1997) and function in combination with Aux/IAA repressors, which dimerize with ARF activators in an auxin-regulated manner (Ulmasov *et al.*, 1999, Guilfoyle and Hagen, 2007). Unlike ARF activators, few reports have demonstrated that ARF repressors are able to interact with other ARF proteins or Aux/IAA proteins (Vernoux *et al.*, 2011). A recent top-notch work revealed a newly identified mechanism where the IAA32 and IAA34 transcriptional repressors are stabilized by the transmembrane kinase 1 (TMK1) at the concave side of the apical hook to regulate ARF gene expression and ultimately inhibit growth (Cao *et al.*, 2019).

In most of the well-established transcription factor annotation procedures, such as those implemented by the PlnTFDB(Perez-Rodriguez *et al.*, 2010), PlantTFDB(Jin *et al.*, 2017), iTAK(Zheng *et al.*, 2016) and TAPScan(Wilhelmsson *et al.*, 2017), ARFs were identified using two signature domains: the B3 (PF02362) domain and the Auxin-resp (PF06507) domain,  although some ARF proteins (e.g. ARF23 in *A. thaliana*) may be truncated and lack the Auxin-resp domain(Guilfoyle and Hagen, 2007). Finet *et al.* established a robust and comprehensive phylogenetic framework for the ARF gene families(Finet *et al.*, 2013), however *ARF* genes from non-flowering plants

70    were under-represented. Comprehensive annotation of transcription factors covering distinctive plant clades

71    demonstrated that a number of plant specific transcription factor families (including ARF) evolved in streptophytic

72    algae (charophytes)(Wilhelmsson *et al.*, 2017), suggesting an earlier origin of ARF than that proposed by Finet and

73    collegues(Finet *et al.*, 2013).

74

75    Compared to conventional gene family studies that focus on one or a limited number of  species of interest(Kalluri *et*

76    *al.*, 2007, Wang *et al.*, 2007, Wang *et al.*, 2012b), phylogenetic studies on a broader scale that include multiple plant

77    lineages were able to generate a more robust insights into the evolutionary process that gave rise to the modern

78    assemblage of a target gene family (Finet *et al.*, 2013, Li *et al.*, 2015). The inclusion of genomic synteny data provides

79    important information that impacts the determination of the evolutionary past of a gene family, especially when the

80    gene family of interest evolved in parallel with ancestral genome duplication events (Gao *et al.*, 2018). The

81    conventional genomic block alignment that connects orthologues, retained on genomic syntenic blocks, worked well

82    for a limited number of species(Cheng *et al.*, 2013, Gao *et al.*, 2018), but a network approach was more effective when

83    multiple genomes were included in the synteny analyses(Zhao *et al.*, 2017, Zhao and Schranz, 2017). A comprehensive

84    genomic synteny network can be constructed using nodes to represent the target genes and associated adjacent genomic

85    blocks and the network edges (connecting lines) to represent syntenic relationships(Zhao and Schranz, 2017, Gao *et al.*,

86    2018). The recently established phylogenomic synteny network methodology was able to integrate and summarize

87    genomic synteny relationships to uncover and place genomic events (e.g. ancient tandem duplications, lineage-specific

88    transposition activities) into the evolutionary past of a target gene family(Zhao *et al.*, 2017, Zhao and Schranz, 2017).

89

90    In this study, we collected more than 3,500 ARF members to generate a comprehensive gene-family phylogeny with

91    the aim of filling evolutionary gaps in the non-flowering plants and splitting the long branches present in the current

92    phylogeny(Finet *et al.*, 2013). We propose an updated model for evolution of the ARF family that covered the major

93    streptophytic clades that was based on the six-group classification system we proposed for the ARF genes in

94    angiosperms. A phylogenomic synteny network analyses of angiosperm genomes revealed the deep positional

95    conservation of *ARF* gene-family members within each of the six groups. Detailed individual synteny network analyses

96    together with phylogenetic reconstructions for the six ARF groups revealed their distinctive evolutionary histories.

97    Ancestral duplication events in angiosperms, and subsequent WGDs in eudicots and monocots have contributed to the

98    expansion of ARF members. Ancestral lineage-specific transpositions in important angiosperm families such as

99    crucifers, legumes and grasses were also unveiled. Together, the results presented here add to our current understanding

100   of the evolutionary process that established *ARF* genes in plants. We also expect this broad-scale evolutionary

101   framework could help direct future functional studies that further explore the interplay between auxin signaling and the

102   evolution of land plants.

103

104

## RESULTS AND DISCUSSION

**Auxin Response Factor Evolved in Streptophytic Alage**

To generate a broad-scale phylogenetic profile for ARF genes in plants, we collected a total of 3,502 *ARF* homologues in the streptophytes (including charophytes and embryophytes). *ARFs* were present in all major clades of streptophytes including charophytes, hornworts, liverworts, mosses, lycophytes, ferns and seed plants (Supplementary Fig. S1). In chlorophytes, the Auxin-resp domain was not detected although some chlorophyte genes did contain the B3 domain. Genes containing both the B3 (PF02362) and the Auxin_resp (PF06507) domains were identified in streptophytic algae (charophytes). This was consistent with the observation that a number of plant-specific transcription factors evolved in streptophytic algae (Wilhelmsson *et al.*, 2017). Charophytes represented a paraphyletic clade encompassing successive sister lineages to the land plants (Leliaert *et al.*, 2012, Wilhelmsson *et al.*, 2017). We identified *ARF* homologues in species that are found in three charophyte orders: Zygnematales, Coleochaetales and Chlorokybales, but *ARF* homologues were not identified in the transcriptomes of Charales and Klebsormidiales. However, the presence of ARF in charophytes was affirmed by the *Chara braunii* (Charales) genome (Nishiyama *et al.*, 2018). The identification of the single *ARF* gene in *Chlorokybus atmophyticus* (Chlorokybaceae) suggests that the origin of ARF genes traces back to the root position of streptophytes (Supplementary Fig. S1). This observation suggests an earlier origin of *ARF* gene than those reported previously (Finet *et al.*, 2013, Wilhelmsson *et al.*, 2017).

**Broad-scale Phylogenetic Profile of ARFs in plants**

Overall, the numbers of *ARF* genes in individual angiosperm genomes are greater than those in the individual genomes of non-flowering plants and the 'recent' polyploids, such as *Glycine max,* possess conspicuously more *ARF* genes than other plants (Supplementary Fig. S2). The inclusion of homologues identified from the 1KP transcriptome database provided a comprehensive atlas for the ARF family phylogeny. Overall, the broad-scale phylogeny of *ARF*s generated in this analysis was closely in parallel with the phylogenetic relationships among plant lineages (Fig. 1) derived from large-scale phylotranscriptomic study(Wickett *et al.*, 2014). The phylogenetic tree generated from the ARF gene collection was rooted by the *ARF* gene from *Chlorokybus atmophyticus*, an early diverging charophyte, and exhibited a consistent tree topology with that reported previously(Finet *et al.*, 2013). The incorporation of transcriptomic data from non-flowering plants enabled long evolutionary branches to be split. The phylogenetic analyses also provided robust evidence that angiosperm *ARFs* could be separated clearly into three major subfamilies (Clade A, B and C; consistent with previously reported groupings(Finet *et al.*, 2013)). The three major subfamilies encompassed the six groups (designated as Group I through VI) in this study (Fig. 1). In comparison to the classification proposed by Finet *et al.*(Finet *et al.*, 2013), the Group-I *ARFs* contained the ARF1 and ARF9 subfamilies (which were likely to have derived from an ancient angiosperm-wide duplication) and Group II through VI correspond to the ARF 2, ARF 3/4, ARF 6/8, ARF 5/7 and ARF 10/16/17 subfamilies (Supplementary Table S1), respectively.

139　Groups I, II and III were clustered in the subfamily Clade-B, Groups IV and V in Clade A and Group Vi in Clade C.

140　Clade C was revealed to be basal and a sister clade to subfamilies A plus B. *ARF*s from the charophytes (Zygnematales

141　and Coleochaetales) were separated into clade-C and clade-B failing to partition into a basal mono- or para-phyletic

142　clade, which suggested an ancient diversification of *ARF* genes within the charophytes. The family phylogeny also

143　revealed that each of the six angiosperm ARF family groups was located with gymnosperm *ARF* genes as the closest

144　sister lineage. The tree branches of gymnosperm *ARF* genes are conspicuously shorter than those for angiosperms (Fig.

145　1), which suggested lower amino acid substitutional rates and higher levels of protein sequence conservation in

146　gymnosperm *ARF* genes likely a result of longer generation times that are common in the gymnosperms(Smith and

147　Donoghue, 2008).

148

149　Clade-A contains *ARF* genes that cover all major embryophyte (land plant) clades and contains ARF genes of group-III

150　together with orthologues from gymnosperms and ferns. The *ARF* genes from seed plants and ferns were separated into

151　two major clades which are sister to each other which resulted in a tree topology that was consistent with two child

152　clades derived from an ancient duplication. While lycophyte *ARF* genes were placed outside of and sister to the large

153　duplication clade shared by ferns and seed plants. *ARF* genes from hornworts were identified as basal-most in Clade-A,

154　followed by genes from mosses and liverworts.

155

156　Clade-B was the most diversified lineage containing the angiosperm group I and II genes and along the gymnosperm

157　orthologous genes delineated a conspicuous seed-plant duplication (the ζ event)(Jiao *et al.*, 2011). However, ARF

158　genes from hornworts, liverworts and ferns were mixed into this large duplication clade (Fig. 1). We hypothesize that

159　they might be derived from convergent evolution, though the possibilities of horizonal gene transfer or sequence

160　contaminations cannot be eliminated. Genes from ferns, mosses, liverworts and lycophytes were placed as successive

161　sister lineages to this duplication clade.

162

163　Clade-C was situated as the basal clade with a relatively simple phylogenetic profile and contains genes from every

164　major plant lineage (from charophytes to angiosperms, Fig. 1). This configuration updated the evolutionary model in

165　which clade-C *ARFs* were absent in gymnosperms(Finet *et al.*, 2013).

166

167　The broad-scale phylogenetic analyses in this study established a robust and unified six-group classification system for

168　angiosperm *ARF* genes, which is consistent with previous phylogenetic and domain architecture studies(Guilfoyle and

169　Hagen, 2007, Finet *et al.*, 2013). The relative phylogenetic positions of other land plant lineages were also clarified

170　(Fig. 1), providing a consistent phylogenetic framework for subsequent synteny network analyses.

171

172　**Evolutionary trajectory of ARFs augmented with current genomic and transcriptomic data**

173　In concordance with the phylogenetic analyses described by Finet *et al.*(Finet *et al.*, 2013), we augmented the

174　evolutionary trajectory of ARF family in plants with gene sequences from the currently available genomic and

175    transcriptomic data. The resulting phylogenetic trajectory path (Fig 2.) suggested that the three ARF subfamilies

176    (clades A, B and C) were likely diversified through an ancient duplication in the charophytes, which is consistent to the

177    evolutionary trajectory proposed previously (Flores-Sandoval *et al.*, 2018, Mutte *et al.*, 2018). Tree uncertainties and

178    unresolved land plant phylogenies were also reflected in the ARF gene-family phylogeny, leaving some of the

179    evolutionary processes elusive.

180

181    All of the ARF transcriptional activators (ARF 5-8 and 19 in *A. thaliana*) were clustered in the clade-A subfamily.

182    Within clade-A, the ARF genes were well-conserved in all land plant lineages and appear to have experienced a

183    conspicuous ancient duplication event that occurred in the ancestor of ferns and seed plants. This ancient duplication

184    generated groups IV (ARF6 and 8 in *A. thaliana*) and V (ARF5, 7 and 19 in *A. thaliana*) in the angiosperms and the

185    corresponding sister groups in gymnosperms and ferns (Fig. 2). The *ARF* genes in bryophytes (including hornworts,

186    liverworts and mosses) and lycophytes (clubmosses) were outside of this duplication. The *ARF* genes in clade-A also

187    exhibited a gene tree topology consistent with the 'hornwort-sister' land plant phylogeny in contrast to the 'bryophytes-

188    monophyletic' phylogeny(Wickett *et al.*, 2014, Morris *et al.*, 2018). The evolutionary well-conserved aspect of the

189    ARF activator genes indicates an early genetic foundation for auxin signaling networks in the embryophytes(Thelander

190    *et al.*, 2018).

191

192    In clade-B, unlike clade-A, *ARF* genes from hornworts were not found densely populated at the basal position of the

193    subtree and some hornwort *ARFs* were found clustered with angiosperm *ARF* genes, making the evolution of clade-B

194    ARFs in hornworts elusive. The trajectory analysis suggests two ancient duplications in clade-B, an embryophyte

195    duplication shared by mosses, liverworts and tracheophytes that occurred before the diversification of groups I/II and

196    group III, and a seed plant duplication that generated the groups: I and II. However, the close sister groups for group III

197    were only found in the gymnosperms and ferns which suggested there were gene losses in mosses, liverworts and

198    lycophytes (Fig. 2).

199

200    The subfamily of clade-C is well-conserved, covering all streptophytic lineages, and generated the simplest

201    phylogenetic profile (Fig. 2), containing the group-VI angiosperm *ARF* genes (the ARF 10, 16, and 17 in *A. thaliana*).

202    Hornwort *ARF* genes were placed as direct sisters to the vascular plants (tracheophytes) and the *ARF genes of* mosses

203    and liverworts were placed at the base of the subtree (Fig. 1), generating discrepancies in the gene tree topology and the

204    phylogeny of early land plant lineages.

205

206    **Phylogenomic Synteny Network Analyses of *ARF* genes**

207    The broad-scale phylogenetic analyses suggested some subtree topologies that are consistent with the occurrence of

208    ancient gene duplications but genomic synteny analyses are required to provide more substantive evidence(Tang *et al.*,

209    2008). The recently established synteny network approach, taking advantages of accumulated plant genomes, was able

210    to provide such substantive evidence for ancient evolutionary processes of a specific gene family(Zhao *et al.*, 2017,

211 Zhao and Schranz, 2017). Applying this approach, we conducted a phylogenomic syntenic network analyses for *ARF*

212 genes using a collection of available plant genomes (Supplementary Fig. S1). Syntenic *ARF* genes (syntelogs) were

213 observed in some non-flowering plants (e.g. a lycophyte and a moss), but all represented in-paralogues which were

214 considered to have derived from lineage-specific duplications and e *ARF* genes identified in angiosperm genomes were

215 the primary target of the analysis and used as anchors to construct the genomic synteny network.

216

217 Among the 1,227 annotated angiosperm *ARF* genes containing valid B3 and Auxin_resp domains (Supplementary

218 Table S2), 1,096 (89.3%) were detected to be located within genomic synteny regions that demonstrated genomic

219 collinear relationships with at least another one *ARF* gene, and a total of 18,511 syntenic connections among *ARF*

220 genes were detected (Figs. 3A and 3B). Consistent with the family phylogeny described previously, most of the

221 genome syntenic connections were observed within each of the six groups. *ARF* genes from distinctive ARF groups

222 were syntenically connected (Fig. 3A), for example, *ARF* genes from group VI were connected to *ARF* genes from

223 group III/IV/V and group III *ARF* genes with group I. The ARF synteny network analyses uncovered a total of 82 inter-

224 group connections (Fig. 3A and Supplementary Table S3).

225

226 In the *ARF* gene synteny network, we detected 96 *ARF* genes that did not pass our ARF identification procedure but

227 were demonstrated to be homologous and syntenic to the annotated *ARF* genes. These syntelogs were further inspected

228 and most contained truncated B3 and/or Auxin_resp domains or lacking either or both of these signature domains.

229 These truncated or pseudogenes that were retained in the syntenic genomic blocks were not incorporated in the

230 phylogenetic analyses, however, we were able to assign and label them into one of the six angiosperm *ARF* gene

231 groups by aligning them to classified angiosperm genes. In this way, both intact (total 1,096) and truncated (total 96)

232 *ARF* genes involved in the synteny network were classified. The classification for these truncated genes were

233 considered reliable because of the distant phylogenetic relationships among the six groups (Fig. 1). This may suggest

234 that using genomic syntenic relationships could be a robust approach for detecting pseudogenes retained in the syntenic

235 genomic blocks and which exhibit significant local sequence identity with intact functional paralogues.

236

237 *ARF* genes from each group were conspicuously found in separate and distinct syntenic communities in the initial

238 synteny network visualization (Fig. 3B). The ARF synteny network was further dissected to find subnetwork

239 communities by the use of clique percolation clustering at k = 3 implemented in CFinder v2.0.6 (Adamcsek *et al.*,

240 2006). A total of 25 communities (numbered 0 through 24) (nodes clustered within a subnetwork usually possess more

241 connections in its community than with nodes in other communities) were obtained (Fig. 4). Among the 1,192 ARF

242 syntelogs that were extracted from the synteny network database, 1,128 (94.6%) were identified in the 25 network

243 communities, other syntelogs that had a single syntenic connection or were not involved in a clique (at k=3) were

244 excluded. For example, among the 22 *ARF* genes in *Arabidopsis thaliana*, 17 members were clustered in 13 synteny

245 network communities (Fig. 4A). The chromosome-level genome assemblies represented the best material for genome

246 synteny analyses, but some plant genome assemblies currently available are still highly fragmented. For example, in the

7

247     *Malus domestica* (apple) genome, only one *ARF* gene was clustered in the synteny network because the genome

248     assembly version we obtained from Phytozome database and that was used in our synteny network construction was

249     fragmented (approximately 881.3 Mb arranged in 122,107 scaffolds) (Fig. 4A). However, the network approach using

250     multiple plant genomes appeared to be error-tolerant and the results were unaffected by the inclusion of a few

251     fragmented genomes(Zhao *et al.*, 2017).

252

253     Species compositions for each of the 25 synteny network communities (Fig. 4A) indicate that network communities 4

254     and 5 are angiosperm-wide, containing *ARF* genes from monocots, eudicots and *Amborella*, Community 23, on the

255     other hand, only contains *ARF* genes from monocots and community 24 is solely confined to *ARF* genes from eudicots.

256     Other communities are lineage specific such as community 21 which only contains *ARF* genes from grasses,

257     communities 13 and 14 that are specific to legumes, and communities 16 and 22 that are specific to the Brassica.

258

259     Subnetwork communities were separately visualized, using node colors to depict different plant lineages and node

260     shapes, to delineate *ARF* genes from the different classification groups (Fig. 4B). Community 0 (labeled as 'VI-16-45')

261     consisted of *ARF* members from group-VI, with a total of 16 nodes and 45 connections within the community. Some

262     syntenic communities contained *ARF* genes from multiple groups. Community 5 was recognized as the largest

263     community with 226 nodes and 4,742 connections, and nodes in this community were primarily *ARF* genes from

264     group-VI and group-IV, with a minority of members from group-III (3 nodes) and group-V (1 nodes). The mixed group

265     communities suggest the existence of ancient tandem duplications(Zhao *et al.*, 2017), where duplicated paralogues

266     were likely lost in the ancestral genome such that ancient tandem paralogues are not seen in most current plant

267     genomes, but synteny network analyses reflect them as multigroup communities. Consistent with this hypothesis,

268     tandem *ARF* genes from distant groups were not present in the genome of a single species used in the analysis

269     (Supplementary Table S1). To illustrate this, the *ARF* gene (scaffold00029187) from *Amborella* was classified as a

270     group-IV member, but it had a syntenic connection with group-VI *ARF* genes from *Oryza sativa* (LOC_Os10g33940),

271     *Oropetium thomaeum* (Oropetium_20150105_02810A) and *Phaseolus vulgaris* (Phvul_003G075800). This could be

272     explained by the occurrence of an ancient tandem gene duplication that was generated prior to the separation of groups

273     VI and IV. Following the speciation of basal angiosperms and eudicots plus monocots, the group-VI member was lost

274     in *Amborella*, and the group-IV member was lost in the ancestor of monocots and eudicots resulting in the syntenic

275     relationship seen between group-VI and group-IV *ARF* genes. The inter-group genomic syntenic connections not only

276     provided evidence for ancient gene duplications followed by lineage-specific gene losses, but also suggested that

277     modern *ARF* genes evolved from a common ancestor present in the streptophytes.

278

279     **Evolutionary characteristics for each of the six groups of ARFs in angiosperms**

280     The global phylogenetic and synteny network analyses generated a robust six-group classification system for *ARF*

281     genes and indicated pervasive intra-group syntenic phylogenetic relationships. To elaborate the evolutionary processes

282     within each of the six groups, individual phylogenetic trees for angiosperm genes in each of the six groups were

283 estimated separately and syntenic connections within each network community were mapped onto the six gene trees

284 (Gamboa-Tuz *et al.*, 2018). Along with the *ARF* geness identified from Phytozome plant genomes, the *ARF* genes from

285 basal angiosperms (also ANA grade) and magnoliids were incorporated in the phylogenetic analyses, however these

286 *ARF* gene sequences were derived from transcriptomes and thus did not provide syntenic information. The number of

287 angiosperm (including eudicots, monocots, magnoliids and ANA grade) *ARF* genes in each of the six groups ranged

288 from 190 (group II) to 318 (group I). Below we describe the primary evolutionary characteristics for the six ARF

289 groups separately.

290

291 **Group-I:** This group represented the largest clade (containing 318 angiosperm *ARF* gene members) of the six groups

292 (Fig. 5A). An evident angiosperm-wide duplication (delineated as groups IA and IB) was identified from the tree

293 topology with the three relevant bootstrap values supporting the duplication node and the two child clades greater than

294 95%. Both IA and IB clades include genes from monocots, eudicots, magnoliids and basal angiosperm lineages. The

295 single *ARF* gene member from *Amborella* was placed as sister to the IA plus IB duplication clade, suggesting that the

296 *ARF* gene duplication likely occurred after the separation of *Amborella* from other angiosperms.

297

298 Network communities associated with group-I primarily included angiosperm-wide communities 4 (116 nodes) and 2

299 (65 nodes) (Fig. 4B), which align to groups IA and IB (Fig. 5A), respectively. Group IA was consistent with the

300 designation ARF9 and group IB and the designation ARF1 in *A. thaliana* as reported previousely(Finet *et al.*, 2013).

301 The number of *ARF* genes included in group IA was conspicuously greater than in group IB, particularly for the *ARF*

302 genes from superrosids. The core-eudicot duplication (also known as gamma event)(Jiao *et al.*, 2012) may have

303 contributed to the family expansion, but some ARF genes from magnoliids were also included in the duplication clade

304 and the bootstrap supporting value for the duplication node was also lower than 70%. Moreover, some lineage-specific

305 network communities for *ARF* genes in group-I were observed, where communities 10, 11, 13, 14 and 16 are small

306 communities containing the *ARF* genes from superrosids (Fig. 4B) and these syntenic communities rendered as

307 monophyletic clades in the phylogenetic analyses (Fig. 5A). The species composition analysis (Fig. 4A) for these

308 lineage-specific communities indicated ancestral transposition activities in the Brassicaceae (communities 10, 11 and

309 16) and legumes (communities 13 and 14).

310

311 **Group-II:** Group-II was the smallest group (containing 190 angiosperm *ARF* genes) and synteny network analyses

312 revealed two primary communities, 19 and 8, as depicted in Fig. 5B. Community 8 contains 26 nodes with *ARF* genes

313 from only eudicots and *Amborella*, clustered with a group of magnoliid genes, that formed a paraphyletic clade at the

314 basal position. While the nodes in community 19 were angiosperm-wide, and *ARF* genes from grasses were separated

315 into two clades, one clade following the *ARF* genes in community 8 and the other clade clustered with the other

316 monocots. However, the *ARF* genes clustered in each of the two grass clades did not share syntenic connections (Fig.

317 5B), and the two basal species (*Aquilegia* and *Amborella*) were included in both communities (Fig. 4B), suggesting the

318 genome context (e.g. regulatory elements and adjacent genes) were altered for the *ARF* genes in the two communities.

9

319    The nodes clustered in community 19 may correspond to an ancient tandem duplication in the ancestor of angiosperms

320    as a clade of *ARF* genes from the grasses were evidently separated from other nodes in community 19, indicative of

321    more intra than inter connections (Fig. 4B).

322

323    **Group-III:** Group-III contains 216 *ARF* genes incorporated primarily in network communities 17 and 24 in the

324    synteny network analyses (Fig. 4B and Fig. 5C). The phylogenetic profile for group-III genes identified them as

325    forming two well-separated monophyletic clades (delineated as IIIA and IIIB in this study). The group-IIIA

326    (community 24) contains *ARF* genes from only eudicots and magnoliids, while community 17 is angiosperm-wide and

327    recognized as group-IIIB. The species composition analysis of group-IIIA encompassed a core-eudicot duplication

328    (gamma event), although a magnoliids *ARF* gene was also in this group, that was shared by superrosids and

329    superasterids. *ARF* genes from basal eudicots are recognized as sister to this duplication clade. Similarly, a duplication

330    clade shared by *ARF* genes from grasses (and one gene from pineapple) was conspicuous and likely contributed to the

331    generation of more *ARF* gene members in group-IIIB in the grasses.

332

333    **Group-IV:** Group-IV contains 282 angiosperm *ARF* genes that were contained in six major network communities 5, 9,

334    18, 15, 3 and 6; community 5 was the largest community containing genes from multiple groups (Fig. 4B). Network

335    communities 5 and 9 are angiosperm-wide and 18 contains *ARF* genes from only eudicots. The remaining three

336    communities (3, 6 and 15) were smaller and none formed a high-confidence monophyletic clade which in turn does not

337    support the possibility of ancestral lineage specific transpositions. By comparing the genomic synteny connections with

338    the phylogenetic profile, two evident clusters of *ARF* genes within this group were recognized (Fig. 5D). Community 5

339    (also communities 6 and 15) was clustered into one group and communities 9 and 18 were clustered into another. Both

340    groups were recognized as angiosperm-wide groups suggesting an angiosperm-wide duplication within group-IV,

341    although the duplication topology cannot be easily deduced from the gene tree. In the community 9 network, a

342    subnetwork of monocot *ARF* genes contained more intra-connections than other communities and were separated from

343    other nodes (Fig. 4B), suggesting the possibility of extra rounds of gene duplications and losses in the evolutionary past

344    of *ARF* genes in group IV in angiosperms.

345

346    **Group-V:** Group-V contains a total of 287 angiosperm *ARF* genes that were clustered in four synteny communities, 12,

347    20, 23 and 1 (Fig. 5E), among which communities 12 and 20 were angiosperm-wide, and communities 23 and 1 contain

348    small numbers of monocots *ARF* genes. By integrating the synteny network and phylogenetic profile analyses, three

349    subgroups could be identified (delineated as VA, VB and VC), and consistent with the community network analyses,

350    the nodes in community 12 were phylogenetically separated into two subgroups (groups-VA and -VC). Nodes in

351    communities 23 and 1 were recognized in one monophyletic clade (group-VB). An ancestral transposition in the

352    ancestor of commelinids (including grasses, pineapple and banana genes, community 23 and community 1) (Figs. 4A

353    and 4B), and an *ARF* gene from *Spirodela polyrhiza* (duckweed, sister to commelinids) was syntenically clustered in

354    community 20.

10

355

356 **Group-VI:** The Group-VI included 295 *ARF* genes integrated into 5 network communities 5, 7, 0, 21 and 22 (Fig. 5F).

357 Community 5 was conspicuously angiosperm-wide, community 7 encompassed primarily *ARF* genes from eudicot and

358 *Amborella*, community 0 contained *ARF* genes from monocots and *Amborella*, and communities 21 and 22 were solely

359 comprised of *ARF* genes from grasses and crucifers, respectively (Fig. 4A). Mapping the syntenic connections on the

360 phylogenetic tree, conspicuous monophyletic clades in grasses (community 21) and crucifers (community 22) were

361 generated and provided phylogenomic evidence for ancestral transposition activities in these two lineages. The *ARF*

362 genes clustered in community 5 were phylogenetically separated into two distinct clades with some *ARF* genes from

363 grasses were placed in a basal position in the group-VI phylogeny. The nodes in community 7 were well-clustered in

364 the family phylogeny.

365

366 In the phylogenomic synteny network analyses we employed the maximum-likelihood gene tree generated by IQ-TREE

367 in which more evolutionary models were implemented. We attempted to reconstruct the *ARF* gene family phylogeny

368 using RAxML and the PROTGAMMAAUTO model (Supplementary Fig. S3), which generated alternative tree

369 topologies, nevertheless, the syntenic community patterns remained steady and the major duplication clades and

370 transposition activities could be consistently captured. Tree uncertainties may make some of the evolutionary processes

371 of *ARF* gene family elusive, but the synteny network approach appears robust and uncovered evolutionary details and

372 provided more clues for future experimental studies. For example, *ARF* genes were recurrently duplicated and

373 transposed in specific lineages which suggests that the functions of these transposed genes might reveal novel

374 regulatory elements that were captured in their altered genomic context. The transpositions that we indicate to have

375 occurred in crucifers, legumes, commelinids and grasses were tightly associated with ancestral polyploidy events(Van

376 de Peer *et al.*, 2017), which generated more possibilities in the gene regulatory network. The ancestral gene duplication

377 together with transpositions could have greatly contributed to the expansion of the auxin regulatory network which

378 would have had important implications in the understanding of the evolutionary processes of current land plants.

379

380 **CONCLUSION**

381 In this study, we generated a broad-scale family phylogeny for *ARF* genes from augmented genome and transcriptomic

382 data, that updated our current understanding of the evolutionary history of this transcription factor in streptophytes.

383 Based on the family phylogeny, we proposed a six-group classification regime for angiosperm *ARF* genes. Group IV

384 contains the ARF activators and these genes are well-conserved in all land plant clades. The Group IV subfamily

385 phylogeny also supported the 'hornwort-sister' hypothesis. Genomic synteny network analyses revealed highly-

386 conserved genomic syntenies among angiosperm *ARF* gene loci and within each of the six *ARF* gene groups. CFinder

387 clique analyses of the *ARF* gene synteny network identified 25 subnetwork communities, which were further projected

388 onto the six subfamily phylogenies. The analyses suggest that ancient duplications and transpositions have greatly

389 contributed to the diversification of *ARF* genes in angiosperms. Ancestral lineage-specific transpositions involving *ARF*

390     genes were unveiled in crucifers, legumes, commelinids and grasses in groups I, V and VI. Future studies focusing on

391     non-angiosperm specific lineages should benefit from the evolutionary framework used in this study, especially when

392     more genomes in these plant lineages become available(Cheng *et al.*, 2018).

393

## MATERIALS AND METHODS

### Collection of Auxin Response Factors

396     To generate a broad-scale family phylogeny homologues of plant *ARF* transcription factor genes were obtained from

397     Phytozome v12.1.6 (https://phytozome.jgi.doe.gov/pz/portal.html) and the OneKP (https://db.cngb.org/onekp/)(Matasci

398     *et al.*, 2014) databases using blastp searches filtered with an e-value threshold of 1e-5. *ARF* gene sequences from fern

399     genomes were collected from FernBase (https://www.fernbase.org)(Li *et al.*, 2018). The protein domain composition of

400     each of the putative ARF protein sequences were determined by querying the NCBI Conserved Domain

401     Database(Marchler-Bauer *et al.*, 2017) and only sequences that contained both definitive functional domains: B3 DNA-

402     binding domain (Pfam accession: PF02362) and Auxin_resp (Pfam accession: PF06507), were included in subsequent

403     analyses (Supplementary Table S2).

404

### Family Phylogeny Construction

406     To generate reliable sequence alignments for the collected *ARF* gene-family members, boundaries of the B3 and

407     Auxin_resp domains were identified by aligning each of the protein sequences onto the two HMM profiles using

408     hmmalign v3.2.1(Eddy, 2008, Eddy, 2011). Alignments of the two domains were separately refined using muscle

409     v3.8.1551 and concatenated to generate a broad-scale sequence alignment for *ARF* genes. Columns in the alignment

410     with more than 20% gaps were removed using Phyutility v2.2.6(Smith and Dunn, 2008).

411

412     IQ-TREE v1.6.8(Nguyen *et al.*, 2015) software was employed to reconstruct the maximum likelihood (ML) gene tree.

413     For the obtained broad-scale amino acid alignment, the JTT+R9 model was the best-fit evolutionary model selected by

414     ModelFinder(Kalyaanamoorthy *et al.*, 2017) under Bayesian Information Criterion. The SH-aLRT test and ultrafast

415     bootstrap (Hoang *et al.*, 2018) analyses with 1000 replicates were conducted in IQ-TREE to obtain the supporting

416     values for each internal node of the tree. The obtained maximum-likelihood gene trees were visualized and edited using

417     FigureTree v1.4.4 (http://tree.bio.ed.ac.uk/software/figtree/) and iTOL v4.3 (https://itol.embl.de)(Letunic and Bork,

418     2016). Maximum-likelihood trees for each of the six angiosperm *ARF* clades using IQ-TREE (including model-

419     selection procedure) were also reconstructed to infer potential duplication nodes by analyzing the detailed clade-

420     specific phylogenies.

421

422     The phylogenetic analyses for each of the six ARF groups were performed using both IQ-TREE v1.6.8 (Nguyen *et al.*,

423     2015) and RAxML v8.2.12 (Stamatakis, 2014). The model selection procedure was performed within IQ-TREE based

424     on the Bayesian information criterion (BIC) and for RAxML analyzes we used the '-m PROTGAMMAAUTO' model

425 with 500 bootstrap replicates. All trees were inspected, but the IQ-TREE algorithm produced better bootstrap support

426 overall for branches (Fig. 5 and Supplementary Fig. S3). Each of the six ARF groups contained multiple synteny

427 network communities and syntenic connections in different communities were plotted using different colors as

428 implemented in the iTOL v4.3(Letunic and Bork, 2016).

429

430 **Genomic synteny network construction**

431 To unveil the genomic syntenic relationships among plants, protein sequences for each of the 52 angiosperm genomes

432 were compared with each other and themselves using Diamond v0.9.22.123 software(Buchfink *et al.*, 2015) with an e-

433 value cutoff at 1e-5. In this way, blastp tables for a total of $52\times51/2+52=2,704$ whole proteome comparisons were

434 generated. Only the top five non-self blastp hits were retained as input for the MCScanX(Wang *et al.*, 2012a) analyses.

435 The *ARF* gene associated syntenic genomic block were extracted (Supplementary Table S3) and visualized in

436 Cytoscape v3.7.0(Shannon *et al.*, 2003) and Gephi v0.9.2(Bastian *et al.*, 2009). Some ARF syntelogs were truncated or

437 demonstrate absence of signature domains and were not included in our phylogenetic analyses. These truncated ARF

438 genes were classified and labelled (clade I through VI) by comparing with those classified as *ARF* genes. The

439 phylogeny of angiosperm species and the associated paleopolyploidy events were redrawn based on a tree reported

440 earlier by Van de Peer *et al.*(Van de Peer *et al.*, 2017) and the APG IV system(Byng *et al.*, 2016) with minor

441 modifications: the hexapolyploidy event in cucurbitaceae(Wang *et al.*, 2018), the fern genome duplications(Li *et al.*,

442 2018), the ancestral duplication events mosses(Devos *et al.*, 2016, Lang *et al.*, 2018) and in Caryophyllales(Yang *et al.*,

443 2018), were included in the tree.

444

445 The ARF syntenic networks were analyzed using CFinder v2.0.6(Adamcsek *et al.*, 2006) utilizing the unweighted CPM

446 algorithm and no time limit. All possible k-clique (from 3 to 21) communities were identified for the complete *ARF*

447 gene syntenic network. We used k=3 as the clique community threshold and in this scenario one *ARF* gene (node)

448 involved in a subnetwork community should have at least two connections (edge) with other nodes in the community.

449 Increasing k values made the communities smaller and more disintegrated but also more connected. For illustration

450 purposes, we used different nodal shapes to represent the members from the six ARF groups and different colors to

451 depict specific plant lineages using the Cytoscape v3.7.0 software(Shannon *et al.*, 2003). For each of the 25

452 communities, the species composition of the syntelogs were counted and a heatmap was generated using the pheatmap

453 v1.0.10 (https://github.com/raivokolde/pheatmap) package implemented in the R statistical environment.

454

455 **ACKNOWLEDGEMENTS**

460

## AUTHOR CONTRIBUTIONS

BG conceived the study, performed the bioinformatic analyses and wrote the manuscript. XL and MC contributed to the data collection and discussion. JZ critically revised and improved the manuscript. All authors read and approve the final manuscript.

## CONFLICTS OF INTEREST

The authors declare no conflicts of interest.

## SUPPORTING INFORMATION

Additional Supporting Information may be found in the online version of this article.

**Figure S1.** Plant lineages screened for ARF homologues.

**Figure S2.** Number of Auxin Response Factor genes identified from each of the plant genomes.

**Figure S3.** Phylogenic and synteny network analyses for each of the six groups of ARFs in angiosperms.

**Table S1.** Annotation and Classification of ARF genes in *Arabidopsis thaliana*.

## REFERENCES

**Adamcsek, B., Palla, G., Farkas, I.J., Derenyi, I. and Vicsek, T.** (2006) CFinder: locating cliques and overlapping modules in biological networks. *Bioinformatics*, **22**, 1021-1023.

**Bastian, M., Heymann, S. and Jacomy, M.J.I.** (2009) Gephi: an open source software for exploring and manipulating networks. *International AAAI Conference on Weblogs and Social Media*, **8**, 361-362.

**Buchfink, B., Xie, C. and Huson, D.H.** (2015) Fast and sensitive protein alignment using DIAMOND. *Nature methods*, **12**, 59-60.

**Byng, J.W., Chase, M.W., Christenhusz, M.J.M., Fay, M.F., Judd, W.S., Mabberley, D.J., Sennikov, A.N., Soltis, D.E., Soltis, P.S., Stevens, P.F., Briggs, B., Brockington, S., Chautems, A., Clark, J.C., Conran, J., Haston, E., Moller, M., Moore, M., Olmstead, R., Perret, M., Skog, L., Smith, J., Tank, D., Vorontsova, M., Weber, A. and Grp, A.P.** (2016) An update of the Angiosperm Phylogeny Group classification for the orders and families of flowering plants: APG IV. *Botanical Journal of the Linnean Society*, **181**, 1-20.

**Cao, M., Chen, R., Li, P., Yu, Y., Zheng, R., Ge, D., Zheng, W., Wang, X., Gu, Y., Gelova, Z., Friml, J., Zhang, H., Liu, R., He, J. and Xu, T.** (2019) TMK1-mediated auxin signalling regulates differential growth of the apical hook. *Nature*.

**Chapman, E.J. and Estelle, M.** (2009) Mechanism of auxin-regulated gene expression in plants. *Annu Rev Genet*, **43**, 265-285.

**Cheng, S., Melkonian, M., Smith, S.A., Brockington, S., Archibald, J.M., Delaux, P.M., Li, F.W., Melkonian, B., Mavrodiev, E.V., Sun, W., Fu, Y., Yang, H., Soltis, D.E., Graham, S.W., Soltis, P.S., Liu, X., Xu, X. and Wong, G.K.** (2018) 10KP: A phylodiverse genome sequencing plan. *Gigascience*, **7**, 1-9.

14

499 **Cheng, S.F., van den Bergh, E., Zeng, P., Zhong, X., Xu, J.J., Liu, X., Hofberger, J., de Bruijn, S.,**
500     **Bhide, A.S., Kuelahoglu, C., Bian, C., Chen, J., Fan, G.Y., Kaufmann, K., Hall, J.C.,**
501     **Becker, A., Brautigam, A., Weber, A.P.M., Shi, C.C., Zheng, Z.J., Li, W.J., Lv, M.J., Tao,**
502     **Y.M., Wang, J.Y., Zou, H.F., Quan, Z.W., Hibberd, J.M., Zhang, G.Y., Zhu, X.G., Xu, X.**
503     **and Schranz, M.E.** (2013) The Tarenaya hassleriana Genome Provides Insight into
504     Reproductive Trait and Genome Evolution of Crucifers. *The Plant cell*, **25**, 2813-2830.
505 **Devos, N., Szovenyi, P., Weston, D.J., Rothfels, C.J., Johnson, M.G. and Shaw, A.J.** (2016)
506     Analyses of transcriptome sequences reveal multiple ancient large-scale duplication events in the
507     ancestor of Sphagnopsida (Bryophyta). *The New phytologist*, **211**, 300-318.
508 **Eddy, S.R.** (2008) A probabilistic model of local sequence alignment that simplifies statistical
509     significance estimation. *PLoS Comput Biol*, **4**, e1000069.
510 **Eddy, S.R.** (2011) Accelerated Profile HMM Searches. *PLoS Comput Biol*, **7**, e1002195.
511 **Finet, C., Berne-Dedieu, A., Scutt, C.P. and Marletaz, F.** (2013) Evolution of the ARF gene family in
512     land plants: old domains, new tricks. *Mol Biol Evol*, **30**, 45-56.
513 **Flores-Sandoval, E., Eklund, D.M., Hong, S.-F., Alvarez, J.P., Fisher, T.J., Lampugnani, E.R.,**
514     **Golz, J.F., Vázquez-Lobo, A., Dierschke, T., Lin, S.-S. and Bowman, John L.** (2018) Class C
515     ARFs evolved before the origin of land plants and antagonize differentiation and developmental
516     transitions in Marchantia polymorpha. *New Phytologist*, **218**, 1612-1630.
517 **Gamboa-Tuz, S.D., Pereira-Santana, A., Zhao, T., Schranz, M.E., Castano, E. and Rodriguez-**
518     **Zapata, L.C.** (2018) New insights into the phylogeny of the TMBIM superfamily across the tree
519     of life: Comparative genomics and synteny networks reveal independent evolution of the BI and
520     LFG families in plants. *Mol Phylogenet Evol*, **126**, 266-278.
521 **Gao, B., Chen, M., Li, X., Liang, Y., Zhu, F., Liu, T., Zhang, D., Wood, A.J., Oliver, M.J. and**
522     **Zhang, J.** (2018) Evolution by duplication: paleopolyploidy events in plants reconstructed by
523     deciphering the evolutionary history of VOZ transcription factors. *BMC Plant Biology*, **18**, 256.
524 **Guilfoyle, T.J. and Hagen, G.** (2007) Auxin response factors. *Current opinion in plant biology*, **10**,
525     453-460.
526 **Hoang, D.T., Chernomor, O., von Haeseler, A., Minh, B.Q. and Vinh, L.S.** (2018) UFBoot2:
527     Improving the Ultrafast Bootstrap Approximation. *Mol Biol Evol*, **35**, 518-522.
528 **Jiao, Y., Wickett, N.J., Ayyampalayam, S., Chanderbali, A.S., Landherr, L., Ralph, P.E., Tomsho,**
529     **L.P., Hu, Y., Liang, H., Soltis, P.S., Soltis, D.E., Clifton, S.W., Schlarbaum, S.E., Schuster,**
530     **S.C., Ma, H., Leebens-Mack, J. and dePamphilis, C.W.** (2011) Ancestral polyploidy in seed
531     plants and angiosperms. *Nature*, **473**, 97-100.
532 **Jiao, Y.N., Leebens-Mack, J., Ayyampalayam, S., Bowers, J.E., McKain, M.R., McNeal, J., Rolf,**
533     **M., Ruzicka, D.R., Wafula, E., Wickett, N.J., Wu, X.L., Zhang, Y., Wang, J., Zhang, Y.T.,**
534     **Carpenter, E.J., Deyholos, M.K., Kutchan, T.M., Chanderbali, A.S., Soltis, P.S., Stevenson,**
535     **D.W., McCombie, R., Pires, J.C., Wong, G.K.S., Soltis, D.E. and dePamphilis, C.W.** (2012)
536     A genome triplication associated with early diversification of the core eudicots. *Genome biology*,
537     **13**.
538 **Jin, J., Tian, F., Yang, D.C., Meng, Y.Q., Kong, L., Luo, J. and Gao, G.** (2017) PlantTFDB 4.0:
539     toward a central hub for transcription factors and regulatory interactions in plants. *Nucleic Acids*
540     *Res*, **45**, D1040-D1045.
541 **Kalluri, U.C., DiFazio, S.P., Brunner, A.M. and Tuskan, G.A.** (2007) Genome-wide analysis of
542     Aux/IAA and ARF gene families in Populus trichocarpa. *Bmc Plant Biology*, **7**.

**Kalyaanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A. and Jermiin, L.S.** (2017) ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature methods*, **14**, 587-589.

**Kieffer, M., Neve, J. and Kepinski, S.** (2010) Defining auxin response contexts in plant development. *Current opinion in plant biology*, **13**, 12-20.

**Lang, D., Ullrich, K.K., Murat, F., Fuchs, J., Jenkins, J., Haas, F.B., Piednoel, M., Gundlach, H., Van Bel, M., Meyberg, R., Vives, C., Morata, J., Symeonidi, A., Hiss, M., Muchero, W., Kamisugi, Y., Saleh, O., Blanc, G., Decker, E.L., van Gessel, N., Grimwood, J., Hayes, R.D., Graham, S.W., Gunter, L.E., McDaniel, S.F., Hoernstein, S.N.W., Larsson, A., Li, F.W., Perroud, P.F., Phillips, J., Ranjan, P., Rokshar, D.S., Rothfels, C.J., Schneider, L., Shu, S., Stevenson, D.W., Thummler, F., Tillich, M., Villarreal Aguilar, J.C., Widiez, T., Wong, G.K., Wymore, A., Zhang, Y., Zimmer, A.D., Quatrano, R.S., Mayer, K.F.X., Goodstein, D., Casacuberta, J.M., Vandepoele, K., Reski, R., Cuming, A.C., Tuskan, G.A., Maumus, F., Salse, J., Schmutz, J. and Rensing, S.A.** (2018) The Physcomitrella patens chromosome-scale assembly reveals moss genome structure and evolution. *The Plant journal : for cell and molecular biology*, **93**, 515-533.

**Leliaert, F., Smith, D.R., Moreau, H., Herron, M.D., Verbruggen, H., Delwiche, C.F. and De Clerck, O.** (2012) Phylogeny and Molecular Evolution of the Green Algae. *Critical Reviews in Plant Sciences*, **31**, 1-46.

**Letunic, I. and Bork, P.** (2016) Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res*, **44**, W242-245.

**Li, F.-W., Melkonian, M., Rothfels, C.J., Villarreal, J.C., Stevenson, D.W., Graham, S.W., Wong, G.K.-S., Pryer, K.M. and Mathews, S.** (2015) Phytochrome diversity in green plants and the origin of canonical plant phytochromes. *Nature Communications*, **6**, 7852.

**Li, F.W., Brouwer, P., Carretero-Paulet, L., Cheng, S., de Vries, J., Delaux, P.M., Eily, A., Koppers, N., Kuo, L.Y., Li, Z., Simenc, M., Small, I., Wafula, E., Angarita, S., Barker, M.S., Brautigam, A., dePamphilis, C., Gould, S., Hosmani, P.S., Huang, Y.M., Huettel, B., Kato, Y., Liu, X., Maere, S., McDowell, R., Mueller, L.A., Nierop, K.G.J., Rensing, S.A., Robison, T., Rothfels, C.J., Sigel, E.M., Song, Y., Timilsena, P.R., Van de Peer, Y., Wang, H., Wilhelmsson, P.K.I., Wolf, P.G., Xu, X., Der, J.P., Schluepmann, H., Wong, G.K. and Pryer, K.M.** (2018) Fern genomes elucidate land plant evolution and cyanobacterial symbioses. *Nat Plants*, **4**, 460-472.

**Marchler-Bauer, A., Bo, Y., Han, L., He, J., Lanczycki, C.J., Lu, S., Chitsaz, F., Derbyshire, M.K., Geer, R.C., Gonzales, N.R., Gwadz, M., Hurwitz, D.I., Lu, F., Marchler, G.H., Song, J.S., Thanki, N., Wang, Z., Yamashita, R.A., Zhang, D., Zheng, C., Geer, L.Y. and Bryant, S.H.** (2017) CDD/SPARCLE: functional classification of proteins via subfamily domain architectures. *Nucleic Acids Res*, **45**, D200-D203.

**Matasci, N., Hung, L.H., Yan, Z.X., Carpenter, E.J., Wickett, N.J., Mirarab, S., Nguyen, N., Warnow, T., Ayyampalayam, S., Barker, M., Burleigh, J.G., Gitzendanner, M.A., Wafula, E., Der, J.P., dePamphilis, C.W., Roure, B., Philippe, H., Ruhfel, B.R., Miles, N.W., Graham, S.W., Mathews, S., Surek, B., Melkonian, M., Soltis, D.E., Soltis, P.S., Rothfels, C., Pokorny, L., Shaw, J.A., DeGironimo, L., Stevenson, D.W., Villarreal, J.C., Chen, T., Kutchan, T.M., Rolf, M., Baucom, R.S., Deyholos, M.K., Samudrala, R., Tian, Z.J., Wu, X.L., Sun, X., Zhang, Y., Wang, J., Leebens-Mack, J. and Wong, G.K.S.** (2014) Data access for the 1,000 Plants (1KP) project. *Gigascience*, **3**, 17.
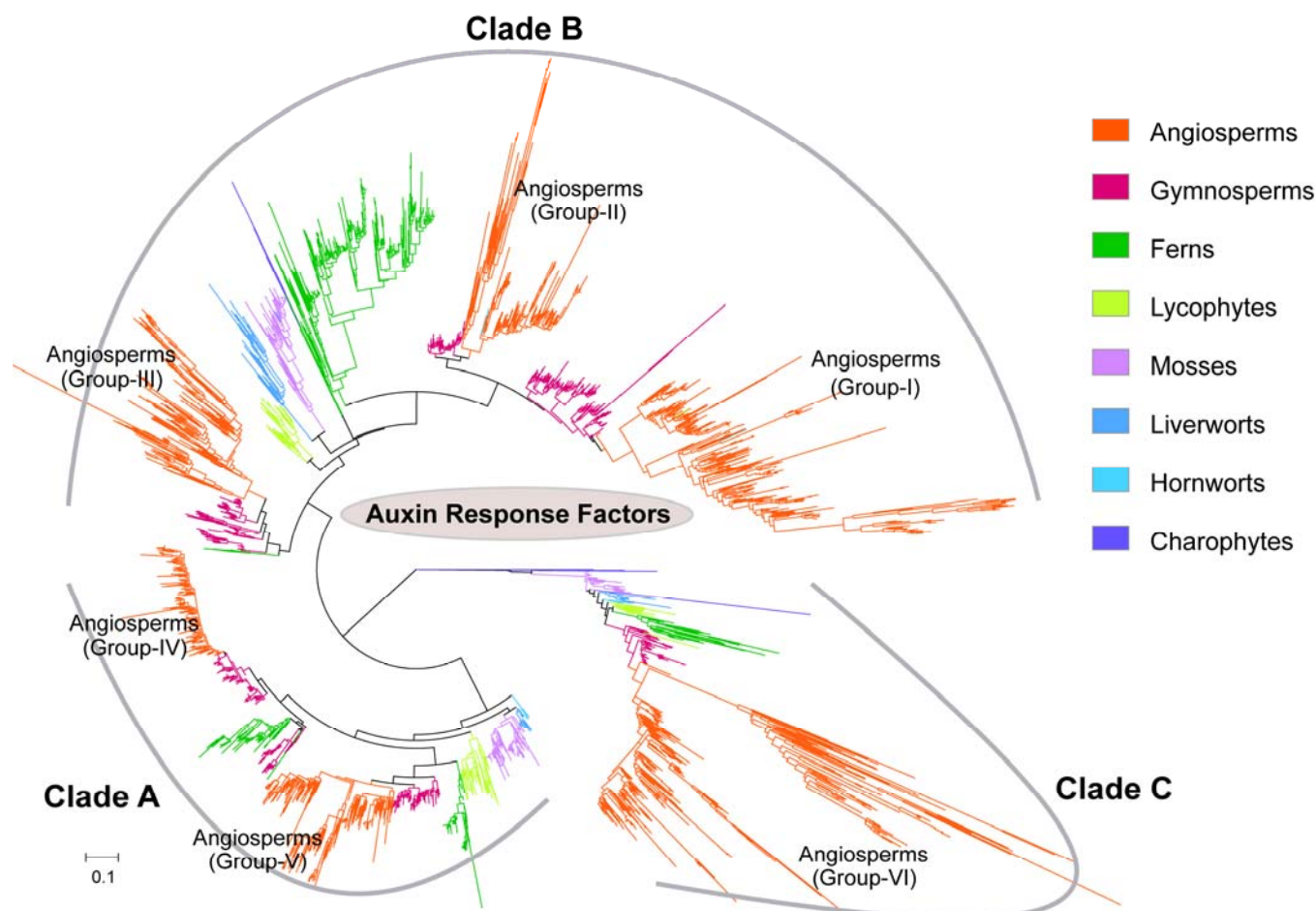
588 **Morris, J.L., Puttick, M.N., Clark, J.W., Edwards, D., Kenrick, P., Pressel, S., Wellman, C.H.,**
589     **Yang, Z., Schneider, H. and Donoghue, P.C.J.** (2018) The timescale of early land plant
590     evolution. *Proceedings of the National Academy of Sciences of the United States of America*,
591     **115**, E2274-E2283.
592 **Mutte, S.K., Kato, H., Rothfels, C., Melkonian, M., Wong, G.K.S. and Weijers, D.** (2018) Origin
593     and evolution of the nuclear auxin response system. *Elife*, **7**.
594 **Nguyen, L.-T., Schmidt, H.A., von Haeseler, A. and Minh, B.Q.** (2015) IQ-TREE: A Fast and
595     Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular*
596     *Biology and Evolution*, **32**, 268-274.
597 **Nishiyama, T., Sakayama, H., de Vries, J., Buschmann, H., Saint-Marcoux, D., Ullrich, K.K.,**
598     **Haas, F.B., Vanderstraeten, L., Becker, D., Lang, D., Vosolsobe, S., Rombauts, S.,**
599     **Wilhelmsson, P.K.I., Janitza, P., Kern, R., Heyl, A., Rumpler, F., Villalobos, L., Clay, J.M.,**
600     **Skokan, R., Toyoda, A., Suzuki, Y., Kagoshima, H., Schijlen, E., Tajeshwar, N., Catarino,**
601     **B., Hetherington, A.J., Saltykova, A., Bonnot, C., Breuninger, H., Symeonidi, A.,**
602     **Radhakrishnan, G.V., Van Nieuwerburgh, F., Deforce, D., Chang, C., Karol, K.G.,**
603     **Hedrich, R., Ulvskov, P., Glockner, G., Delwiche, C.F., Petrasek, J., Van de Peer, Y., Friml,**
604     **J., Beilby, M., Dolan, L., Kohara, Y., Sugano, S., Fujiyama, A., Delaux, P.M., Quint, M.,**
605     **Theissen, G., Hagemann, M., Harholt, J., Dunand, C., Zachgo, S., Langdale, J., Maumus,**
606     **F., Van Der Straeten, D., Gould, S.B. and Rensing, S.A.** (2018) The Chara Genome:
607     Secondary Complexity and Implications for Plant Terrestrialization. *Cell*, **174**, 448-464 e424.
608 **Perez-Rodriguez, P., Riano-Pachon, D.M., Correa, L.G., Rensing, S.A., Kersten, B. and Mueller-**
609     **Roeber, B.** (2010) PlnTFDB: updated content and new features of the plant transcription factor
610     database. *Nucleic Acids Res*, **38**, D822-827.
611 **Rensing, S.A., Lang, D., Zimmer, A.D., Terry, A., Salamov, A., Shapiro, H., Nishiyama, T.,**
612     **Perroud, P.F., Lindquist, E.A., Kamisugi, Y., Tanahashi, T., Sakakibara, K., Fujita, T.,**
613     **Oishi, K., Shin, I.T., Kuroki, Y., Toyoda, A., Suzuki, Y., Hashimoto, S., Yamaguchi, K.,**
614     **Sugano, S., Kohara, Y., Fujiyama, A., Anterola, A., Aoki, S., Ashton, N., Barbazuk, W.B.,**
615     **Barker, E., Bennetzen, J.L., Blankenship, R., Cho, S.H., Dutcher, S.K., Estelle, M., Fawcett,**
616     **J.A., Gundlach, H., Hanada, K., Heyl, A., Hicks, K.A., Hughes, J., Lohr, M., Mayer, K.,**
617     **Melkozernov, A., Murata, T., Nelson, D.R., Pils, B., Prigge, M., Reiss, B., Renner, T.,**
618     **Rombauts, S., Rushton, P.J., Sanderfoot, A., Schween, G., Shiu, S.H., Stueber, K.,**
619     **Theodoulou, F.L., Tu, H., Van de Peer, Y., Verrier, P.J., Waters, E., Wood, A., Yang, L.,**
620     **Cove, D., Cuming, A.C., Hasebe, M., Lucas, S., Mishler, B.D., Reski, R., Grigoriev, I.V.,**
621     **Quatrano, R.S. and Boore, J.L.** (2008) The Physcomitrella genome reveals evolutionary
622     insights into the conquest of land by plants. *Science*, **319**, 64-69.
623 **Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N., Schwikowski,**
624     **B. and Ideker, T.** (2003) Cytoscape: a software environment for integrated models of
625     biomolecular interaction networks. *Genome Res*, **13**, 2498-2504.
626 **Smith, S.A. and Donoghue, M.J.** (2008) Rates of molecular evolution are linked to life history in
627     flowering plants. *Science*, **322**, 86-89.
628 **Smith, S.A. and Dunn, C.W.** (2008) Phyutility: a phyloinformatics tool for trees, alignments and
629     molecular data. *Bioinformatics*, **24**, 715-716.
630 **Stamatakis, A.** (2014) RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
631     phylogenies. *Bioinformatics*, **30**, 1312-1313.
632 **Tang, H., Bowers, J.E., Wang, X., Ming, R., Alam, M. and Paterson, A.H.** (2008) Synteny and
633     collinearity in plant genomes. *Science*, **320**, 486-488.

17

634 **Thelander, M., Landberg, K. and Sundberg, E.** (2018) Auxin-mediated developmental control in the
635   moss Physcomitrella patens. *J Exp Bot*, **69**, 277-290.
636 **Ulmasov, T., Hagen, G. and Guilfoyle, T.J.** (1997) ARF1, a transcription factor that binds to auxin
637   response elements. *Science*, **276**, 1865-1868.
638 **Ulmasov, T., Hagen, G. and Guilfoyle, T.J.** (1999) Dimerization and DNA binding of auxin response
639   factors. *Plant Journal*, **19**, 309-319.
640 **Van de Peer, Y., Mizrachi, E. and Marchal, K.** (2017) The evolutionary significance of polyploidy.
641   *Nature reviews. Genetics*, **18**, 411-424.
642 **Vernoux, T., Brunoud, G., Farcot, E., Morin, V., Van den Daele, H., Legrand, J., Oliva, M., Das,**
643   **P., Larrieu, A., Wells, D., Guedon, Y., Armitage, L., Picard, F., Guyomarc'h, S., Cellier, C.,**
644   **Parry, G., Koumproglou, R., Doonan, J.H., Estelle, M., Godin, C., Kepinski, S., Bennett,**
645   **M., De Veylder, L. and Traas, J.** (2011) The auxin signalling network translates dynamic input
646   into robust patterning at the shoot apex. *Molecular Systems Biology*, **7**.
647 **Wang, D., Pei, K., Fu, Y., Sun, Z., Li, S., Liu, H., Tang, K., Han, B. and Tao, Y.** (2007) Genome-
648   wide analysis of the auxin response factors (ARF) gene family in rice (Oryza sativa). *Gene*, **394**,
649   13-24.
650 **Wang, J., Sun, P., Li, Y., Liu, Y., Yang, N., Yu, J., Ma, X., Sun, S., Xia, R., Liu, X., Ge, D., Luo, S.,**
651   **Liu, Y., Kong, Y., Cui, X., Lei, T., Wang, L., Wang, Z., Ge, W., Zhang, L., Song, X., Yuan,**
652   **M., Guo, D., Jin, D., Chen, W., Pan, Y., Liu, T., Yang, G., Xiao, Y., Sun, J., Zhang, C., Li,**
653   **Z., Xu, H., Duan, X., Shen, S., Zhang, Z., Huang, S. and Wang, X.** (2018) An Overlooked
654   Paleotetraploidization in Cucurbitaceae. *Mol Biol Evol*, **35**, 16-26.
655 **Wang, Y., Tang, H., Debarry, J.D., Tan, X., Li, J., Wang, X., Lee, T.H., Jin, H., Marler, B., Guo,**
656   **H., Kissinger, J.C. and Paterson, A.H.** (2012a) MCScanX: a toolkit for detection and
657   evolutionary analysis of gene synteny and collinearity. *Nucleic Acids Res*, **40**, e49.
658 **Wang, Y.J., Deng, D.X., Shi, Y.T., Miao, N., Bian, Y.L. and Yin, Z.T.** (2012b) Diversification,
659   phylogeny and evolution of auxin response factor (ARF) family: insights gained from analyzing
660   maize ARF genes. *Molecular biology reports*, **39**, 2401-2415.
661 **Wickett, N.J., Mirarab, S., Nguyen, N., Warnow, T., Carpenter, E., Matasci, N., Ayyampalayam,**
662   **S., Barker, M.S., Burleigh, J.G., Gitzendanner, M.A., Ruhfel, B.R., Wafula, E., Der, J.P.,**
663   **Graham, S.W., Mathews, S., Melkonian, M., Soltis, D.E., Soltis, P.S., Miles, N.W., Rothfels,**
664   **C.J., Pokorny, L., Shaw, A.J., DeGironimo, L., Stevenson, D.W., Surek, B., Villarreal, J.C.,**
665   **Roure, B., Philippe, H., dePamphilis, C.W., Chen, T., Deyholos, M.K., Baucom, R.S.,**
666   **Kutchan, T.M., Augustin, M.M., Wang, J., Zhang, Y., Tian, Z., Yan, Z., Wu, X., Sun, X.,**
667   **Wong, G.K. and Leebens-Mack, J.** (2014) Phylotranscriptomic analysis of the origin and early
668   diversification of land plants. *Proceedings of the National Academy of Sciences of the United*
669   *States of America*, **111**, E4859-4868.
670 **Wilhelmsson, P.K.I., Muhlich, C., Ullrich, K.K. and Rensing, S.A.** (2017) Comprehensive Genome-
671   Wide Classification Reveals That Many Plant-Specific Transcription Factors Evolved in
672   Streptophyte Algae. *Genome Biology and Evolution*, **9**, 3384-3397.
673 **Yang, Y., Moore, M.J., Brockington, S.F., Mikenas, J., Olivieri, J., Walker, J.F. and Smith, S.A.**
674   (2018) Improved transcriptome sampling pinpoints 26 ancient and more recent polyploidy events
675   in Caryophyllales, including two allopolyploidy events. *The New phytologist*, **217**, 855-870.
676 **Zhao, T., Holmer, R., de Bruijn, S., Angenent, G.C., van den Burg, H.A. and Schranz, M.E.** (2017)
677   Phylogenomic Synteny Network Analysis of MADS-Box Transcription Factor Genes Reveals
678   Lineage-Specific Transpositions, Ancient Tandem Duplications, and Deep Positional
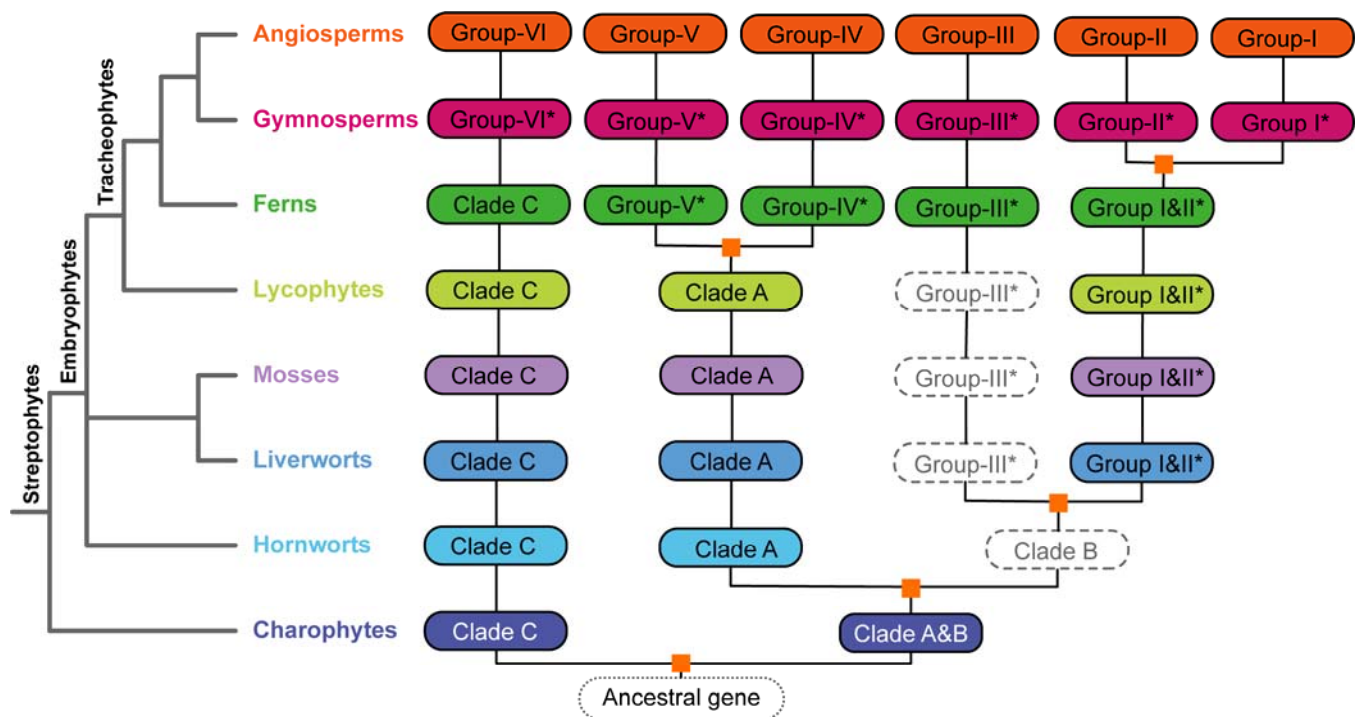679   Conservation. *The Plant cell*, **29**, 1278-1292.

680    **Zhao, T. and Schranz, M.E.** (2017) Network approaches for plant phylogenomic synteny analysis.
681        *Current opinion in plant biology*, **36**, 129-134.
682    **Zheng, Y., Jiao, C., Sun, H., Rosli, H.G., Pombo, M.A., Zhang, P., Banf, M., Dai, X., Martin, G.B.,**
683        **Giovannoni, J.J., Zhao, P.X., Rhee, S.Y. and Fei, Z.** (2016) iTAK: A Program for Genome-
684        wide Prediction and Classification of Plant Transcription Factors, Transcriptional Regulators,
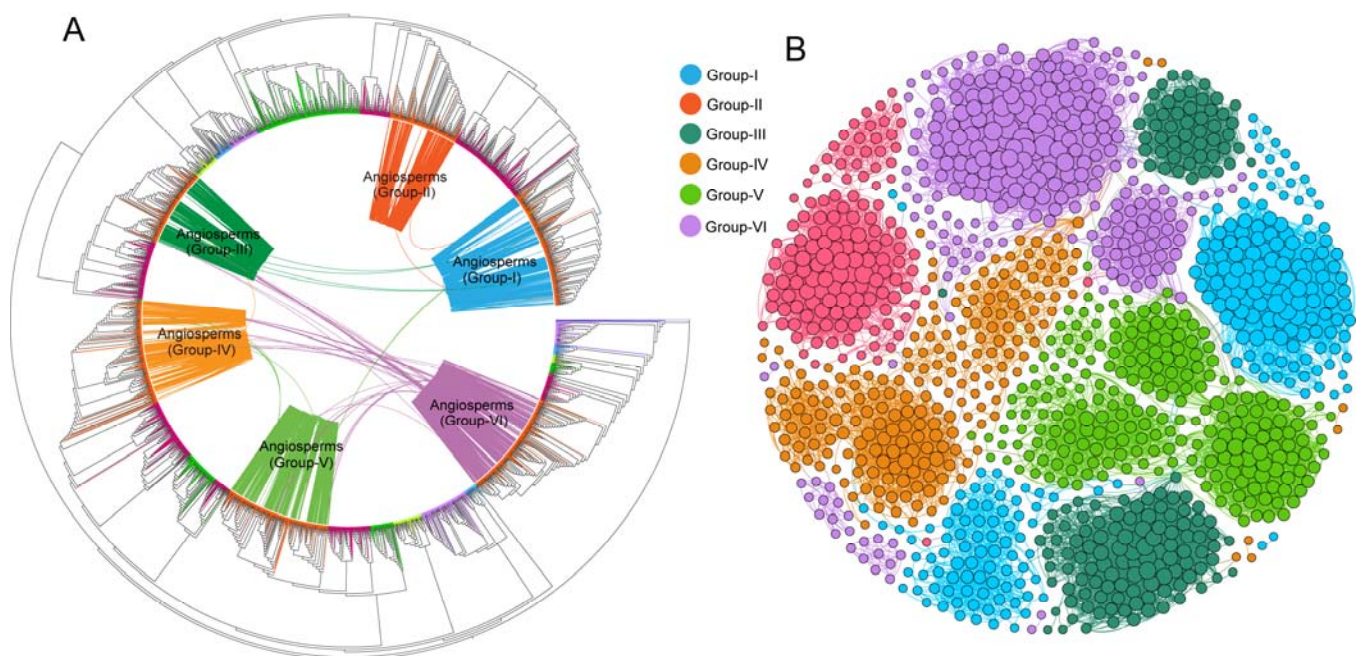685        and Protein Kinases. *Mol Plant*, **9**, 1667-1670.

686    **Figures and Figure legends:**



687

688    **Figure 1 - The broad-scale family phylogeny of *ARF* genes in plants.** The broad-scale family phylogeny of ARF
689    genes in different plant lineages estimated using IQ-TREE maximum-likelihood algorithm. Branches representing *ARF*
690    genes from different plant clades were colored and six conspicuous groups for angiosperm *ARF* genes were obtained
691    and labeled on the tree (Groups I through VI). Gene tree structure and the three subfamilies (denoted as clades A, B and
692    C) were consistent with that reported in (Finet *et al.* 2013) with more *ARF* genes identified from plant genomic or
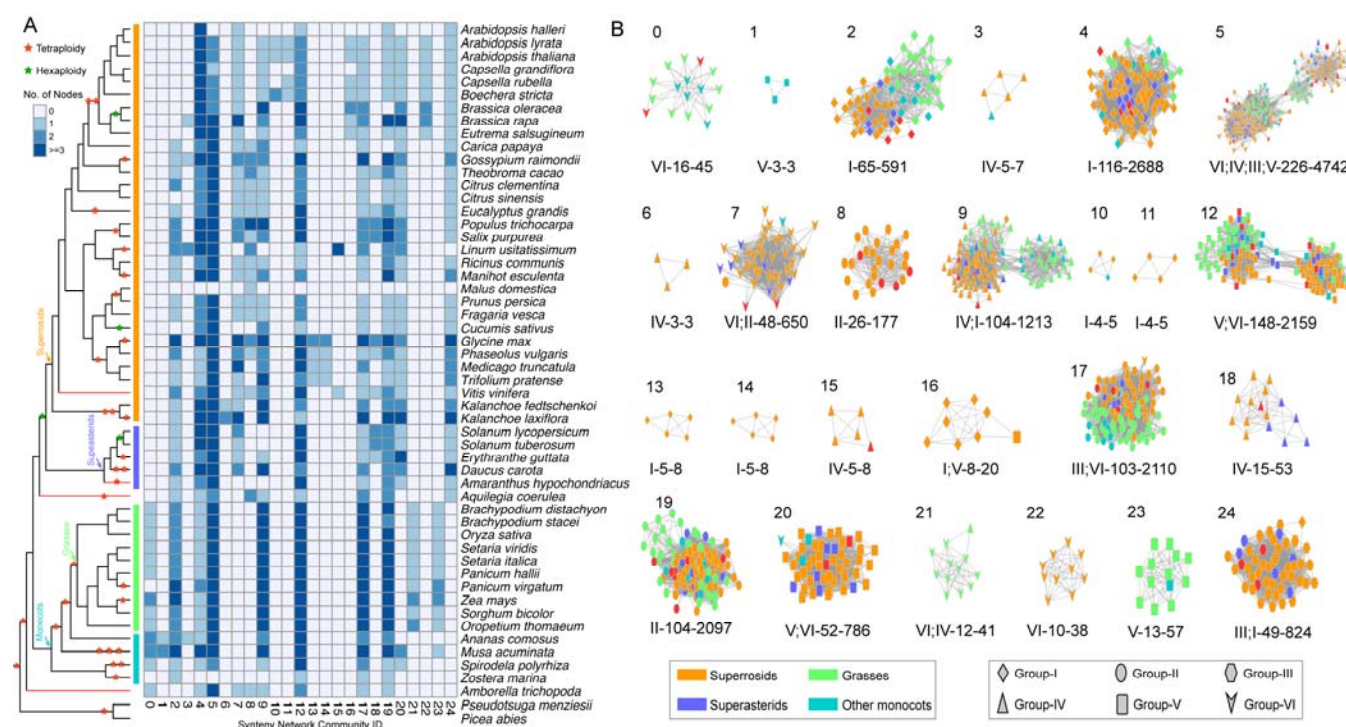693    transcriptomic datasets.

694

**Figure 2 - Ancient evolutionary trajectory of *ARF* genes in streptophytes.** Annotation of auxin response factor genes suggested an early origin and diversification in streptophytic algae. The eight major plant lineages were represented with different colors, solid rounded rectangles indicate presence of ARF genes in corresponding plant lineages, dashed ones suggested absence of data and potential gene losses. Gene groups indicated with an asterisk represented close sister lineages to corresponding angiosperm *ARFs*. Inferred ancient gene duplications were depicted as golden squares. Group-I and Group-II ARFs was suggested to be derived from the seed plant duplication event.



701

20

**Figure 3 - Genomic synteny analyses of ARF genes among angiosperms genomes. (A)** Maximum-likelihood gene tree for the ARF gene family with genomic syntenic relationships between the genes. Each connecting line located inside the inverted circular gene tree (implemented in iTOL) indicates a syntenic relationship between two *ARF* genes (syntelogs). The connecting lines are colored in congruence with the six angiosperm ARF groups. **(B)** Synteny network of the ARF gene family using detected syntenic relations extracted from the genome synteny network database, using nodes representing *ARF* loci and edges (connecting lines) representing syntenic relationships. Colors of the nodes represented the six groups of *ARF* genes in angiosperms and size of each node indicates its connectivity (bigger nodes have more connections). The synteny network were clustered and visualized using the 'Fruchterman Reingold' layout implemented in Gephi.



**Figure 4 - Detailed network representations for ARF synteny network communities among angiosperm genomes. (A)** Species composition for each of the 25 network communities. Blue-colored cells depict the presence of ARF syntelogs in the different species. The 25 network communities were identified using CFinder at k=3. **(B)** Detailed visualization for each of the ARF synteny network communities. Nodes in different colors represented different plant lineages, and the node shapes represented different ARF subfamilies (Groups I through VI). Selected basal lineages including *Vitis vinifera* (sister to other rosids), *Amarsanthus hypochondriacus* (sister to asterids), and *Amborella trichopoda* (the basal angiosperm) were depicted as red nodes in the communities. Each network community was presented following the 'ARF clades'-'number of nodes'-'number of connections' nomenclature system and some communities contain genes from multiple ARF clades. One synteny community may contain genes from various groups.

**Figure 5 - Phylogenetic and synteny network analyses for each of the six groups of ARFs in angiosperms.** Maximum-likelihood trees for each of the six ARF groups were constructed, genes from different species groups were colored using different colors and genes detected in syntenic genomic blocks (syntelogs) were connected using curved lines. The syntenic connections belonging to different synteny network communities were plotted using different colors. Synteny network communities were numbered according to that depicted in figure 2B. Inferred ancestral transposition activities were indicated by red arrows.