# A genome for gnetophytes and early evolution of seed plants

Tao Wan[1,2,‡], Zhi-Ming Liu[3,‡], Ling-Fei Li[1,‡], Andrew R. Leitch[4,‡], Ilia J. Leitch[5,‡], Rolf Lohaus[6,7], Zhong-Jian Liu[8,‡], Hai-Ping Xin[2,9], Yan-Bing Gong[10], Yang Liu[1], Wen-Cai Wang[4], Ling-Yun Chen[2,11], Yong Yang[12], Laura J. Kelly[4], Ji Yang[13], Jin-Ling Huang[14,15], Zhen Li[6,7], Ping Liu[1], Li Zhang[1], Hong-Mei Liu[1], Hui Wang[1], Shu-Han Deng[3], Meng Liu[3], Ji Li[3], Lu Ma[4], Yan Liu[3], Yang Lei[3], Wei Xu[3], Ling-Qing Wu[3], Fan Liu[2], Qian Ma[10], Xin-Ran Yu[3], Zhi Jiang[3], Guo-Qiang Zhang[8], Shao-Hua Li[16], Rui-Qiang Li[3], Shou-Zhou Zhang[1], Qing-Feng Wang[2,11*], Yves Van de Peer[6,7,17*], Jin-Bo Zhang[3*] & Xiao-Ming Wang[1*]


[1]Key Laboratory of Southern Subtropical Plant Diversity, Fairy Lake Botanical Garden, Shenzhen & Chinese Academy of Science 518004, PR China. [2]Sino-Africa Joint Research Centre, Chinese Academy of Science, Wuhan 430074, PR China. [3]Novogene Bioinformatics Institute, Beijing 100083, PR China. [4]School of Biological and Chemical Sciences, Queen Mary University of London, Mile End Road, London E1 4NS, UK. [5]Jodrell Laboratory, Royal Botanic Gardens, Kew, Surrey TW9 3AB, UK. [6]Department of Plant Biotechnology and Bioinformatics, Ghent University, Ghent 9052, Belgium. [7]Centre for Plant Systems Biology, VIB, Ghent 9052, Belgium. [8]Shenzhen Key Laboratory for Orchid Conservation and Utilization, National Orchid Conservation Centre of China and Orchid Conservation and Research Centre, Shenzhen 518114, PR China. [9]Key Laboratory of Plant Germplasm Enhancement and Specialty Agriculture, Wuhan Botanical Garden, Chinese Academy of Sciences, Wuhan 430074, PR China. [10]State Key Laboratory of Hybrid Rice, College of Life Sciences, Wuhan University, Wuhan 430072, PR China. [11]Key Laboratory of Aquatic Botany and Watershed Ecology, Wuhan Botanical Garden, Chinese Academy of Sciences, Wuhan 430074, PR China. [12]State Key Laboratory of Systematic and Evolutionary Botany, Institute of Botany, Chinese Academy of Sciences, Beijing 100093, PR China. [13]Education Key Laboratory for Biodiversity Science and Ecological Engineering, Fudan University, Shanghai 200438, PR China. [14]Institute of Plant Stress Biology, State Key Laboratory of Cotton Biology, Henan University, Kaifeng, 475001, PR China. [15]Department of Biology, East Carolina

University, Greenville, NC27858, USA. [16]Beijing Key Laboratory of Grape Sciences and

Enology, Institute of Botany, Chinese Academy of Sciences, Beijing 100093, PR China.

[17]Genomics Research Institute, University of Pretoria, Private bag X20, Pretoria 0028, South

Africa.

[+]These authors contributed equally to this work. Correspondence should be addressed to

Q.F.W. (qfwang@wbgcas.cn), Y.V.d.P. (yves.vandepeer@psb.vib-ugent.be), J.B.Z.

(zhangjinbo@novogene.com) or X.M.W. (719921868@qq.com).

**Abstract**

Gnetophytes are an enigmatic gymnosperm lineage comprising three genera, *Gnetum*, *Welwitschia* and *Ephedra*, which are morphologically distinct from all other seed plants. Their distinctiveness has triggered much debate as to their origin, evolution, and phylogenetic placement amongst seed plants. To increase our understanding of the evolution of gnetophytes, and their relation to other gymnosperms and seed plants, we report here a high-quality draft genome sequence for *Gnetum montanum* - the first for any gnetophyte. By using a novel genome assembly strategy to deal with high levels of heterozygosity, we assembled > 4 Gb of sequence encoding 27,491 protein-coding genes. Comparative analysis of the *G. montanum* genome with other gymnosperm genomes unveiled some remarkable and distinctive genomic features, such as a diverse assemblage of retrotransposons with evidence for elevated frequencies of elimination rather than accumulation, considerable differences in intron architecture, including both length distribution and proportions of (retro) transposon elements, and distinctive patterns of proliferation of functional protein domains. Furthermore, a few gene families showed *Gnetum*-specific copy number expansions (e.g. CesA) or contractions (e.g. LEA), which could be connected with *Gnetum*'s distinctive morphological innovations associated with their adaptation to warm, mesic environments. Overall, the *G. montanum* genome enables a better resolution of ancestral genomic features within seed plants, and the identification of genomic characters that distinguish *Gnetum* from other gymnosperms.


**Introduction**

The seed plants today are represented by five distinct lineages: the species-rich angiosperms (flowering plants, *c*. 352,000 species) and four gymnosperm lineages (which together comprise *c*. 1,000 species and encompass cycads, *Ginkgo biloba*, conifers and gnetophytes). It is apparent from their long fossil record (dating back to the Late Devonian *c*. 360 million years ago (Mya)) that considerably greater seed plant diversity existed in the past[1]. Nevertheless, widespread extinctions among many

66  gymnosperm lineages means that today's gymnosperms are only a relic of their

67  former diversity, and this has presented a major challenge for reconstructing

68  evolutionary relationships between the extant lineages[2]. Probably the most

69  controversial outstanding question in plant evolution is the phylogenetic position of

70  gnetophytes[3] (comprising the genera *Gnetum*, *Welwitschia* and *Ephedra*, Fig. 1) in

71  relation to the other seed plant lineages. Apparent morphological similarities with

72  angiosperms, such as vessel-like water conducting cells, double fertilization, and leaf

73  morphologies with reticulate venation, have historically led to the proposition that

74  gnetophytes form a group that is sister to angiosperms (termed the 'Anthophyte

75  hypothesis')[4,5]. That hypothesis has, however, largely been rejected by molecular

76  phylogenetic data and a deeper understanding of the developmental pathways that

77  lead to similar morphological features. Nevertheless, the use of molecular data has

78  also been problematic in inferring the exact phylogenetic position of gnetophytes,

79  with topologies differing depending on the type of sequence data (e.g. plastid versus

80  nuclear genes, nucleotide versus amino acid data) and analytical approach used (e.g.

81  maximum parsimony, maximum likelihood, Bayesian, multispecies coalescent based

82  methods)[6-8]. Consequently, several possible hypotheses have been put forward that

83  place gnetophytes as sister to: (i) Pinaceae ('Gnepine' hypothesis); (ii) cupressophytes

84  ('Gnecup' hypothesis); (iii) all conifers ('Gnetifer' hypothesis); (iv) all other

85  gymnosperms; or (v) all seed plants[9]. Currently, the emerging consensus, based on

86  both older and more recent studies, and recently released data from the 1KP initiative

87  (see https://sites.google.com/a/ualberta.ca/onekp/, and Wickett et al. (8)), indicates

88  that gnetophytes are sister to, or within, the conifers.

89  So far, the availability of whole genome sequences for gymnosperms has been limited

90  to conifers (specifically to Pinaceae)[10-13] and *G. biloba*[14], with no whole genome

91  assemblies available for the two remaining major seed plant lineages - cycads and

92  gnetophytes. This deficiency, together with the conflicting phylogenetic evidence for

93  relationships among these groups, is impeding our understanding of genome evolution

94  across all seed plants. Here, we present a high-quality draft genome of *G. montanum*,

95  the first for gnetophytes. The availability of this genome, as well as survey sequence

96  data and transcriptome data from other vascular plants (including novel data from

97  gnetophytes *Ephedra* and *Welwitschia*), enables us to compare genomic characters

98  with *G. biloba*, conifers, angiosperms and non-seed plants. Comparisons within

99  gymnosperms, and between gymnosperms and angiosperms, highlight the unique

100 nature of the *Gnetum* genome, providing new insights into patterns of genome

101 divergence across seed plants.

102

### *Genome assembly and annotation*

104 The genome of *G. montanum* (2n = 44) is small compared with other gymnosperms

105 (flow cytometry: 4.2 Gb / 1C; k-mer analysis: 4.11 Gb), and is highly heterozygous

106 and rich in repeats (Supplementary Fig. 1a-c, and Supplementary Note 1). To

107 overcome problems caused by repeats and heterozygosity, we generated deep

108 coverage (~302 ×, Supplementary Table 1) Illumina sequence data and applied a

109 novel genome assembly strategy (Supplementary Note 2, Supplementary Fig. 2) to

110 assemble 4.07 Gb of sequence (contig N50 size = 25.02 kb, scaffold N50 size =

111 475.17 kb, Supplementary Table 2), to which > 99% of genome reads, > 90% ESTs

112 and > 99% of BACs were mapped (Supplementary Fig. 1d, e, Supplementary Table 3

113 and Note 3).

114 A total of 27,491 protein-coding genes were predicted from this assembly

115 (Supplementary Table 4 and Note 4), 97% of which were supported by orthology (>

116 50% coverage of high-scoring segment pair, Supplementary Fig. 3a) with existing

117 protein sequences and/or RNA-seq data from multiple tissues (Supplementary Table

118 5). A BUSCO analysis to assess the quality of the genome and annotation

119 completeness suggested that 81% of the genes have been recovered (Supplementary

120 Table 6). Unlike conifer genomes, which contain numerous pseudogenes[15] (e.g. 8,328

121 in *Picea abies*, 13,550 in *Pinus taeda*), many fewer were found in the *G. montanum*

122 genome (3,122, Supplementary Note 5). The read depth distribution across genic

123 regions (Supplementary Fig. 3b) suggested little sequence redundancy caused by

124    heterozygosity (see Supplementary Fig. 3c for further confirmation of gene assembly

125    quality).

126

127    ***Repetitive sequence dynamics***

128    Repetitive sequences have been shown to account for the major component of all

129    gymnosperm genomes that have been sequenced to date[11-14], with diverse and ancient

130    transposable elements (TEs), especially LTR retrotransposons (LTR-RTs), being

131    particularly prevalent. Overall, the repetitive element content of *G. montanum* was

132    also high (85.9%) and dominated by LTR-RTs (especially *gypsy*-like elements), which

133    comprised 77.4% of the genome (Supplementary Table 8 and Supplementary Note 6).

134    The genome assembly of *G. montanum* is likely to be sufficient to represent most of

135    the LTR-RTs, since their length is typically around 25 kb[16], whilst 90% of the

136    scaffolds are larger than 34 kb. Phylogenetic reconstructions of the reverse

137    transcriptase domains of LTR-RTs in *G. montanum* and *P. taeda* revealed that most of

138    the *gypsy*- and *copia*-like elements in *G. montanum* were restricted to just a few clades,

139    representing only a small minority of the diversity encountered in *P. taeda*

140    (Supplementary Fig. 4, Supplementary Note 6).

141    Comparative analyses of repeats identified by RepeatExplorer using survey sequence

142    data from multiple gnetophytes (*G. montanum, G. gnemon*, *W. mirabilis* and

143    *E. altissima*) and *P. taeda* revealed substantial differences in the abundance of the

144    major repeat classes (Supplementary Fig. 5a, Supplementary Table 9 and

145    Supplementary Notes 1, 7). Further, the majority of individual repeat types (repeat

146    clusters in RepeatExplorer) were shown to be species-specific (i.e. containing

147    Illumina reads from just one species, data not shown). The species-specific nature of

148    the repeat profiles probably reflects the long estimated divergence times between

149    species (e.g. the two *Gnetum* species likely diverged between *c.* 25 Mya and 75

150    Mya)[17,18].

151    Previously, it was reported from conifers and *G. biloba* that LTR-RTs have

152  accumulated steadily over the last *c*. 25 Mya, especially between 16-24 Mya, a

153  process contributing to their large genome sizes[11,12,14]. This interpretation is consistent

154  with the data here (Supplementary Table 10), which shows that most LTR-RTs in

155  conifers are intact (solo LTR / intact LTR ratio ranged from 0.16:1 to 0.72:1,

156  Supplementary Table 10). It is notable that the solo LTR / intact LTR ratio was

157  substantially higher in *G. montanum* (~1.94:1), which together with its small genome

158  and similar profile of accumulation (Supplementary Fig. 5b), suggests higher

159  frequencies of LTR-RT elimination than amplification compared with *G. biloba* and

160  conifers.

161  Most angiosperm genomes analysed to date have far fewer ancient repeats and less

162  divergent LTR-RT subsets than conifers and *G. biloba*, presumably due to more

163  efficient elimination and replacement processes operating within these angiosperm

164  genomes[19] (e.g. in *Oryza sativa* the half-life of LTR-RTs is estimated to be less than

165  five million years[20], leading to "genome turnover"[21]). However, an exception to this

166  pattern has been observed in *Amborella trichopoda*. The genome of this species is

167  considered to have retained many features that were likely present in the ancestral

168  angiosperm genome[22]. It is notable that its repeat content[13] and lower abundance of

169  intact LTR-RTs (i.e. solo LTR / intact LTR ratio = 2.43/1.0; Supplementary Table 10)

170  is similar to that observed in *G. montanum*. These observations suggest that neither *A.*

171  *trichopoda* nor *G. montanum* genomes have experienced recent, extensive (retro)

172  transposon activity, although they continue to eliminate repetitive sequences. Both

173  these species seem to differ from conifers and *G. biloba* with respect to the dynamics

174  of repeat accumulation[11,12,14], and from other angiosperms in terms of the levels of

175  repeat amplification/removal.

176

177  ***Intron morphologies***

178  Although intron size has been positively correlated with genome size across

179  eukaryotes as a whole[23], this trend does not translate well across broad and some

180  narrow taxonomic distances in seed plants (Fig. 2a). Previous studies of *G. biloba*[14]

181  and conifers[11,12] have reported larger introns than angiosperms, probably arising from

182  the long-term, steady amplification of LTR-RTs (Fig. 2b), as also observed here,

183  where LTR-RTs account for 51% and 59% of the large intron sequences in *P. taeda*

184  and *G. biloba*, respectively (Fig. 2a, Supplementary Table 12). The evolution of these

185  large introns may have arisen from similar repeat accumulation processes that are

186  operating across the genome as a whole.

187  When comparing these observations with introns of *G. montanum*, it is apparent that

188  their introns are substantially smaller (minimum, mean and maximum intron lengths)

189  than those of *P. taeda* and *G. biloba* (Fig. 2a, see also statistics test in Supplementary

190  Table 11). In addition, the repeat composition of *G. montanum*'s introns is dominated

191  by both long interspersed nuclear elements (LINEs) as well as LTR-RTs, rather than

192  predominantly LTR-RTs, as in conifers and *G. biloba* (Fig. 2b, Supplementary Table

193  12). The correlation between smaller intron sizes and smaller genome size in *G.*

194  *montanum* compared with conifers and *G. biloba* may reflect the repeat dynamic

195  processes operating across its genome as a whole. In contrast, the variable length

196  distributions of introns in angiosperms suggest that the evolution of repeats in their

197  introns do not necessarily reflect the repeat dynamics observed across the rest of their

198  genomes[24]. In the highly dynamic repetitive genome of *Z. mays*, the profile of repeats

199  across the genome[25] and within the whole intron set (Supplementary Fig. 6a) both

200  suggests many recent insertions. However, in *A. trichopoda*, the intron sizes are

201  overall larger, and the genome size smaller than in *Z. mays* (Fig. 2a, b). In addition, an

202  analysis of introns in *A. trichopoda* and *G. montanum* highlighted a closer similarity

203  to each other (in terms of length distributions, repeat composition and divergence)

204  than either species has to conifers and *G. biloba*, despite a 4.8-fold difference in their

205  genome sizes (Fig. 2a, 2b, Supplementary Table 12).

206  Previous comparisons of orthologous introns have led to the suggestion that the

207  expansion of introns occurred early in the evolutionary history of conifers[12].

208  Comparisons of orthologous introns (with identical adjacent exons) between *P. taeda*

209  and *G. biloba* showed that introns identified as being long (> 6 kb) in *P. taeda* were

210    also typically long in their orthologues in *G. biloba*, containing, in both cases,

211    abundant LTR-RTs (both *gypsy-* and *copia*-like elements, Fig. 2c). These features

212    were likely to have been present in their most recent common ancestor (MRCA).

213    Using similar approaches to analyse the length and repeat content of 4,348

214    orthologous introns of *G. montanum* shared with *P. taeda* (Supplementary Note 8)

215    highlighted notable differences. Whilst the length of exons remained similar, a

216    substantial fraction of orthologous genes had longer introns in *P. taeda*

217    (Supplementary Fig. 6b). The introns identified as 'short' in *P. taeda* comprised *c.* 4%

218    repeats, rising to *c.* 56% in 'long' introns, largely through the accumulation of

219    LTR-RTs (especially *copia* elements) (Fig. 2d, Supplementary Table 13). In contrast,

220    introns in *G. montanum* that are orthologous to the 'long' introns of *P. taeda* (36% of

221    introns analysed) showed high proportions of LINEs. As with comparisons of all

222    introns, pairwise comparisons of orthologous introns in *G. montanum* and *A.*

223    *trichopoda* again showed some similarities in their introns, with both species having

224    abundant LINEs (Fig. 2e). Collectively, these data reveal a different repeat dynamic

225    within introns of *G. montanum* compared with the other gymnosperms.

226

227    *('Lack of') Whole genome duplication* **(WGD)**

228    All angiosperms are reported to have undergone at least one round of ancient WGD,

229    and in many lineages WGDs are recurrent and ongoing[26]. In addition, a WGD event

230    has been proposed at the base of all seed plants *c.* 341 Mya (= *zeta* WGD[27]), although

231    the underlying evidence for these two ancient WGD events has been recently

232    questioned[28]. In gymnosperms, WGDs have been reported for conifers, *G. biloba* and

233    cycad (a likely shared WGD)[14,29,30]. Although recent polyploidy seems common in

234    extant *Ephedra*[31]*,* evidence for ancient WGDs in gnetophytes is missing

235    (Supplementary Note 9 and Supplementary Fig. 7), except for a WGD in *Welwitchia*

236    which likely occurred after the divergence of its lineage from that leading to *Ephedra*

237    (Supplementary Fig. 7)[29]. If indeed the ancient zeta WGD is shared by all seed plants,

238    the absence of evidence for this event in gnetophytes is best explained by their faster

239  rates of gene evolution compared with other gymnosperms[32,33], erasing all evidence of

240  this more than 300 million year old event (Supplementary Note 9 and Supplementary

241  Fig. 7).

242

243  ***Organization of functional protein domains***

244  To characterize the patterns of functional diversification in gene domains across land

245  plants, we used principal component analysis (PCA) to analyse the number of pfam

246  domains (conserved protein domains) in multiple species (Supplementary Note 10,

247  Supplementary Table 13). Our approach showed that angiosperms formed a discrete

248  cluster that was separate from the gymnosperms (Fig. 3a), with *G. montanum* being an

249  outlier. Indeed, heatmaps compiled from the pfam data that contributed most (top 10%)

250  to PCA1 and PCA2 showed that *G. montanum* formed a clade with the lycophyte *S.*

251  *moellendorffii* and the moss *Physcomitrella patens* (Fig. 3b), whilst the

252  non-gnetophyte gymnosperms formed a separate clade (Fig. 3b).

253  Given the distinct distributions of *G. montanum*, non-gnetophyte gymnosperms and

254  angiosperms in the PCA analysis, the data suggest that significant functional

255  diversification of the conserved protein domains has occurred since these major

256  lineages split. It may be surprising given the long divergence times (*c.* 300 Mya)[2], that

257  *G. biloba* and conifers retain similar conserved domain organizations (with similar

258  eigenvector values). This could reflect their relatively low substitution rates (on

259  average $7 \times$ lower) compared with angiosperms[33].

260  An analysis of the pfam domain expansions that contributed most to the PCA1 and

261  PCA2 distributions amongst angiosperms (except *A. trichopoda*). included genes

262  associated with flower and organ development (Supplementary Table 15). In contrast,

263  non-gnetophyte gymnosperms showed large-scale specific expansions of pfam

264  domains in genes associated with defence and secondary metabolism, as previously

265  suggested (Supplementary Table 16)[10,11]. The clustering of *G. montanum* with

266  non-seed plants in the heatmap (Fig. 3b) was a surprise, and may indicate the

267 approach has identified proteins that have diverged very little since the MRCA of seed

268 plants. Nevertheless, such an explanation is at odds with the hypothesis that the genes

269 of gnetophytes have diverged rapidly, given their comparatively high substitution rate

270 compared with other gymnosperms[33].

271

### *Growth form (shrubs and lianas) and leaf morphology*

273 Gnetophytes differ from other extant gymnosperms in growth form, with the unusual

274 and distinct form of *Welwitschia*, the shrub habit of *Ephedra* and the shrub and liana

275 habit and specialized leaf morphologies of *Gnetum*[34]. Cellulose synthase (*Ces*A) and

276 cellulose synthase-like (*Csl*) genes are considered to play a role in influencing the

277 biomechanical properties of the cell[35], hence potentially the distinctive growth forms

278 of gnetophytes are associated with the divergence of these genes. To explore this

279 hypothesis, *Ces*A and *Csl* family members were examined in *G. montanum* and

280 compared with those in other seed plants. The total number of *Ces*A and *Csl* family

281 members ranged about 3-fold amongst the seed plants analysed (*P. abies*, *P. taeda*, *A.*

282 *trichopoda*, *A. thaliana* and *O. sativa*). However, only *G. montanum* showed a large

283 expansion of the *Csl*B/H gene subfamily (to 20 genes, Supplementary Table 17),

284 involving tandem duplications (Supplementary Fig. 9), and accounting for two-thirds

285 of its total *Csl* gene repertoire. Furthermore, transcriptome analysis showed that these

286 *Csl*B/H genes were differentially expressed in leaves, stems and roots of *G. montanum*,

287 supporting an association with distinct growth forms and leaf morphologies

288 (Supplementary Fig. 9). In contrast, all other species analysed, including *Welwitschia*

289 and *Ephedra,* were seen to have only 1-6 *Csl*B/H genes (at least based on

290 transcriptome analysis) (Supplementary Note 11, Supplementary Table 16,

291 Supplementary Fig. 8).

292 Another gene family associated with leaf morphology and development is the *WOX*

293 (*WUSCHEL-related homeobox*) family[36]. Recent studies have shown that the

294 conserved family members WOX3 and WOX4, which play a role in leaf

295 development, show diffuse *WOX3* expression at the leaf bases of *Arabidopsis* and

296  *Gnetum*, with such patterns being associated with the distinctive reticulate venation

297  observed in their leaves[37]. Two unusual paralogues, GgWOXX and GgWOXY, were

298  previously reported to occur only in gnetophytes[37], and this is confirmed here in

299  phylogenetic reconstructions of gene family members (Supplementary Note 12,

300  Supplementary Fig. 10). These paralogues are unlikely to have arisen by

301  *Gnetum*-specific gene amplifications, as this would group them with other *Gnetum*

302  paralogues. Alternatively, these genes may correspond to ancestral seed plant

303  sequences that have been lost in other plant lineages. Potentially the different patterns

304  of gene loss, retention and amplification compared with other gymnosperms may be

305  associated with their distinctive growth forms.

306

307  *Vessels*

308  The presence of vessel-like water-conducting cells, morphologically distinct from

309  tracheids, is another feature that sets gnetophytes apart from other gymnosperms.

310  However, there has been long-standing debate as to whether gnetophyte "vessels" are

311  homologous to the "vessels" of angiosperms. In angiosperms,

312  VASCULAR-RELATED NAC-DOMAIN (VND) proteins *VND1*-7 are members of

313  the *NAC* domain class of transcription factors, *VND7* being a master regulator of

314  vessel formation in *Arabidopsis thaliana*[38], and *VND1-6* being upstream regulators of

315  *VND7*[39]. Although five *NAC* domain genes were identified in the genome of *G.*

316  *montanum*, no orthologues of *VND*7 or *VND1-3* in the sister clade were identified,

317  consistent with previous analyses of other gymnosperms[12], and suggesting that these

318  proteins are restricted to angiosperms (Supplementary Fig. 11). Nevertheless, *Gnetum*

319  does share the *VND4-6* clade with angiosperms and other gymnosperms. Furthermore,

320  *A. trichopoda,* which lacks angiosperm vessels, also lacks orthologues of *VND1-3*, but

321  it does have *VND7* (Supplementary Fig. 11), indicating that the ability to form vessels

322  may have occurred after angiosperms diverged. Taken together, these data suggest a

323  greater dependency of vessel development on *VND1-3* than is apparent from

324  experiments on *A. thaliana*. The most parsimonious explanation of our data is that

325   angiosperm vessel formation requires genes from the *VND7* clade (and potentially its

326   sister clade *VND1-3*), and that gymnosperms, including gnetophytes, which lack

327   sequences from both these clades cannot form structures that are homologous to

328   angiosperm vessels. Such an interpretation supports Carlquist's[40] morphological

329   interpretations of vessels. It is therefore most likely that different molecular

330   mechanisms underpin the origin and development of vessels in *Gnetum* and

331   angiosperms. Indeed, these new molecular data support the hypothesis based on

332   morphological studies that *Gnetum* vessels are actually more closely related to conifer

333   tracheids than angiosperm vessels and that vessels in the two groups are convergent

334   characters[40].

335

336   ***Water stress***

337   Extant species of *Gnetum* are unusual amongst gymnosperms in being restricted to

338   warm, mesic habitats[41], this contrasts to conifers that are adapted to cold and

339   water-stressed environments. An analysis of genes involved in water and cold stress

340   revealed some substantial differences between conifers and *Gnetum*. The Late

341   Embryogenesis Abundant protein (LEA) gene family encodes crucial proteins that are

342   involved in protecting plants from desiccation or osmotic stresses associated with low

343   temperature[42,43]. An analysis of LEA family members suggests that some members

344   have been reduced in number in *Gnetum* or expanded in conifers (e.g. LEA-3), or lost

345   completely in *Gnetum* (i.e. LEA-4, 5, 6). In addition, dehydrins, which play a role in

346   the response to cold/drought[44], had only two members in *G. montanum*, compared

347   with 38 in *P. abies,* 28 in *P. taeda* and 3-15 in angiosperms (Supplementary Table 19).

348   Further analysis of the *G. montanum* genome also revealed relatively few gene family

349   members of the AP2 domain containing protein families, which are involved in the

350   cold stress response[45,46], and GPX and GST families, involved in the oxidant stress

351   response[47,48]. Taken together, these data appear consistent with the hypothesis that the

352   ecological shift to a warm, wet forest habitat is associated with a relaxation of

353   selection pressure on genes associated with water stress and low temperature.

354

355 *Conclusion*

356 Here, we have described the assembly, annotation, and comparative analysis of the

357 first gnetophyte genome, namely that of *G. montanum*. Its genome is particularly

358 enigmatic given a phylogenetic position within or sister to conifers. It also carries

359 genomic peculiarities that may reflect its morphological and ecological uniqueness

360 amongst gymnosperms. Comparisons of these genome features with the genomes of

361 conifers and *G. biloba* provide opportunities to predict the nature and direction of

362 genomic change accompanying the evolution of the lineage leading to *Gnetum* (Fig.

363 4). Assuming that gnetophytes do indeed form a clade that is sister to, or within, the

364 conifers, the following genomic features can be predicted to have been present in the

365 MRCA of the gymnosperms, as observed in *G. biloba*[14] and conifers[11,12]: (1) A large

366 genome size (1C > 10 Gb) comprised predominantly of a heterogeneous set of large

367 numbers of LTR-RTs associated with low levels of repeat deletion[14]; (2) Long introns

368 predominantly shaped by insertions of LTR-RTs (*gypsy* and *copia* elements); (3) Pfam

369 domains that show a profile distinct from angiosperms; If this is so, and assuming a

370 common ancestry of gnetophytes and conifers, these genomic characters, or their

371 signatures, have subsequently been lost or diverged considerably in the lineage

372 leading to *Gnetum*. This most likely involved the following genomic processes: (1)

373 Genome downsizing, leading to the relatively (for a gymnosperm) small genomes of

374 *Gnetum* species (1C= 2.25-4.11 Gb). This is supported by the high ratio of solo LTR /

375 intact LTR-RTs observed in the genome of *Gnetum* compared with conifers, and is

376 indicative of the activity of recombination-based processes, which can eliminate DNA

377 from the genome. Similar processes leading to genome downsizing have also been

378 reported in many angiosperms, resulting in small genomes despite the occurrence of

379 multiple rounds of polyploidy detected in many lineages[49]; (2) Reduction in the size

380 of introns in *G. montanum* and a replacement of many of the LTR-RTs repeats with

381 LINEs to give rise to introns that are more similar to those of, for instance, *A.*

382 *trichopoda* than to other gymnosperms; (3) Elevated rates of sequence divergence

383   causing the erosion of a hypothesised shared seed-plant WGD event and leading to a

384   pattern of Pfam domains, which is distinct from the remaining gymnosperms; (4)

385   Expansion and contraction of specific gene families associated with adaptation to new

386   ecologies.

387

388 **Methods summary**

389 The sequenced *G. montanum* is a single mature female individual growing naturally

390 in Fairy Lake Botanical Garden, Shenzhen, China. Genome sequences were generated

391 using an Illumina platform and assembled with a novel hierarchical assembly strategy.

392 Gene annotations were determined by integrating results from both *de novo* prediction

393 approaches and alignment-based methods based on orthology and transcriptomic data.

394 RNA-seq was performed using an Illumina platform. All methods and bioinformatic

395 analyses are detailed in the Supplementary Information.

396

397 **Data availability**

398 The *G. montanum* genome project has been deposited at the NCBI under the

399 BioProject number PRJNA339497. The whole genome sequencing data were

400 deposited in the Sequence Read Archive (SRA) database under the accession number

401 SRX2052734, SRX2098865, SRX2099144, SRX2114825, SRX2114827,

402 SRX2134147, SRX2134160, SRX2134177, SRX2134180, SRX2134596, and

403 SRX2134624. And the *G. montanum* assemblies, gene sequences, and annotation data

404 are also available at the DRYAD website. The data or related program scripts that

405 support the findings of this study are available from the corresponding author upon

406 request.

**Acknowledgements**

**Author contributions**

T.W. and X.M.W. conceived and initiated the study, managing the gnetophytes (*Gnetum*, *Welwitschia*, *Ephedra*) genome sequencing project. T.W. designed the major scientific objectives and led the manuscript preparation together with A.R.L., I.J.L., J.B.Z, L.J.K and Y.V.d.P. The collaboration between groups was close in all aspects of the project. T.W., Z.M.L., L.F.L., A.R.L., I.J.L. and Z.J.L. are joint first authors, H.P.X., Y.B.G., Y.L., L.Y.C. and W.C.W. are joint second authors. Z.M.L., J.B.Z., J.L., Y.L. performed the genome assembly and annotation; H.P.X., L.F.L., L.Y.C., L.M., X.R.Y. contributed to the RNA-seq and corresponding analysis. A.R.L., I.J.L., and W.C.W. coordinated the *RepeatExplorer* analysis in gnetophytes and contributed to the design of the analysis for investigating the dynamics of genome evolution. Z.M.L., J.B.Z., L.F.L., F.L., H.M.L., T.W., A.R.L., I.J.L., W.X. and Y.L. participated in the analyses of LTR-RTs and comparisons of introns. R.L., T.W., Y.V.d.P., Z.L., Z.J.L. and Z.M.L. were involved in the WGD determination; M.L., L.F.L., J.B.Z., J.Y., T.W., L.Z., Y.B.G. and Y.H.D. conducted PCA analysis of pfam domains. J.B.Z., T.W., J.L., L.F.L., L.J.K., Y.L. and Z.M.L. performed the analysis investigating the divergence of gene families. J.L.H., P.L., Q.M., Y.L. and G.Q.Z. contributed to the analysis of pseudogenes; Q.F.W., S.H.L. and S.Z.Z. helped with the collecting of *Welwitschia* and *Ephedra*. Y.Y. provided experimental information on the taxonomic identity of the species used for genome sequencing and collated the distribution records of gnetophytes.

445 **Additional Information**

446 Information on reprints and permissions is available at http://www.nature.com/reprints.

447 Correspondence and requests for materials should be addressed to Q.F.W.

448 (qfwang@wbgcas.cn) or Y.V.d.P. (yves.vandepeer@psb.vib-ugent.be) or J.B.Z.

449 (zhangjinbo@novogene.com) or X.M.W. (719921868@qq.com).

450

451 **Competing interests**

452 The authors declare no competing financial interests.

453

454 **References**

455 1    Rothwell, G. W. & Scheckler, S. E. *Biology of Ancestral Gymnosperms.*
456      *Origin and Evolution of Gymnosperms* (Columbia Uni. Press, 1988).

457 2    Lu, Y., Ran, J. H., Guo, D. M., Yang, Z. Y. & Wang, X. Q. Phylogeny and
458      divergence times of gymnosperms inferred from single-copy nuclear genes.
459      *PLoS One* **9**, e107679 (2014).

460 3    Doyle, J. A. Molecular and Fossil Evidence on the Origin of Angiosperms.
461      *Annu. Rev. Earth Planet. Sci.* **40**, 301-326 (2012).

462 4    Doyle, J. A. & Donoghue, M. J. Seed plant phylogeny and the origin of
463      angiosperms: An experimental cladistic approach. *Bot. Rev.* **52**, 321-431
464      (1986).

465 5    Crane, P. R. Phylogenetic Analysis of Seed Plants and the Origin of
466      Angiosperms. *Annals of the Missouri Botanical Garden* **72**, 716-793 (1985).

467 6    Mathews, S. Phylogenetic relationships among seed plants: Persistent
468      questions and the limits of molecular data. *Am. J. Bot.* **96**, 228-236 (2009).

469 7    Wang, X. Q. & Ran, J. H. Evolution and biogeography of gymnosperms. *Mol.*
470      *Phylogenet. Evol.* **75**, 24-40 (2014).

471 8    Wickett, N. J. *et al.* Phylotranscriptomic analysis of the origin and early
472      diversification of land plants. *Proc. Natl Acad. Sci. USA* **111**, E4859-4868
473      (2014).

474 9    Li, Z. *et al.* Single-Copy genes as molecular markers for phylogenomic studies
475      in seed plants. *Genome Biol. Evol.* **9**, 1130-1147 (2017).

476 10   Warren, R. L. *et al.* Improved white spruce (*Picea glauca*) genome assemblies
477      and annotation of large gene families of conifer terpenoid and phenolic
478      defense metabolism. *Plant J.* **83**, 189-212 (2015).

479 11   Neale, D. B. *et al.* Decoding the massive genome of loblolly pine using
480      haploid DNA and novel assembly strategies. *Genome Biol.* **15**, R59 (2014).

481 12   Nystedt, B. *et al.* The Norway spruce genome sequence and conifer genome
482      evolution. *Nature* **497**, 579-584 (2013).

483 13   Stevens, K. A. *et al.* Sequence of the sugar pine megagenome. *Genetics* **204**,
484      1613-1626 (2016).

485 14   Guan, R. *et al.* Draft genome of the living fossil *Ginkgo biloba. GigaScience*
486      **5**, 49 (2016).

487 15   Garcia-Gil, M. R. Evolutionary aspects of functional and pseudogene
488      members of the phytochrome gene family in Scots pine. *J. Mol. Evol.* **67**,
489      222-232 (2008).

490 16   Wicker, T. *et al.* A unified classification system for eukaryotic transposable
491      elements. *Nat. Rev. Genet.* **8**, 973-982 (2007).

492 17   Won, H. & Renner, S. S. Dating dispersal and radiation in the gymnosperm
493      Gnetum (Gnetales)--clock calibration when outgroup relationships are
494      uncertain. *Syst. Biol.* **55**, 610-622 (2006).

495 18   Hou, C., Humphreys, A. M., Thureborn, O. & Rydin, C. New insights into the
496      evolutionary history of *Gnetum* (Gnetales). *Taxon* **64**, 239-253 (2015).

497    19    Kejnovsky, E., Leitch, I. J. & Leitch, A. R. Contrasting evolutionary dynamics
498          between angiosperm and mammalian genomes. *Trends Ecol. Evol.* **24**,
499          572-582 (2009).

500    20    Ma, J., Devos, K. M. & Bennetzen, J. L. Analyses of LTR-retrotransposon
501          structures reveal recent and rapid genomic DNA loss in rice. *Genome Res.* **14**,
502          860-869 (2004).

503    21    Lim, K. Y. *et al.* Sequence of events leading to near-complete genome
504          turnover in allopolyploid *Nicotiana* within five million years. *New Phytol.*
505          **175**, 756-763 (2007).

506    22    Amborella Genome Project. The Amborella genome and the evolution of
507          flowering plants. *Science* **342**, 1241089 (2013).

508    23    Vinogradov, A. E. Intron-genome size relationship on a large evolutionary
509          scale. *J. Mol. Evol.* **49**, 376-384 (1999).

510    24    Wendel, J. F. *et al.* Intron size and genome size in plants. *Mol. Biol. Evol.* **19**,
511          2346-2352 (2002).

512    25    Schnable, P. S. *et al.* The B73 maize genome: complexity, diversity, and
513          dynamics. *Science* **326**, 1112-1115 (2009).

514    26    Van de Peer, Y., Mizrachi, E. & Marchal, K. The evolutionary significance of
515          polyploidy. *Nat. Rev. Genet.* (2017).

516    27    Jiao, Y. *et al.* Ancestral polyploidy in seed plants and angiosperms. *Nature*
517          **473**, 97-100 (2011).

518    28    Ruprecht, C. *et al.* Revisiting ancestral polyploidy in plants. *Sci. Adv.* **3**,
519          e1603195 (2017).

520    29    Li, Z. *et al.* Early genome duplications in conifers and other seed plants. *Sci.*
521          *Adv.* **1**, e1501084 (2015).

522    30    Roodt, D. *et al.* Evidence for an ancient whole genome duplication in the
523          cycad lineage. *PLoS One* **12**, e0184454 (2017).

524    31    Wu, H. *et al.* A high frequency of allopolyploid speciation in the
525          gymnospermous genus *Ephedra* and its possible association with some
526          biological and ecological features. *Mol. Ecol.* **25**, 1192-1210 (2016).

527    32    Hajibabaei, M., Xia, J. & Drouin, G. Seed plant phylogeny: gnetophytes are
528          derived conifers and a sister group to Pinaceae. *Mol. Phylogenet. Evol.* **40**,
529          208-217 (2006).

530    33    De La Torre, A. R., Li, Z., Van de Peer, Y. & Ingvarsson, P. K. Contrasting
531          rates of molecular evolution and patterns of selection among gymnosperms
532          and flowering plants. *Mol. Biol. Evol.* **34**, 1363-1377 (2017).

533    34    Frohlich, M. W. & Chase, M. W. After a dozen years of progress the origin of
534          angiosperms is still a great mystery. *Nature* **450**, 1184-1189 (2007).

535    35    Popper, Z. A. *et al.* Evolution and diversity of plant cell walls: from algae to
536          flowering plants. *Annu. Rev. Plant Biol.* **62**, 567-590 (2011).

537    36    Nakata, M. *et al.* Roles of the middle domain-specific *WUSCHEL-RELATED*
538          *HOMEOBOX* genes in early development of leaves in *Arabidopsis*. *Plant Cell*
539          **24**, 519-535 (2012).

540    37    Nardmann, J. & Werr, W. Symplesiomorphies in the *WUSCHEL* clade suggest

that the last common ancestor of seed plants contained at least four independent stem cell niches. *New Phytol.* **199**, 1081-1092 (2013).

38    Yamaguchi, M. *et al.* VASCULAR-RELATED NAC-DOMAIN7 directly regulates the expression of a broad range of genes for xylem vessel formation. *Plant J.* **66**, 579-590 (2011).

39    Endo, H. *et al.* Multiple classes of transcription factors regulate the expression of *VASCULAR-RELATED NAC-DOMAIN7*, a master switch of xylem vessel differentiation. *Plant Cell Physiol.* **56**, 242-254 (2015).

40    Carlquist, S. Wood, bark and stem anatomy of New World species of *Gnetum*. *Bot. J. Linn. Soc.* **120**, 1-19 (1996).

41    Ickert-Bond, S. M. & Renner, S. S. The Gnetales: Recent insights on their morphology, reproductive biology, chromosome numbers, biogeography, and divergence times. *J. Syst. Evol.* **54**, 1-16 (2016).

42    Hundertmark, M. & Hincha, D. K. LEA (Late Embryogenesis Abundant) proteins and their encoding genes in *Arabidopsis thaliana*. *BMC Genomics* **9**, 1-22 (2008).

43    Gao, J. & Lan, T. Functional characterization of the late embryogenesis abundant (LEA) protein gene family from Pinus tabuliformis (Pinaceae) in Escherichia coli. *Sci. Rep.* **6**, 19467 (2016).

44    Richard, S., Morency, M. J., Drevet, C., Jouanin, L. & Seguin, A. Isolation and characterization of a dehydrin gene from white spruce induced upon wounding, drought and cold stresses. *Plant Mol. Biol.* **43**, 1-10 (2000).

45    Chinnusamy, V., Zhu, J. & Zhu, J. K. Cold stress regulation of gene expression in plants. *Trends Plant Sci.* **12**, 444-451 (2007).

46    Du, C. *et al.* Dynamic transcriptome analysis reveals AP2/ERF transcription factors responsible for cold stress in rapeseed (*Brassica napus* L.). *Mol. Genet. Genomics* **291**, 1053-1067 (2016).

47    Roxas, V. P., Smith Jr, R. K., Allen, E. R. & Allen, R. D. Overexpression of glutathione S-transferase/glutathioneperoxidase enhances the growth of transgenic tobacco seedlings during stress. *Nat. Biotechnol.* **15**, 988 (1997).

48    Zhao, J. *et al.* Global transcriptional profiling of a cold-tolerant rice variety under moderate cold stress reveals different cold stress response mechanisms. *Physiol. Plant* **154**, 381-394 (2015).

49    Wendel, J. F. The wondrous cycles of polyploidy in plants. *Am. J. Bot.* **102**, 1753-1756 (2015).

577    **Figure Legends**

578

579    **Fig. 1 | Morphological variation and geographical distribution of gnetophytes**

580    **and some other gymnosperms.** Top row from left to right, female cones of *Gnetum*

581    *montanum*, male cones of *Welwitschia mirabilis* and female cones of *Ephedra*

582    *equisetina* (Bar = 5 cm). Below, pantropical distribution of the three gnetophyte

583    genera, compared with three conifer species that are most abundant at higher latitudes

584    and altitudes. The range of genomes sizes (1C-values) found in the three genera

585    comprising gnetophytes and the three conifer species are also shown (data taken from

586    http://data.kew.org/cvalues/ and unpublished data).

587

588    **Fig. 2 | Comparative analysis of seed plant intron morphologies**. (**a**) Intron length

589    distributions and genome sizes (1C-values, depicted by the relative circle size) are

590    shown for nine representative seed plants. (**b**) Distribution of sequence divergence for

591    four types of transposable elements (TEs) in introns of *A. trichopoda*, *G. montanum*,

592    *P. taeda*, and *G. biloba*. The data show that TEs in *G. montanum* and *A. trichopoda*

593    are more diverse than in *P. taeda* and *G. biloba.* The latter two species also show a

594    peak at around 10% sequence divergence probably reflecting a pulse of LTR-RT

595    expansions. (**c**), (**d**) and (**e**), Comparison of orthologous introns between *P. taeda* (Pta)

596    vs. *G. biloba* (Gbi) (**c**), *P. taeda* (Pta) vs. *G. montanum* (Gmo) (**d**) *G. montanum*

597    (Gmo) vs. *A. trichopoda* (Atr) (**e**). Two orthologous intron sets that differed more than

598    two-fold in length were examined, i.e. 'short' introns = 0.5-3 kb and 'long' introns ≥

599    6 kb. Orthologous introns that were 'long' in one species were also found to be 'long'

600    in the other species of the pair. Analysis of the TEs in orthologous introns showed the

601    'long' introns of *G. montanum* and *A. trichopoda* carried a high proportion of LINEs,

602    contributing to intron expansion. In contrast, *gypsy* and *copia* LTR-RT elements

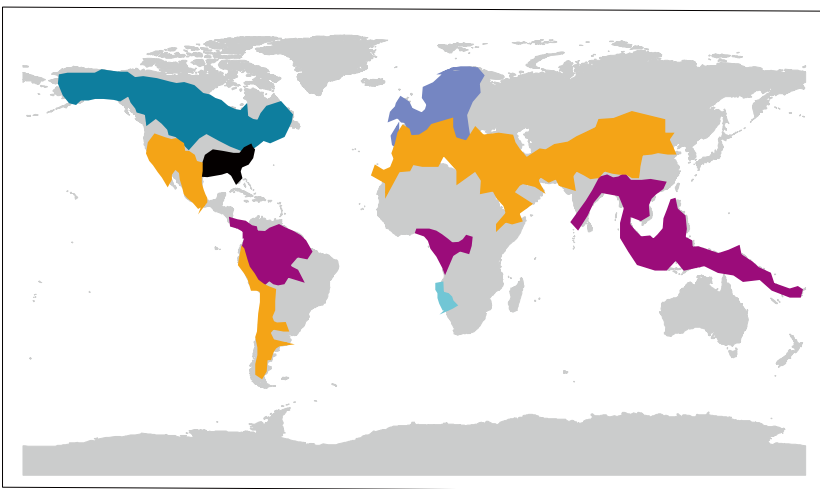603    contributed most to intron expansion in *P. taeda* and *G. biloba*.

604

605 **Fig. 3 | Genome-wide analysis to show the contrasting diversification of**
606 **functional protein domains across land plants**. (**a**) PCA analysis of the occurrence
607 and number of pfam domains in multiple orthologous genes across land plants.
608 Plotting PC1 against PC2 reveals that monocots and eudicots cluster together, as do
609 conifers with *G. biloba*, whilst the remaining species are separate from these clusters.
610 (**b**) Heatmaps reveal the ancestral coding repertories shared by *S. moellendorffii* and *G.*
611 *montanum*. Different patterns of expansion and contraction of the pfam domains are
612 seen for other gymnosperms and angiosperms (see **Supplementary Table 7** for
613 species name list and corresponding abbreviations).

614

615 **Fig. 4 | Prediction of patterns of genome divergence across seed plants.** The origin
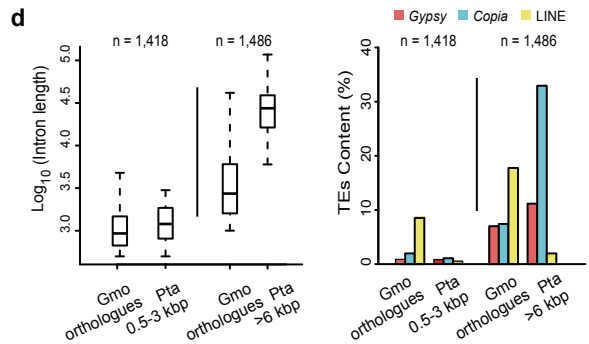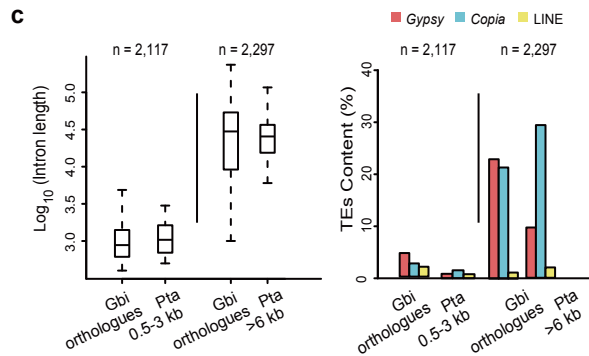616 and evolution of distinctive genomic features observed in *G. montanum* genome are
617 inferred, assuming a phylogenetic placement of gnetophytes as sister to, or within
618 conifers. The predicted features shared by respective lineages are marked by coloured
619 circles. Likely whole genome duplication (WGD) events (red stars) and a putative
620 WGD event (grey star) are shown.

621

622

Legend:
- **Gnetum** (2-4 Gb)
- **Ephedra** (8-18 Gb)
- **Welwitschia** (8 Gb)
- **Picea glauca** (22 Gb)
- **Picea abies** (20 Gb)
- **Pinus taeda** (22 Gb)

**a**

PC2 ( 13.28% )

PC1 ( 24.74% )

- ■ Angiosperms
- ■ Other gymnosperms
- ■ Gnetophytes
- ■ Non-seed plants

**b**

Atr Ath Sly Vvi Peq Ptr Gma Zma Mac Osa Ppa Smo Gmo Pab Pta Gbi

number

- -2.0
- -1.5
- -1.0
- -0.5
- 0.0
- 0.5
- 1.0
- 1.5
- 2.0

Top 10% of PC1 & PC2 (454 items)

**Assemblage of diverse repeats**

**Intron evolution (LINE insertions involved)**

**High levels of TE amplification and deletion**

**Expansion of genes involved in flower and/or organ development**

**Seed plants ancestor**

Dicots

Monocots

*Amborella*

*G. montanum*

Conifers

Ginkgo

Cycads

Shared features

Limited recent TE amplifications

Elevated levels of LTR-RT excision

Similar intron morphology

Increased rates of evolution

Reduction of intron size

Unusual functional domain organization

Expansion/contraction of specific gene families (e.g. CesA, LEA)

**Large genome size**

**Abundant LTR-RTs**

**Low level of TE deletions**

**Long introns (LTR-RT insertion)**

**Gene expansion associated with defense**