

# A genome for gnetophytes and early evolution of seed plants

Tao Wan<sup>1,2,18</sup>, Zhi-Ming Liu<sup>3,18</sup>, Ling-Fei Li<sup>1,18</sup>, Andrew R. Leitch<sup>4,18</sup>, Ilia J. Leitch<sup>5,18</sup>, Rolf Lohaus<sup>6,7</sup>, Zhong-Jian Liu<sup>8,18</sup>, Hai-Ping Xin<sup>2,9</sup>, Yan-Bing Gong<sup>10</sup>, Yang Liu<sup>1</sup>, Wen-Cai Wang<sup>4</sup>, Ling-Yun Chen<sup>2,11</sup>, Yong Yang<sup>12</sup>, Laura J. Kelly<sup>4</sup>, Ji Yang<sup>13</sup>, Jin-Ling Huang<sup>14,15</sup>, Zhen Li<sup>6,7</sup>, Ping Liu<sup>1</sup>, Li Zhang<sup>1</sup>, Hong-Mei Liu<sup>1</sup>, Hui Wang<sup>1</sup>, Shu-Han Deng<sup>3</sup>, Meng Liu<sup>3</sup>, Ji Li<sup>3</sup>, Lu Ma<sup>4</sup>, Yan Liu<sup>3</sup>, Yang Lei<sup>3</sup>, Wei Xu<sup>3</sup>, Ling-Qing Wu<sup>3</sup>, Fan Liu<sup>2</sup>, Qian Ma<sup>10</sup>, Xin-Ran Yu<sup>3</sup>, Zhi Jiang<sup>3</sup>, Guo-Qiang Zhang<sup>8</sup>, Shao-Hua Li<sup>16</sup>, Rui-Qiang Li<sup>3</sup>, Shou-Zhou Zhang<sup>1</sup>, Qing-Feng Wang<sup>10,2,11\*</sup>, Yves Van de Peer<sup>6,7,17\*</sup>, Jin-Bo Zhang<sup>10,3\*</sup> and Xiao-Ming Wang<sup>1\*</sup>

**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 among seed plants. To increase our understanding of the evolution of gnetophytes, and their relation to other 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 (for example, cellulose synthase) or contractions (for example, Late Embryogenesis Abundant protein), 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.**

The seed plants today are represented by five distinct lineages: the species-rich angiosperms (flowering plants, approximately 352,000 species) and four gymnosperm lineages (which together comprise approximately 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 approximately 360 million years ago (Ma)) that considerably greater seed plant diversity existed in the past<sup>1</sup>. Nevertheless, widespread extinctions among many gymnosperm lineages mean that today's gymnosperms are only a relic of their former diversity, and this has presented a major challenge for reconstructing evolutionary relationships between the extant lineages<sup>2</sup>. Probably the most

controversial outstanding question in plant evolution is the phylogenetic position of gnetophytes<sup>3</sup> (comprising the genera *Gnetum*, *Welwitschia* and *Ephedra*, Fig. 1) in relation to the other seed plant lineages. Apparent morphological similarities with angiosperms, such as vessel-like water-conducting cells, double fertilization and leaf morphologies with reticulate venation, have historically led to the proposition that gnetophytes form a group that is sister to angiosperms (termed the 'Anthophyte hypothesis')<sup>4,5</sup>. That hypothesis has, however, largely been rejected by molecular phylogenetic data and a deeper understanding of the developmental pathways that lead to similar morphological features. Nevertheless, the use of molecular data has also been problematic in inferring the exact

<sup>1</sup>Key Laboratory of Southern Subtropical Plant Diversity, Fairy Lake Botanical Garden, Shenzhen & Chinese Academy of Science, Shenzhen, China.

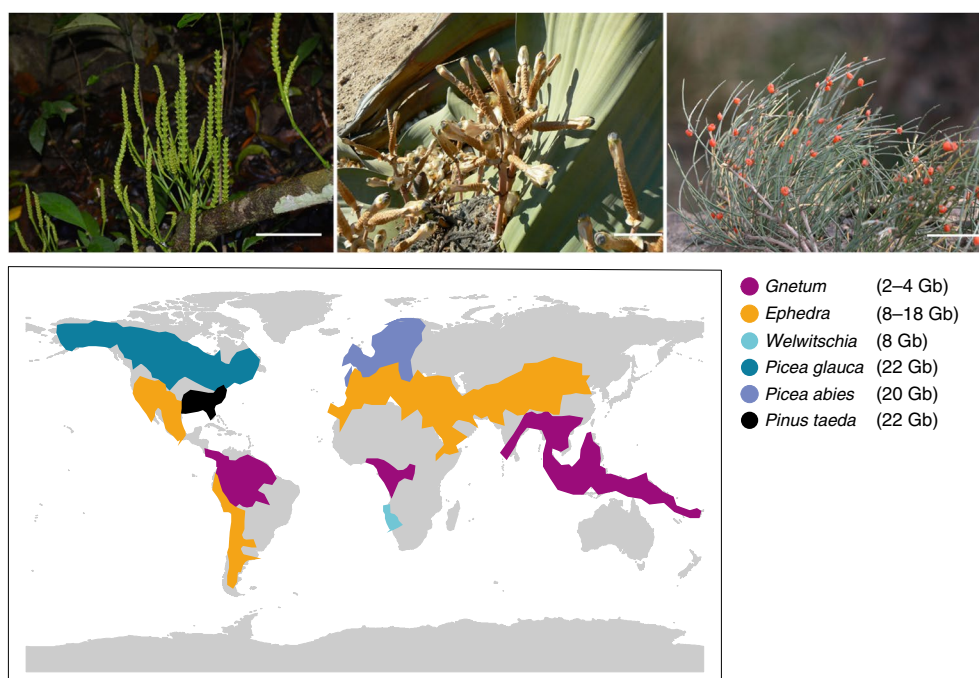
<sup>2</sup>Sino-Africa Joint Research Centre, Chinese Academy of Science, Wuhan, China. <sup>3</sup>Novogene Bioinformatics Institute, Beijing, China. <sup>4</sup>School of Biological and Chemical Sciences, Queen Mary University of London, London, UK. <sup>5</sup>Jodrell Laboratory, Royal Botanic Gardens, Kew, UK. <sup>6</sup>Department of Plant Biotechnology and Bioinformatics, Ghent University, Ghent, Belgium. <sup>7</sup>Centre for Plant Systems Biology, VIB, Ghent, Belgium. <sup>8</sup>Shenzhen Key Laboratory for Orchid Conservation and Utilization, National Orchid Conservation Centre of China and Orchid Conservation and Research Centre, Shenzhen, China.

<sup>9</sup>Key Laboratory of Plant Germplasm Enhancement and Specialty Agriculture, Wuhan Botanical Garden, Chinese Academy of Sciences, Wuhan, China. <sup>10</sup>State Key Laboratory of Hybrid Rice, College of Life Sciences, Wuhan University, Wuhan, China. <sup>11</sup>Key Laboratory of Aquatic Botany and Watershed Ecology, Wuhan Botanical Garden, Chinese Academy of Sciences, Wuhan, China. <sup>12</sup>State Key Laboratory of Systematic and Evolutionary Botany, Institute of Botany, Chinese Academy of Sciences, Beijing, China. <sup>13</sup>Education Key Laboratory for Biodiversity Science and Ecological Engineering, Fudan University, Shanghai, China. <sup>14</sup>Institute of Plant Stress Biology, State Key Laboratory of Cotton Biology, Henan University, Kaifeng, China.

<sup>15</sup>Department of Biology, East Carolina University, Greenville NC, USA. <sup>16</sup>Beijing Key Laboratory of Grape Sciences and Enology, Institute of Botany, Chinese Academy of Sciences, Beijing, China. <sup>17</sup>Genomics Research Institute, University of Pretoria, Pretoria, South Africa. <sup>18</sup>These authors contributed equally: Tao Wan, Zhi-Ming Liu, Ling-Fei Li, Andrew R. Leitch, Ilia J. Leitch and Zhong-Jian Liu. \*e-mail: [qfwang@wbgcas.cn](mailto:qfwang@wbgcas.cn); [yves.vandeppeer@psb.vib-ugent.be](mailto:yves.vandeppeer@psb.vib-ugent.be); [zhangjinbo@novogene.com](mailto:zhangjinbo@novogene.com); [719921868@qq.com](mailto:719921868@qq.com)

<sup>19</sup>Key Laboratory of Plant Germplasm Enhancement and Specialty Agriculture, Wuhan Botanical Garden, Chinese Academy of Sciences, Wuhan, China. <sup>20</sup>State Key Laboratory of Hybrid Rice, College of Life Sciences, Wuhan University, Wuhan, China. <sup>21</sup>Key Laboratory of Aquatic Botany and Watershed Ecology, Wuhan Botanical Garden, Chinese Academy of Sciences, Wuhan, China. <sup>22</sup>State Key Laboratory of Systematic and Evolutionary Botany, Institute of Botany, Chinese Academy of Sciences, Beijing, China. <sup>23</sup>Education Key Laboratory for Biodiversity Science and Ecological Engineering, Fudan University, Shanghai, China. <sup>24</sup>Institute of Plant Stress Biology, State Key Laboratory of Cotton Biology, Henan University, Kaifeng, China.

<sup>25</sup>Department of Biology, East Carolina University, Greenville NC, USA. <sup>26</sup>Beijing Key Laboratory of Grape Sciences and Enology, Institute of Botany, Chinese Academy of Sciences, Beijing, China. <sup>27</sup>Genomics Research Institute, University of Pretoria, Pretoria, South Africa. <sup>28</sup>These authors contributed equally: Tao Wan, Zhi-Ming Liu, Ling-Fei Li, Andrew R. Leitch, Ilia J. Leitch and Zhong-Jian Liu. \*e-mail: [qfwang@wbgcas.cn](mailto:qfwang@wbgcas.cn); [yves.vandeppeer@psb.vib-ugent.be](mailto:yves.vandeppeer@psb.vib-ugent.be); [zhangjinbo@novogene.com](mailto:zhangjinbo@novogene.com); [719921868@qq.com](mailto:719921868@qq.com)



**Fig. 1 | Morphological variation and geographical distribution of gnetophytes and some other gymnosperms.** Top, left to right, female cones of *G. montanum*, male cones of *W. mirabilis* and female cones of *E. equisetina*. Scale bars, 5 cm. Bottom, pantropical distribution of the three gnetophyte genera, compared with three conifer species that are most abundant at higher latitudes and altitudes. The range of genome sizes (1C-values) found in the three genera comprising gnetophytes and the three conifer species are also shown (data taken from <http://data.kew.org/cvalues/> and unpublished data).

phylogenetic position of gnetophytes, with topologies differing depending on the type of sequence data (for example, plastid versus nuclear genes, nucleotide versus amino acid data) and analytical approach used (for example, maximum parsimony, maximum likelihood, Bayesian, multispecies coalescent based methods)<sup>6–8</sup>. Consequently, several possible hypotheses have been put forward that place gnetophytes as sister to (1) Pinaceae (‘Gnepine’ hypothesis); (2) cupressophytes (‘Gnecup’ hypothesis); (3) all conifers (‘Gnetifer’ hypothesis); (4) all other gymnosperms; or (5) all seed plants<sup>9</sup>. Currently, the emerging consensus, based on both older and more recent studies, and recently released data from the IKP initiative (see <https://sites.google.com/a/ualberta.ca/onekpl/>, and Wickett et al.<sup>8</sup>), indicates that gnetophytes are sister to, or within, the conifers.

So far, the availability of whole genome sequences for gymnosperms has been limited to conifers (specifically to Pinaceae)<sup>10–13</sup> and *G. biloba*<sup>14</sup>, with no whole genome assemblies available for the two remaining major seed plant lineages—cycads and gnetophytes. This deficiency, together with the conflicting phylogenetic evidence for relationships among these groups, is impeding our understanding of genome evolution across all seed plants. Here, we present a high-quality draft genome of *Gnetum montanum*, the first for gnetophytes. The availability of this genome, as well as survey sequence data and transcriptome data from other vascular plants (including novel data from gnetophytes *Ephedra* and *Welwitschia*), enables us to compare genomic characters with *G. biloba*, conifers, angiosperms and non-seed plants. Comparisons within gymnosperms, and between gymnosperms and angiosperms, highlight the unique nature of the *Gnetum* genome, providing new insights into patterns of genome divergence across seed plants.

### Genome assembly and annotation

The genome of *G. montanum* ( $2n = 44$ ) is small compared with other gymnosperms (flow cytometry, 4.2 Gb/1C; k-mer analysis, 4.11 Gb), and is highly heterozygous and rich in repeats (Supplementary

Fig. 1a–c and Supplementary Information). To overcome problems caused by repeats and heterozygosity, we generated deep coverage (~302×, Supplementary Table 1) Illumina sequence data and applied a novel genome assembly strategy (Supplementary Information and Supplementary Fig. 2) to assemble 4.07 Gb of sequence (contig N50 size = 25.02 kb, scaffold N50 size = 475.17 kb, Supplementary Table 2), to which >99% of genome reads, >90% expressed sequence tags (ESTs) and >99% of bacterial artificial chromosomes (BACs) were mapped (Supplementary Fig. 1d,e, Supplementary Table 3 and Note 3).

A total of 27,491 protein-coding genes were predicted from this assembly (Supplementary Table 4 and Supplementary Information), 97% of which were supported by orthology (>50% coverage of a high-scoring segment pair, Supplementary Fig. 3a) with existing protein sequences and/or RNA-seq data from multiple tissues (Supplementary Table 5). A BUSCO (Benchmarking Universal Single-Copy Orthologs) analysis to assess the quality of the genome and annotation completeness suggested that 81% of the genes have been recovered (Supplementary Table 6). Unlike conifer genomes, which contain numerous pseudogenes<sup>15</sup> (for example, 8,328 in *Picea abies*, 13,550 in *Pinus taeda*), many fewer were found in the *G. montanum* genome (3,122, Supplementary Information). The read depth distribution across genic regions (Supplementary Fig. 3b) suggested little sequence redundancy caused by heterozygosity (see Supplementary Fig. 3c for further confirmation of gene assembly quality).

### Repetitive sequence dynamics

Repetitive sequences have been shown to account for the major component of all gymnosperm genomes that have been sequenced to date<sup>11–14</sup>, with diverse and ancient transposable elements (TEs), especially LTR retrotransposons (LTR-RTs), being particularly prevalent. Overall, the repetitive element content of *G. montanum* was also high (85.9%) and dominated by LTR-RTs (especially *gypsy*-like

elements), which constituted 77.4% of the genome (Supplementary Table 8 and Supplementary Information). The genome assembly of *G. montanum* is likely to be sufficient to represent most of the LTR-RTs, since their length is typically around 25 kb<sup>16</sup>, and 90% of the scaffolds are larger than 34 kb. Phylogenetic reconstructions of the reverse transcriptase domains of LTR-RTs in *G. montanum* and *P. taeda* revealed that most of the *gypsy*- and *copla*-like elements in *G. montanum* were restricted to just a few clades, representing only a small minority of the diversity encountered in *P. taeda* (Supplementary Fig. 4 and Supplementary Information).

Comparative analyses of repeats identified by RepeatExplorer using survey sequence data from multiple gnetophytes (*G. montanum*, *Gnetum gnemon*, *Welwitschia mirabilis* and *Ephedra altissima*) and *P. taeda* revealed substantial differences in the abundance of the major repeat classes (Supplementary Fig. 5a, Supplementary Table 9 and Supplementary Information). Further, the majority of individual repeat types (repeat clusters in RepeatExplorer) were shown to be species specific (containing Illumina reads from just one species, data not shown). The species-specific nature of the repeat profiles probably reflects the long estimated divergence times between species (for example, the two *Gnetum* species are likely to have diverged between approximately 25 Ma and 75 Ma)<sup>17,18</sup>.

Previously, it was reported from conifers and *G. biloba* that LTR-RTs have accumulated steadily over the last approximately 25 Ma, especially between 16 and 24 Ma, a process contributing to their large genome sizes<sup>11,12,14</sup>. This interpretation is consistent with the data here (Supplementary Table 10), which show that most LTR-RTs in conifers are intact (solo LTR/intact LTR ratio ranged from 0.16:1 to 0.72:1, Supplementary Table 10). It is notable that the solo LTR/intact LTR ratio was substantially higher in *G. montanum* (~1.94:1), which together with its small genome and similar profile of accumulation (Supplementary Fig. 5b) suggest higher frequencies of LTR-RT elimination than amplification compared with *G. biloba* and conifers.

Most angiosperm genomes analysed to date have far fewer ancient repeats and less divergent LTR-RT subsets than conifers and *G. biloba*, presumably because of more efficient elimination and replacement processes operating within these angiosperm genomes<sup>19</sup> (for example, in *Oryza sativa* the half-life of LTR-RTs is estimated to be less than five million years<sup>20</sup>, leading to 'genome turnover'<sup>21</sup>). However, an exception to this pattern has been observed in *Amborella trichopoda*. The genome of this species is considered to have retained many features that were likely to have been present in the ancestral angiosperm genome<sup>22</sup>. It is notable that its repeat content<sup>13</sup> and lower abundance of intact LTR-RTs (solo LTR/intact LTR ratio = 2.43/1.0; Supplementary Table 10) is similar to that observed in *G. montanum*. These observations suggest that neither *A. trichopoda* nor *G. montanum* genomes have experienced recent, extensive (retro) transposon activity, although they continue to eliminate repetitive sequences. Both these species seem to differ from conifers and *G. biloba* with respect to the dynamics of repeat accumulation<sup>11,12,14</sup>, and from other angiosperms in terms of the levels of repeat amplification/removal.

### Intron morphologies

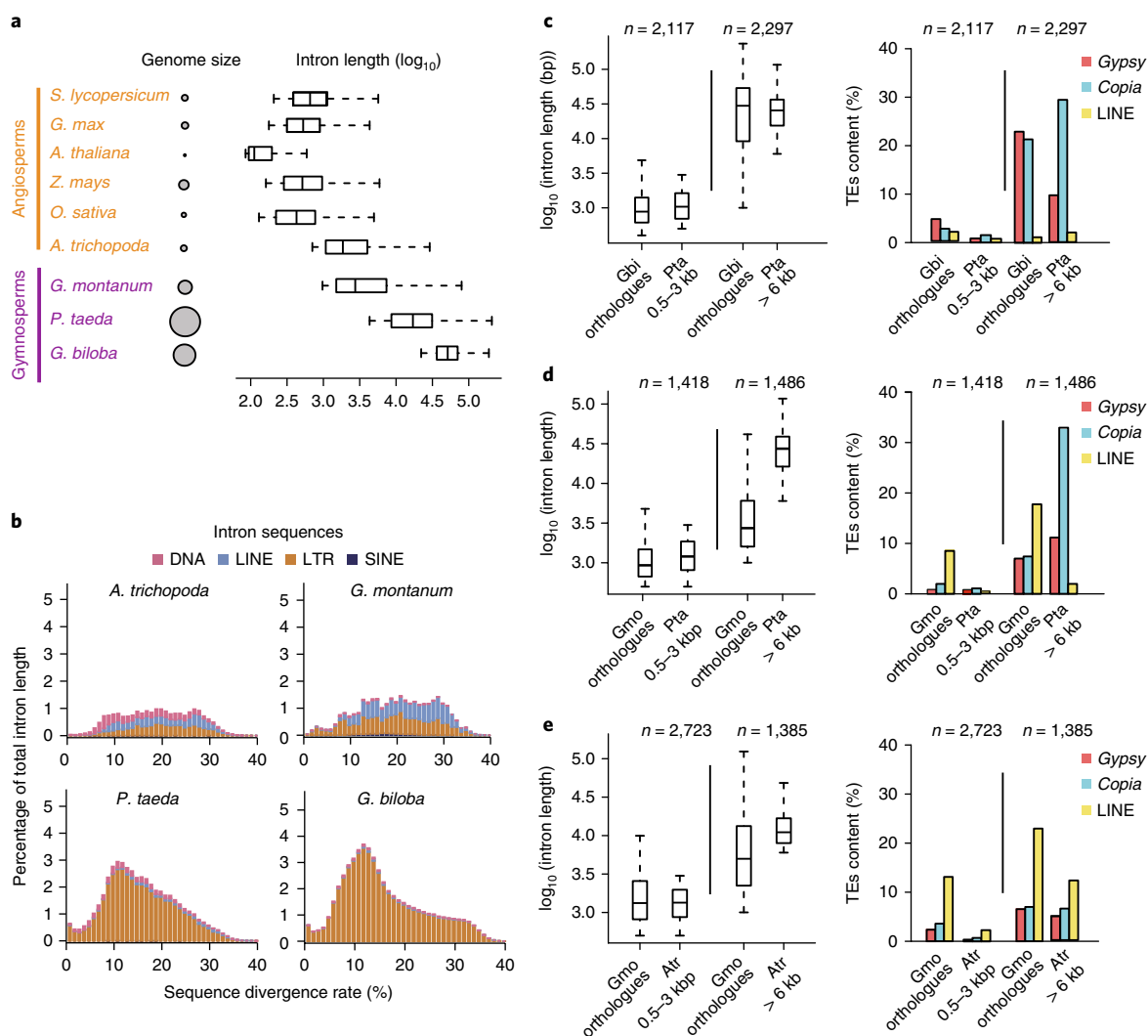
Although intron size has been positively correlated with genome size across eukaryotes as a whole<sup>23</sup>, this trend does not translate well across broad and some narrow taxonomic distances in seed plants (Fig. 2a). Previous studies of *G. biloba*<sup>14</sup> and conifers<sup>11,12</sup> have reported larger introns than angiosperms, probably arising from the long-term, steady amplification of LTR-RTs (Fig. 2b), as also observed here, where LTR-RTs account for 51% and 59% of the large intron sequences in *P. taeda* and *G. biloba*, respectively (Fig. 2a and Supplementary Table 12). The evolution of these large introns may have arisen from similar repeat accumulation processes that are operating across the genome as a whole.

When comparing these observations with introns of *G. montanum*, it is apparent that their introns are substantially smaller (minimum, mean and maximum intron lengths) than those of *P. taeda* and *G. biloba* (Fig. 2a, see also the statistics test in Supplementary Table 11). In addition, the repeat composition of *G. montanum*'s introns is dominated by both long interspersed nuclear elements (LINEs) and LTR-RTs, rather than predominantly LTR-RTs, as in conifers and *G. biloba* (Fig. 2b and Supplementary Table 12). The correlation between smaller intron sizes and smaller genome size in *G. montanum* compared with conifers and *G. biloba* may reflect the repeat dynamic processes operating across its genome as a whole. In contrast, the variable length distributions of introns in angiosperms suggest that the evolution of repeats in their introns do not necessarily reflect the repeat dynamics observed across the rest of their genomes<sup>24</sup>. In the highly dynamic repetitive genome of *Zea mays*, the profile of repeats across the genome<sup>25</sup> and within the whole intron set (Supplementary Fig. 6a) both suggest many recent insertions. However, in *A. trichopoda*, the intron sizes are larger overall, and the genome size smaller than in *Z. mays* (Fig. 2a,b). In addition, an analysis of introns in *A. trichopoda* and *G. montanum* highlighted a closer similarity to each other (in terms of length distributions, repeat composition and divergence) than either species has to conifers and *G. biloba*, despite a 4.8-fold difference in their genome sizes (Fig. 2a,b and Supplementary Table 12).

Previous comparisons of orthologous introns have led to the suggestion that the expansion of introns occurred early in the evolutionary history of conifers<sup>12</sup>. Comparisons of orthologous introns (with identical adjacent exons) between *P. taeda* and *G. biloba* showed that introns identified as being long (>6 kb) in *P. taeda* were also typically long in their orthologues in *G. biloba*, containing, in both cases, abundant LTR-RTs (both *gypsy*- and *copla*-like elements, Fig. 2c). These features were likely to have been present in their most recent common ancestor (MRCA). Using similar approaches to analyse the length and repeat content of 4,348 orthologous introns of *G. montanum* shared with *P. taeda* (Supplementary Information) highlighted notable differences. The length of exons remained similar, but a substantial fraction of orthologous genes had longer introns in *P. taeda* (Supplementary Fig. 6b). The introns identified as 'short' in *P. taeda* comprised approximately 4% repeats, rising to approximately 56% in 'long' introns, largely through the accumulation of LTR-RTs (especially *copla* elements) (Fig. 2d and Supplementary Table 13). In contrast, introns in *G. montanum* that are orthologous to the 'long' introns of *P. taeda* (36% of introns analysed) showed high proportions of LINEs. As with comparisons of all introns, pairwise comparisons of orthologous introns in *G. montanum* and *A. trichopoda* again showed some similarities in their introns, with both species having abundant LINEs (Fig. 2e). Collectively, these data reveal a different repeat dynamic within introns of *G. montanum* compared with the other gymnosperms.

### ('Lack of') Whole genome duplication

All angiosperms are reported to have undergone at least one round of ancient whole genome duplication (WGD), and in many lineages WGDs are recurrent and ongoing<sup>26</sup>. In addition, a WGD event has been proposed at the base of all seed plants approximately 341 Ma (*zeta* WGD<sup>27</sup>), although the underlying evidence for these two ancient WGD events has been recently questioned<sup>28</sup>. In gymnosperms, WGDs have been reported for conifers, *G. biloba* and cycads (a likely shared WGD)<sup>14,29,30</sup>. Although recent polyploidy seems common in extant *Ephedra*<sup>31</sup>, evidence for ancient WGDs in gnetophytes is missing (Supplementary Information and Supplementary Fig. 7), except for a WGD in *Welwitschia* which is likely to have occurred after the divergence of its lineage from that leading to *Ephedra* (Supplementary Fig. 7)<sup>29</sup>. If indeed the ancient *zeta* WGD is shared by all seed plants, the absence of evidence for this event in gnetophytes is best explained by their faster rates of gene evolution



**Fig. 2 | Comparative analysis of seed plant intron morphologies.** **a**, Intron length distributions and genome sizes (1C-values, depicted by the relative circle size) are shown for nine representative seed plants. *S. lycopersicum*, *Solanum lycopersicum*; *G. max*, *Glycine max*. **b**, Distribution of sequence divergence for four types of TEs in introns of *A. trichopoda*, *G. montanum*, *P. taeda* and *G. biloba*. The data show that TEs in *G. montanum* and *A. trichopoda* are more diverse than in *P. taeda* and *G. biloba*. The last two species also show a peak at around 10% sequence divergence probably reflecting a pulse of LTR-RT expansions. **c–e**, Comparison of orthologous introns between *P. taeda* (Pta) versus *G. biloba* (Gbi) (**c**), *P. taeda* versus *G. montanum* (Gmo) (**d**), and *G. montanum* versus *A. trichopoda* (Atr) (**e**). Two orthologous intron sets that differed more than twofold in length were examined: ‘short’ introns, 0.5–3 kb; and ‘long’ introns, 6 kb. Orthologous introns that were long in one species were also found to be long in the other species of the pair. Analysis of the TEs in orthologous introns showed the long introns of *G. montanum* and *A. trichopoda* carried a high proportion of LINES, contributing to intron expansion. In contrast, *gypsy* and *copia* LTR-RT elements contributed most to intron expansion in *P. taeda* and *G. biloba*.

than other gymnosperms<sup>32,33</sup>, erasing all evidence of this more than 300 million year old event (Supplementary Information and Supplementary Fig. 7).

### Organization of functional protein domains

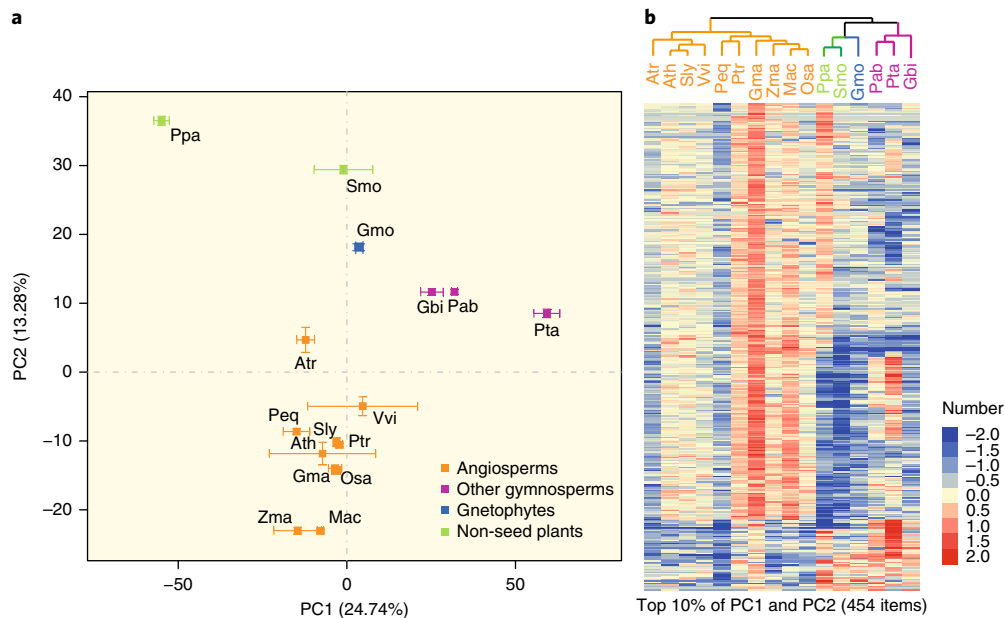
To characterize the patterns of functional diversification in gene domains across land plants, we used principal component analysis (PCA) to analyse the number of pfam domains (conserved protein domains) in multiple species (Supplementary Information and Supplementary Table 13). Our approach showed that angiosperms formed a discrete cluster that was separate from the gymnosperms (Fig. 3a), with *G. montanum* being an outlier. Indeed, heatmaps compiled from the pfam data that contributed most (top 10%) to PCA1 and PCA2 showed that *G. montanum* formed a clade with the lycophyte *Selaginella moellendorffii* and the moss *Physcomitrella patens*

(Fig. 3b), but the non-gnetophyte gymnosperms formed a separate clade (Fig. 3b).

Given the distinct distributions of *G. montanum*, non-gnetophyte gymnosperms and angiosperms in the PCA analysis, the data suggest that significant functional diversification of the conserved protein domains has occurred since these major lineages split. It may be surprising given the long divergence times (approximately 300 Ma)<sup>2</sup>, that *G. biloba* and conifers retain similar conserved domain organizations (with similar eigenvector values). This could reflect their relatively low substitution rates (on average seven times lower) compared with angiosperms<sup>33</sup>.

An analysis of the pfam domain expansions that contributed most to the PCA1 and PCA2 distributions among angiosperms (except *A. trichopoda*) included genes associated with flower and organ development (Supplementary Table 15). In contrast, non-gnetophyte





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

gymnosperms showed large-scale specific expansions of pfam domains in genes associated with defence and secondary metabolism, as previously suggested (Supplementary Table 16)<sup>10,11</sup>. The clustering of *G. montanum* with non-seed plants in the heatmap (Fig. 3b) was a surprise, and may indicate the approach has identified proteins that have diverged very little since the MRCA of seed plants. Nevertheless, such an explanation is at odds with the hypothesis that the genes of gnetophytes have diverged rapidly, given their comparatively high substitution rate compared with other gymnosperms<sup>33</sup>.

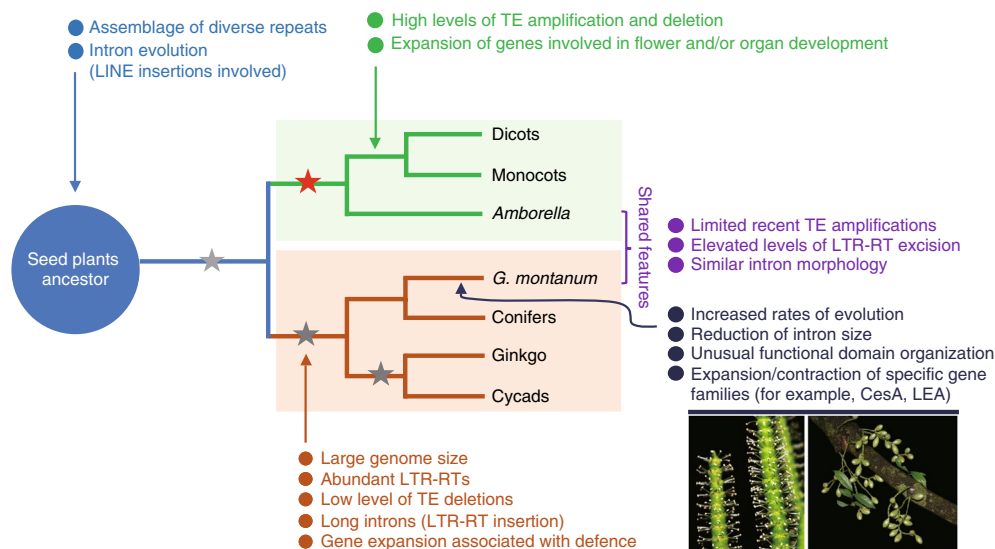
### Growth form (shrubs and lianas) and leaf morphology

Gnetophytes differ from other extant gymnosperms in growth form, with the unusual and distinct form of *Welwitschia*, the shrub habit of *Ephedra* and the shrub and liana habit and specialized leaf morphologies of *Gnetum*<sup>34</sup>. Cellulose synthase (*CesA*) and cellulose synthase-like (*Csl*) genes are considered to play a role in influencing the biomechanical properties of the cell<sup>35</sup>, hence potentially the distinctive growth forms of gnetophytes are associated with the divergence of these genes. To explore this hypothesis, *CesA* and *Csl* family members were examined in *G. montanum* and compared with those in other seed plants. The total number of *CesA* and *Csl* family members ranged about threefold among the seed plants analysed (*P. abies*, *P. taeda*, *A. trichopoda*, *A. thaliana* and *O. sativa*). However, only *G. montanum* showed a large expansion of the *CslB/H* gene subfamily (to 20 genes, Supplementary Table 17), involving tandem duplications (Supplementary Fig. 9), and accounting for two-thirds of its total *Csl* gene repertoire. Furthermore, transcriptome analysis showed that these *CslB/H* genes were differentially expressed in leaves, stems and roots of *G. montanum*, supporting an association with distinct growth forms and leaf morphologies (Supplementary Fig. 9). In contrast, all other species analysed, including *Welwitschia* and *Ephedra*, were seen to have only one to six *CslB/H* genes (at least based on transcriptome analysis) (Supplementary Information, Supplementary Table 16 and Supplementary Fig. 8).

Another gene family associated with leaf morphology and development is the *WOX* (*WUSCHEL-related homeobox*) family<sup>36</sup>. Recent studies have shown that the conserved family members *WOX3* and *WOX4*, which play a role in leaf development, show diffuse *WOX3* expression at the leaf bases of *Arabidopsis* and *Gnetum*, with such patterns being associated with the distinctive reticulate venation observed in their leaves<sup>37</sup>. Two unusual paralogues, *GgWOXX* and *GgWOXY*, were previously reported to occur only in gnetophytes<sup>37</sup>, and this is confirmed here in phylogenetic reconstructions of gene family members (Supplementary Information and Supplementary Fig. 10). These paralogues are unlikely to have arisen by *Gnetum*-specific gene amplifications, as this would group them with other *Gnetum* paralogues. Alternatively, these genes may correspond to ancestral seed plant sequences that have been lost in other plant lineages. Potentially the different patterns of gene loss, retention and amplification compared with other gymnosperms may be associated with their distinctive growth forms.

### Vessels

The presence of vessel-like water-conducting cells, morphologically distinct from tracheids, is another feature that sets gnetophytes apart from other gymnosperms. However, there has been long-standing debate whether gnetophyte 'vessels' are homologous to the 'vessels' of angiosperms. In angiosperms, VASCULAR-RELATED NAC-DOMAIN (VND) proteins *VND1-7* are members of the NAC domain class of transcription factors, *VND7* being a master regulator of vessel formation in *A. thaliana*<sup>38</sup>, and *VND1-6* being upstream regulators of *VND7*<sup>39</sup>. Although five NAC domain genes were identified in the genome of *G. montanum*, no orthologues of *VND7* or *VND1-3* in the sister clade were identified, consistent with previous analyses of other gymnosperms<sup>12</sup>, and suggesting that these proteins are restricted to angiosperms (Supplementary Fig. 11). Nevertheless, *Gnetum* does share the *VND4-6* clade with angiosperms and other gymnosperms. Furthermore, *A. trichopoda*, which lacks angiosperm vessels, also lacks orthologues of *VND1-3*,



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

but it does have *VND7* (Supplementary Fig. 11), indicating that the ability to form vessels may have occurred after angiosperms diverged. Taken together, these data suggest a greater dependency of vessel development on *VND1-3* than is apparent from experiments on *A. thaliana*. The most parsimonious explanation of our data is that angiosperm vessel formation requires genes from the *VND7* clade (and potentially its sister clade *VND1-3*), and that gymnosperms, including gnetophytes, which lack sequences from both these clades cannot form structures that are homologous to angiosperm vessels. Such an interpretation supports Carlquist's<sup>40</sup> morphological interpretations of vessels. It is therefore most likely that different molecular mechanisms underpin the origin and development of vessels in *Gnetum* and angiosperms. Indeed, these new molecular data support the hypothesis based on morphological studies that *Gnetum* vessels are actually more closely related to conifer tracheids than angiosperm vessels and that vessels in the two groups are convergent characters<sup>10</sup>.

### Water stress

Extant species of *Gnetum* are unusual among gymnosperms in being restricted to warm, mesic habitats<sup>41</sup>; this contrasts to conifers that are adapted to cold and water-stressed environments. An analysis of genes involved in water and cold stress revealed some substantial differences between conifers and *Gnetum*. The late embryogenesis abundant protein (LEA) gene family encodes crucial proteins that are involved in protecting plants from desiccation or osmotic stresses associated with low temperature<sup>42,43</sup>. An analysis of LEA family members suggests that some members have been reduced in number in *Gnetum* or expanded in conifers (for example, LEA-3), or lost completely in *Gnetum* (LEA-4, 5, 6). In addition, dehydrins, which play a role in the response to cold/drought<sup>44</sup>, had only two members in *G. montanum*, compared with 38 in *P. abies*, 28 in *P. taeda* and 3–15 in angiosperms (Supplementary Table 19). Further analysis of the *G. montanum* genome also revealed relatively few gene family members of the AP2 domain containing protein families, which are involved in the cold stress response<sup>45,46</sup>, and glutathione peroxidase and glutathione S-transferase families, involved in the oxidant stress response<sup>47,48</sup>. Taken together, these data appear consistent with the hypothesis that the ecological shift to a warm, wet forest habitat is

associated with a relaxation of selection pressure on genes associated with water stress and low temperature.

### Conclusion

Here, we have described the assembly, annotation and comparative analysis of the first gnetophyte genome, namely that of *G. montanum*. Its genome is particularly enigmatic given a phylogenetic position within or sister to conifers. It also carries genomic peculiarities that may reflect its morphological and ecological uniqueness amongst gymnosperms. Comparisons of these genome features with the genomes of conifers and *G. biloba* provide opportunities to predict the nature and direction of genomic change accompanying the evolution of the lineage leading to *Gnetum* (Fig. 4). Assuming that gnetophytes do indeed form a clade that is sister to, or within, the conifers, the following genomic features can be predicted to have been present in the MRCA of the gymnosperms, as observed in *G. biloba*<sup>14</sup> and conifers<sup>11,12</sup>: (1) a large genome size (1C > 10 Gb) comprised predominantly of a heterogeneous set of large numbers of LTR-RTs associated with low levels of repeat deletion<sup>14</sup>; (2) long introns predominantly shaped by insertions of LTR-RTs (*gypsy* and *copia* elements); (3) pfam domains that show a profile distinct from angiosperms. If this is so, and assuming a common ancestry of gnetophytes and conifers, these genomic characters, or their signatures, have subsequently been lost or diverged considerably in the lineage leading to *Gnetum*. This most likely involved the following genomic processes: (1) genome downsizing, leading to the relatively (for a gymnosperm) small genomes of *Gnetum* species (1C = 2.25–4.11 Gb). This is supported by the high ratio of solo LTR/intact LTR-RTs observed in the genome of *Gnetum* compared with conifers, and is indicative of the activity of recombination-based processes, which can eliminate DNA from the genome. Similar processes leading to genome downsizing have also been reported in many angiosperms, resulting in small genomes despite the occurrence of multiple rounds of polyploidy detected in many lineages<sup>49</sup>; (2) reduction in the size of introns in *G. montanum* and a replacement of many of the LTR-RTs repeats with LINES to give rise to introns that are more similar to those of, for instance, *A. trichopoda* than to other gymnosperms; (3) elevated rates of sequence divergence causing the erosion of a hypothesised shared seed-plant WGD event and leading to a pattern of pfam domains, which is distinct from the

remaining gymnosperms; (4) expansion and contraction of specific gene families associated with adaptation to new ecologies.

## Methods

The sequenced *G. montanum* is a single mature female individual growing naturally in Fairy Lake Botanical Garden, Shenzhen, China. Genome sequences were generated using an Illumina platform and assembled with a novel hierarchical assembly strategy. Gene annotations were determined by integrating results from both de novo prediction approaches and alignment-based methods based on orthology and transcriptomic data. RNA-seq was performed using an Illumina platform. All methods and bioinformatic analyses are detailed in the Supplementary Information.

**Life Sciences Reporting Summary.** Further information on experimental design is available in the Life Sciences Reporting Summary.

**Data availability.** The *G. montanum* genome project has been deposited at the NCBI under the BioProject number PRJNA339497. The whole genome sequencing data were deposited in the Sequence Read Archive (SRA) database under the accession number SRX2052734, SRX2098865, SRX2099144, SRX2114825, SRX2114827, SRX2134147, SRX2134160, SRX2134177, SRX2134180, SRX2134596 and SRX2134624. The *G. montanum* assemblies, gene sequences and annotation data are also available at the DRYAD website. The data or related program scripts that support the findings of this study are available from the corresponding author upon request.

Received: 2 June 2017; Accepted: 27 December 2017;

Published online: 29 January 2018

## References

- Rothwell, G. W. & Scheckler, S. E. in *Origin and Evolution of Gymnosperms* (ed. Beck, C. B.) 85–134 (Columbia University Press, New York, 1988).
- Lu, Y., Ran, J. H., Guo, D. M., Yang, Z. Y. & Wang, X. Q. Phylogeny and divergence times of gymnosperms inferred from single-copy nuclear genes. *PLoS ONE* **9**, e107679 (2014).
- Doyle, J. A. Molecular and fossil evidence on the origin of angiosperms. *Annu. Rev. Earth Planet. Sci.* **40**, 301–326 (2012).
- Doyle, J. A. & Donoghue, M. J. Seed plant phylogeny and the origin of angiosperms: an experimental cladistic approach. *Bot. Rev.* **52**, 321–431 (1986).
- Crane, P. R. Phylogenetic analysis of seed plants and the origin of angiosperms. *Ann. Mo. Bot. Gard.* **72**, 716–793 (1985).
- Mathews, S. Phylogenetic relationships among seed plants: persistent questions and the limits of molecular data. *Am. J. Bot.* **96**, 228–236 (2009).
- Wang, X. Q. & Ran, J. H. Evolution and biogeography of gymnosperms. *Mol. Phylogenet. Evol.* **75**, 24–40 (2014).
- Wickett, N. J. et al. Phylotranscriptomic analysis of the origin and early diversification of land plants. *Proc. Natl Acad. Sci. USA* **111**, 4859–4868 (2014).
- Li, Z. et al. Single-copy genes as molecular markers for phylogenomic studies in seed plants. *Genome Biol. Evol.* **9**, 1130–1147 (2017).
- Warren, R. L. et al. Improved white spruce (*Picea glauca*) genome assemblies and annotation of large gene families of conifer terpenoid and phenolic defense metabolism. *Plant J.* **83**, 189–212 (2015).
- Neale, D. B. et al. Decoding the massive genome of loblolly pine using haploid DNA and novel assembly strategies. *Genome Biol.* **15**, R59 (2014).
- Nystedt, B. et al. The Norway spruce genome sequence and conifer genome evolution. *Nature* **497**, 579–584 (2013).
- Stevens, K. A. et al. Sequence of the sugar pine megagenome. *Genetics* **204**, 1613–1626 (2016).
- Guan, R. et al. Draft genome of the living fossil *Ginkgo biloba*. *GigaScience* **5**, 49 (2016).
- Garcia-Gil, M. R. Evolutionary aspects of functional and pseudogene members of the phytochrome gene family in Scots pine. *J. Mol. Evol.* **67**, 222–232 (2008).
- Wicker, T. et al. A unified classification system for eukaryotic transposable elements. *Nat. Rev. Genet.* **8**, 973–982 (2007).
- Won, H. & Renner, S. S. Dating dispersal and radiation in the gymnosperm *Gnetum* (Gnetales)—clock calibration when outgroup relationships are uncertain. *Syst. Biol.* **55**, 610–622 (2006).
- Hou, C., Humphreys, A. M., Thureborn, O. & Rydin, C. New insights into the evolutionary history of *Gnetum* (Gnetales). *Taxon* **64**, 239–253 (2015).
- Leitch, A. R. & Leitch, I. J. Ecological and genetic factors linked to contrasting genome dynamics in seed plants. *New Phytol.* **194**, 629–646 (2012).
- Ma, J., Devos, K. M. & Bennetzen, J. L. Analyses of LTR-retrotransposon structures reveal recent and rapid genomic DNA loss in rice. *Genome Res.* **14**, 860–869 (2004).
- Lim, K. Y. et al. Sequence of events leading to near-complete genome turnover in allopolyploid *Nicotiana* within five million years. *New Phytol.* **175**, 756–763 (2007).
- Amborella Genome Project. The *Amborella* genome and the evolution of flowering plants. *Science* **342**, 1241089 (2013).
- Vinogradov, A. E. Intron-genome size relationship on a large evolutionary scale. *J. Mol. Evol.* **49**, 376–384 (1999).
- Wendel, J. F. et al. Intron size and genome size in plants. *Mol. Biol. Evol.* **19**, 2346–2352 (2002).
- Schnable, P. S. et al. The B73 maize genome: complexity, diversity, and dynamics. *Science* **326**, 1112–1115 (2009).
- Van de Peer, Y., Mizrahi, E. & Marchal, K. The evolutionary significance of polyploidy. *Nat. Rev. Genet.* **18**, 411–424 (2017).
- Jiao, Y. et al. Ancestral polyploidy in seed plants and angiosperms. *Nature* **473**, 97–100 (2011).
- Ruprecht, C. et al. Revisiting ancestral polyploidy in plants. *Sci. Adv.* **3**, e1603195 (2017).
- Li, Z. et al. Early genome duplications in conifers and other seed plants. *Sci. Adv.* **1**, e1501084 (2015).
- Roodt, D. et al. Evidence for an ancient whole genome duplication in the cycad lineage. *PLoS ONE* **12**, e0184454 (2017).
- Wu, H. et al. A high frequency of allopolyploid speciation in the gymnospermous genus *Ephedra* and its possible association with some biological and ecological features. *Mol. Ecol.* **25**, 1192–1210 (2016).
- Hajibabaei, M., Xia, J. & Drouin, G. Seed plant phylogeny: gnetophytes are derived conifers and a sister group to Pinaceae. *Mol. Phylogenet. Evol.* **40**, 208–217 (2006).
- De La Torre, A. R., Li, Z., Van de Peer, Y. & Ingvarsson, P. K. Contrasting rates of molecular evolution and patterns of selection among gymnosperms and flowering plants. *Mol. Biol. Evol.* **34**, 1363–1377 (2017).
- Frohlich, M. W. & Chase, M. W. After a dozen years of progress the origin of angiosperms is still a great mystery. *Nature* **450**, 1184–1189 (2007).
- Popper, Z. A. et al. Evolution and diversity of plant cell walls: from algae to flowering plants. *Annu. Rev. Plant Biol.* **62**, 567–590 (2011).
- Nakata, M. et al. Roles of the middle domain-specific *WUSCHEL-RELATED HOMEBOX* genes in early development of leaves in *Arabidopsis*. *Plant Cell* **24**, 519–535 (2012).
- 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).
- 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).
- 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).
- Carlquist, S. Wood, bark and stem anatomy of New World species of *Gnetum*. *Bot. J. Linn. Soc.* **120**, 1–19 (1996).
- 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).
- Hundertmark, M. & Hincha, D. K. LEA (late embryogenesis abundant) proteins and their encoding genes in *Arabidopsis thaliana*. *BMC Genom.* **9**, 1–22 (2008).
- Gao, J. & Lan, T. Functional characterization of the late embryogenesis abundant (LEA) protein gene family from *Pinus tabulaeformis* (Pinaceae) in *Escherichia coli*. *Sci. Rep.* **6**, 19467 (2016).
- 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).
- Chinnusamy, V., Zhu, J. & Zhu, J. K. Cold stress regulation of gene expression in plants. *Trends Plant Sci.* **12**, 444–451 (2007).
- Du, C. et al. Dynamic transcriptome analysis reveals AP2/ERF transcription factors responsible for cold stress in rapeseed (*Brassica napus* L.). *Mol. Genet. Genom.* **291**, 1053–1067 (2016).
- Roxas, V. P., Smith, R. K. Jr., Allen, E. R. & Allen, R. D. Overexpression of glutathione S-transferase/glutathioneperoxidase enhances the growth of transgenic tobacco seedlings during stress. *Nat. Biotechnol.* **15**, 988–991 (1997).
- 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).
- Wendel, J. F. The wondrous cycles of polyploidy in plants. *Am. J. Bot.* **102**, 1753–1756 (2015).

## Acknowledgements

Genome sequencing, assembly and annotation were conducted by the Novogene Bioinformatics Institute, Beijing, China; mutual contracts were No. NHT140016 and NVT140016004. This work was supported by funding from the Scientific Project

of Shenzhen Urban Administration (201519) and a Major Technical Research Project of the Innovation of Science and Technology Commission of Shenzhen (JSGG20140515164852417). Additional funding was provided in particular by the Scientific Research Program of Sino-Africa Joint Research Center (SAJL201607). We thank X.Q. Wang, G.W. Hu, Z.D. Chen and Y.H. Guo for comments on gnetophyte phylogenetic relationships and ecological issues; H. Wu and X.P. Ning for discussion of related organ development; K.K. Wan and S. Sun for additional help on the analysis of repeats. We also thank X.Y. for support of funding coordination. Y.V.d.P. acknowledges the Multidisciplinary Research Partnership 'Bioinformatics: from nucleotides to networks' Project (no. 01MR0310W) of Ghent University, and funding from the European Union Seventh Framework Programme (FP7/2007-2013) under European Research Council Advanced Grant Agreement 322739-DOUBLEUP.

### 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.L., A.R.L., I.J.L. and Z.J.L. are joint first authors, H.P.X., Y.B.G., Yang Liu, L.Y.C. and W.C.W. are joint second authors. Z.M.L., J.B.Z., J.L., Yan Liu performed the genome assembly and annotation; H.P.X., L.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.L., F.L., H.M.L., T.W., A.R.L., I.J.L., W.X. and Yan Liu 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.L., J.B.Z., J.Y., T.W., L.Z., Y.B.G. and S.H.D. conducted PCA analysis of pfam domains. J.B.Z., T.W., J.L., L.L., L.J.K., Y. Lei and Z.M.L. performed the analysis investigating the divergence of gene families. J.L.H., P.L., Q.M., Yang Liu 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.

### Competing interests

The authors declare no competing financial interests.

### Additional information

**Supplementary information** accompanies this paper at <https://doi.org/10.1038/s41477-017-0097-2>.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Correspondence and requests for materials** should be addressed to Q.-F.W. or Y.V. or J.-B.Z. or X.-M.W.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.



## Life Sciences Reporting Summary

Nature Research wishes to improve the reproducibility of the work we publish. This form is published with all life science papers and is intended to promote consistency and transparency in reporting. All life sciences submissions use this form; while some list items might not apply to an individual manuscript, all fields must be completed for clarity.

For further information on the points included in this form, see Reporting Life Sciences Research. For further information on Nature Research policies, including our data availability policy, see Authors & Referees and the Editorial Policy Checklist.

## ▶ Experimental design

## 1. Sample size

Describe how sample size was determined.

The sequenced individual Gnetophyte plants were collected from natural wild individuals without any specific selection. The detailed information is offered in "method summary" of the main text and section of supplementary Note 1 in the Supplementary Information

## 2. Data exclusions

Describe any data exclusions.

*If no data were excluded from the analyses, state this OR if data were excluded, describe the exclusions and the rationale behind them, indicating whether exclusion criteria were pre-established.*

## 3. Replication

Describe whether the experimental findings were reliably reproduced.

*For each experiment, note whether any attempts at replication failed OR state that all attempts at replication were successful.*

## 4. Randomization

Describe how samples/organisms/participants were allocated into experimental groups.

*Describe how samples were allocated to groups. If allocation was not random, describe how covariates were controlled. If this is not relevant to your study, explain why.*

## 5. Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis.

*Describe the extent of blinding used during data acquisition and analysis. If blinding was not possible, describe why OR explain why blinding was not relevant to your study.*

Note: all studies involving animals and/or human research participants must disclose whether blinding and randomization were used.

## 6. Statistical parameters

For all figures and tables that use statistical methods, confirm that the following items are present in relevant figure legends (or the Methods section if additional space is needed).

- | n/a                                 | Confirmed                |  |
|-------------------------------------|--------------------------|--|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | The <u>exact</u> sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement (animals, litters, cultures, etc.)                                    |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | A description of how samples were collected, noting whether measurements were taken from distinct samples or whether the same sample was measured repeatedly.  |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | A statement indicating how many times each experiment was replicated   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | The statistical test(s) used and whether they are one- or two-sided (note: only common tests should be described solely by name; more complex techniques should be described in the Methods section) |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | A description of any assumptions or corrections, such as an adjustment for multiple comparisons  |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | The test results (e.g. $p$ values) given as exact values whenever possible and with confidence intervals noted   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | A summary of the descriptive statistics, including central tendency (e.g. median, mean) and variation (e.g. standard deviation, interquartile range)   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | Clearly defined error bars   |

*See the web collection on statistics for biologists for further resources and guidance.*

## ► Software

Policy information about availability of computer code

### 7. Software

Describe the software used to analyze the data in this study.

*Provide a description of all commercial and custom code used to analyze the data in this study, specifying the version used.*

For all studies, we encourage code deposition in a community repository (e.g. GitHub). Authors must make computer code available to editors and reviewers upon request. The *Nature Methods* guidance for providing algorithms and software for publication may be useful for any submission.

## ► Materials and reagents

Policy information about availability of materials

### 8. Materials availability

Indicate whether there are restrictions on availability of unique materials or if these materials are only available for distribution by a for-profit company.

*Describe any restrictions on availability of unique materials used in the study OR confirm that all unique materials used are readily available from the authors or from standard commercial sources (and specify these sources) OR state that no unique materials were used.*

### 9. Antibodies

Describe the antibodies used and how they were validated for use in the system under study (i.e. assay and species).

*For all antibodies, as applicable, provide supplier name, catalog number, clone name, and lot number. Also describe the validation of each primary antibody for the species and application, noting any validation statements on the manufacturer's website, relevant citations, antibody profiles in online databases, or data provided in the manuscript OR state that no antibodies were used.*

### 10. Eukaryotic cell lines

a. State the source of each eukaryotic cell line used.

*Provide information on cell line source(s) OR state that no eukaryotic cell lines were used.*

b. Describe the method of cell line authentication used.

*Describe the authentication procedures for each cell line used OR declare that none of the cell lines used have been authenticated OR state that no eukaryotic cell lines were used.*

c. Report whether the cell lines were tested for mycoplasma contamination.

*Confirm that all cell lines tested negative for mycoplasma contamination OR describe the results of the testing for mycoplasma contamination OR declare that the cell lines were not tested for mycoplasma contamination OR state that no eukaryotic cell lines were used.*

d. If any of the cell lines used in the paper are listed in the database of commonly misidentified cell lines maintained by ICLAC, provide a scientific rationale for their use.

*Provide a rationale for the use of commonly misidentified cell lines OR state that no commonly misidentified cell lines were used.*

## ► Animals and human research participants

Policy information about studies involving animals; when reporting animal research, follow the ARRIVE guidelines

### 11. Description of research animals

Provide details on animals and/or animal-derived materials used in the study.

*For laboratory animals, report species, strain, sex and age OR for animals observed in or captured from the field, report species, sex and age where possible OR state that no animals were used.*

Policy information about studies involving human research participants

### 12. Description of human research participants

Describe the covariate-relevant population characteristics of the human research participants.

*Provide all relevant information on human research participants, such as age, gender, genotypic information, past and current diagnosis and treatment categories, etc. OR state that the study did not involve human research participants.*