1   **Title**: Comparative transcriptomics in ferns reveals key innovations and divergent
2   evolution of the secondary cell walls

3

4   **Authors**: Zahin Mohd Ali[1,6], Qiao Wen Tan[1], Peng Ken Lim[1], Hengchi Chen[2,3], Lukas
5   Pfeifer[4], Irene Julca[1,5], Jia Min Lee[1], Birgit Classen[4], Sophie de Vries[6], Jan de Vries[6],
6   Teng Seah Koh[7], Li Li Chin[7], Fanny Vinter[8], Camille Alvarado[8], Amandine Layens[8],
7   Eshchar Mizrachi[9], Mohammed Saddik Motawie[10], Bodil Joergensen[10], Peter Ulvskov[10],
8   Yves van de Peer[2,3,9,11], Boon Chuan Ho[7, 12], Richard Sibout[8], Marek Mutwil[1*]

9

10  1. School of Biological Sciences, Nanyang Technological University, Singapore,
11     Singapore
12  2. Department of Plant Biotechnology and Bioinformatics, Ghent University, 9052
13     Ghent, Belgium
14  3. VIB Center for Plant Systems Biology, VIB, 9052 Ghent, Belgium
15  4. Pharmaceutical Institute, Department of Pharmaceutical Biology, Christian-
16     Albrechts-University of Kiel, Kiel, Germany
17  5. University of Lausanne (UNIL), Genopode 2024.3, 1015 Lausanne, Switzerland
18  6. Department of Applied Bioinformatics, Institute for Microbiology and Genetics,
19     Goettingen Center for Molecular Biosciences (GZMB), Campus Institute Data
20     Science (CIDAS), University of Goettingen, Goldschmidtstrasse 1, D-37077,
21     Göttingen, Germany
22  7. Singapore Botanic Gardens, National Parks Board, 1 Cluny Road, Singapore,
23     259569, Republic of Singapore
24  8. INRAE, UR BIA, F-44316 Nantes, France
25  9. Department of Biochemistry, Genetics and Microbiology, Forestry and Agricultural
26     Biotechnology Institute (FABI), University of Pretoria, Pretoria 0002, South Africa
27  10. Department of Plant and Environmental Sciences, University of Copenhagen,
28     Frederiksberg, Denmark
29  11. College of Horticulture, Academy for Advanced Interdisciplinary Studies, Nanjing
30     Agricultural University, Nanjing, China.

31    12. Department of Biological Sciences, National University of Singapore, Republic of

32        Singapore

33

34    *correspondence: mutwil@ntu.edu.sg

35

36    **Abstract**

37    Despite ferns being crucial to understanding plant evolution, their large and complex

38    genomes has kept their genetic landscape largely uncharted, with only a handful of

39    genomes sequenced and sparse transcriptomic data. Addressing this gap, we generated

40    extensive RNA-sequencing data for multiple organs across 22 representative species

41    over the fern phylogeny, assembling high-quality transcriptomes. These data facilitated

42    the construction of a time-calibrated fern phylogeny covering all major clades, revealing

43    numerous whole-genome duplications and highlighting the uniqueness of fern genetics,

44    with half of the uncovered gene families being fern-specific. Our investigation into fern

45    cell walls through biochemical and immunological analyses identified occurrences of the

46    lignin syringyl unit and its independent evolution in ferns. Moreover, the discovery of an

47    unusual sugar in fern cell walls hints at a divergent evolutionary path in cell wall

48    biochemistry, potentially driven by gene duplication and sub-functionalization. We provide

49    an online database preloaded with genomic and transcriptomic data for ferns and other

50    land plants, which we used to identify an independent evolution of lignocellulosic gene

51    modules in ferns. Our data provide a framework for the unique evolutionary path that ferns

52    have navigated since they split from the last common ancestor of euphyllophytes more

53    than 360 million years ago.

54

55    **Introduction**

56    Since they diverged from a shared ancestor with seed plants more than 360 million years

57    ago, ferns have played a significant role in life on Earth [1]. They occupy various niches in

58    different ecosystems, acting as pioneer species, key ecological players, invasive entities,

59    and contributors to agriculture. They are the second most diverse group of vascular plants

60    after angiosperms, with over 10,500 existing species [2-6]. Ferns exhibit great

61    morphological and physiological diversity, and have evolved equally diverse strategies to

62  cope with environmental challenges [7], such as adaptations to low-light environments [8]. The secondary metabolites produced by ferns and the genes responsible for their biosynthesis are of great interest for environmental clean-up efforts, agriculture, and the discovery of new pharmaceuticals [9–11].

66  Despite fern's ecological importance, the phylogenetic relationship of major clades in Monilophyta (ferns) remains elusive [12]. Based on a Maximum Likelihood (ML) tree of a concatenated matrix of 146 low-copy nuclear genes, Qi and coauthors [13] inferred Marattiales to be sister to Polypodiidae (i.e. the leptosporangiate ferns) as proposed in Pteridophyte Phylogeny Group (PPG) I (2016). Conversely, Shen and coauthors [12] inferred Marattiales to be sister to Ophioglossidae (consisting of Psilotales and Ophioglossales), based on a coalescent-based tree of two low-copy nuclear gene sets of 69 transcriptomes. Nitta and coauthors [14] inferred Gleicheniales as a monophyletic clade based on a ML tree of a concatenated matrix of 79 plastome loci as opposed to the paraphyletic inference of Shen and coauthors [12] and Qi and coauthors [13]. It is also unclear whether horsetails are a sister group to the last common ancestor of all the remaining ferns or other fern clades, such as Marattiales [5,12].

78  Given ferns' critical evolutionary position as the sister group to seed plants, investigating their genomes, coding sequences, and gene families offers unparalleled insights into the evolution of plants [15], especially the key aspects of vasculature and cell walls. The evolution of vasculature and secondary cell walls precipitated a 10-fold increase in plant species numbers (http://www.theplantlist.org/) and shaped the Earth's geo- and biosphere [16]. Ferns thus harbour key information for the evolution of vascular plant form and function [17].

85  However, ferns are infamous for their exceptionally large genomes (on average 12.3 billion base pairs), with one of the largest genome of any living organism - 160 billion base pairs - found in ferns[18]. They also have exceptionally high numbers of chromosomes (averaging at 40.5, with a peak at 720)[19], which are believed to result from multiple instances of whole-genome duplication [20–22] and a relatively slow genome downsizing process [23]. Among plants, ferns show the highest rate of polyploidy-driven speciation [24], a direct relationship between genome size, chromosome number and the age of long terminal repeat-retrotransposon (LTR-RT) insertions [25–27], and a high rate of whole

93    genome duplications (WGDs) among several fern lineages [23,28]. However, the genetic

94    and genomic evidence for widespread whole-genome duplication in ferns remains largely

95    unexplored [29–31].

96         Thus, comparative studies that investigate the evolution of ferns have been

97    hampered by their large, complex genomes, limiting our understanding of fern genome

98    evolution and the genetic underpinnings of the evolution of vasculature and cell walls. To

99    date, only few fern genomes and transcriptomes are available [20,32–35], and no studies that

100   conducted a comprehensive comparison of gene inventories, transcriptional programs

101   and biochemical properties of their cell walls have been reported.

102        To address this, we generated 405 RNA-sequencing samples to generate coding

103   sequences and gene expression atlases for 22 fern species, capturing major

104   representatives of ten fern orders. We investigated ancient polyploidy, the distribution of

105   fern-specific gene families, how gene age correlates with organ-specific expression, and

106   predicted the functions of fern-specific genes. To better understand how fern cell walls

107   have evolved, we performed a comprehensive histological and biochemical analysis of

108   fern tissues and propose a novel biosynthetic pathway of lignocellulose. We further

109   detected a wide-spread occurrence of the unusual hemicellulose mixed-linkage glucan,

110   and show that it evolved independently in ferns, and propose candidate

111   glucosyltransferases responsible for its synthesis. We also detected a novel type of a

112   methylated sugar, a 2-O-Methyl-D-glucopyranose. We also show that ferns likely

113   independently evolved secondary cell walls through several duplication events in the

114   cellulose synthase family. Finally, we make our fern genomic and transcriptomic data

115   easily accessible with the CoNekt database (https://conekt.plant.tools/).

116        Our data, findings, and tools shed light on the evolution of cell walls, lignin,

117   specialised metabolism, and organ development in ferns and other land plants. We

118   envision that similar large, comparative studies will elucidate the evolution of plants and
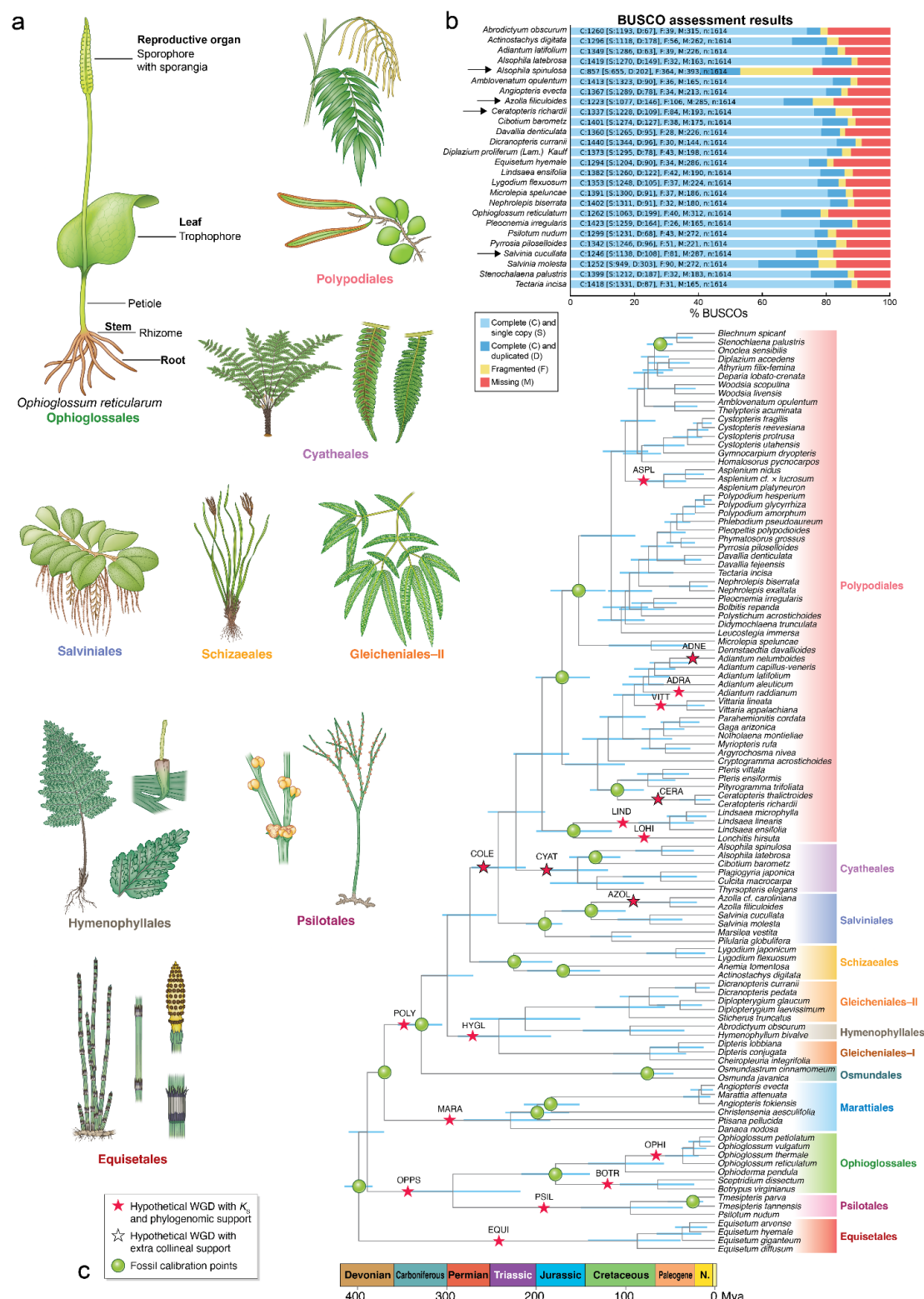
119   other organisms.

120

**Figure 1. Sampling, transcriptome assembly, and species tree of major representatives of ferns.** a) *Ophioglossum reticulatum* with sampled organs labelled, together with samples representatives of the other fern ordersb) Completeness of

5

124  transcriptome assembly measured by Benchmarking Universal Single-Copy Orthologs
125  (BUSCO). *Alsophila spinulosa, Ceratopteris richardii, Azolla filiculoides* and *Salvinia*
126  *cucullata* have available genomes, while the remaining values are for the transcriptome
127  assemblies reported here. C) The evolutionary timescale of Monilophyta based on the
128  inferred consensus fern cladogram. The species tree shows the inferred consensus
129  phylogenetic topology with branch lengths representing absolute divergence time
130  estimated by Bayesian molecular dating analysis. The horizontal coordinates of each
131  internal node denotes the posterior mean divergence time while the bars represent the
132  95% Highest Posterior Density. Hypothetical WGDs are indicated at corresponding
133  phylogenetic nodes as red stars with four-letter identifiers. Black outline around a red star
134  indicates events with additional collinear support. Fossil calibrations are indicated at
135  corresponding phylogenetic nodes with green circles. The clade strips indicating
136  affiliations at order level are shown as vertical bars with distinct colours. The geological
137  timeline refers to the International Commission on Stratigraphy (ICS) v2023/09.
138

## Results

### Construction of fern coding sequences by transcriptome assembly

141  To capture the diversity of ferns, we selected 22 candidate species representing ten

142  orders, which were photographed and dissected on site (Figure S1), with organs

143  categorised with localities and vouchers attached (Table S1). We collected 405 RNA-seq

144  samples (Table S2), capturing 25 specific organs at different developmental stages,

145  categorised into four major organs - leaves, roots, stems, and reproductive organs - for

146  simplified comparison (Figure 1a)(Table S1). Our transcriptome assembly pipeline

147  combined TRINITY and k-mer SOAPdenovo-Trans assemblies concatenated with

148  EvidentialGene (Figure S2)[36,37]. We removed any potential non-fern mRNA contaminants

149  and any sequences with aberrant GC content due to assembly artefacts, low transcripts

150  per million (TPM) values, or sequence similarity higher to non-fern species than to ferns

151  (see methods, Figure S2). The assembly yielded 30,000–100,000 coding sequences

152  (CDSs) per species with high Benchmarking Universal Single-Copy Orthologs (BUSCO)

153  scores (Figure 1b, Table S3) that rivalled the scores of the four sequenced genomes

154  (Figure 1b, black arrows).

155

### Reconstruction of the evolutionary timeline of ferns

157  Given the standing phylogenetic discordance, we reconstructed a phylogenetic tree of

158  108 fern species (22 from this study, 7 sequenced genomes, and 79 from the 1000 Plant

159  Transcriptomes Initiative (1KP) and other studies)[13,20,32–35,38–40], covering the whole

160    backbone of Monilophyta (Table S4, Supplemental Methods 1). Four datasets, each with

161    a different outgroup (horsetails, seed plants, lycopods or bryophytes), were first used in

162    nucleotide with three different methods including ASTRAL-Pro2 [41], concatenation-based

163    method and STAG [42](Supplemental Methods 1) on the 107 ferns dataset and then

164    reanalyzed in both nucleotide and peptide on the 108 ferns (adding the latest *Marsilea*

165    *vestita* genome [39]) using the favored method ASTRAL-Pro2. We recovered a well-

166    supported fern backbone phylogeny in which Marattiales was inferred to be sister to

167    Polypodiidae (i.e. leptosporangiate ferns) and Gleicheniales was inferred as a

168    paraphyletic clade (Supplemental Methods 1, Figure SM5-6). A closer phylogenetic

169    relationship between Marattiales and Polypodiidae than between Marattiales and

170    Ophioglossidae was supported in all datasets with Local Posterior Probability (LPP) as

171    1.00, except the one with bryophytes as outgroup based on peptide alignment.

172    Phylogenies derived from this dataset favored Marattiales and Ophioglossidae as sister

173    groups, with LPP as 0.82 (Supplemental Methods 1, Figure SM6). The monophyly of

174    Gleicheniales was not supported, whereas the Gleicheniales−II (Gleicheniaceae) was

175    closer to Hymenophyllales than Gleicheniales−I (Dipteridaceae) in all datasets, except

176    the one with bryophytes as outgroup based on nucleotide alignment, whose conflicting

177    branching pattern was only supported by LPP as 0.46 (Supplemental Methods 1, Figure

178    SM5).

179    In addition, these analyses provided strong support for a scenario wherein

180    horsetails are a sister group to the last common ancestor of all the remaining ferns in the

181    phylogeny inferred from every combinational setting of non-horsetails outgroups and

182    methods (Supplemental Methods 1). We selected phylogeny derived from the nucleotide

183    dataset with horsetails as outgroup using ASTRAL-Pro2 as the consensus tree and used

184    in all our subsequent analyses (Figure 1c). This selection was based on the general

185    consistency across datasets with varied outgroups, and on the ASTRAL-Pro2 derived

186    phylogeny with STAG method (Supplemental Methods 1).

187    To estimate the absolute divergence time of the 108 ferns, we used Bayesian

188    molecular dating under the independent rate and LG general amino acid substitution

7

189  model [43] with 18 soft fossil constraints (indicated in Figure 1c as green dots, Table S5).

190  The 95% Highest Posterior Density (HPD) and posterior mean of the stem ages of major

191  fern clades was summarised (Supplemental Methods 1, Table S6). The resulting high-

192  confidence fossil-calibrated tree thus resolved a long-standing discussion on fern

193  phylogenetics.

194

195  **Identification of 18 separate whole genome duplication events in ferns**

196  It has been proposed that ferns have a large number of chromosomes due to repeated

197  rounds of whole-genome duplication (WGD)[44]. Here, we used $K_S$-age distributions and

198  phylogenomic methods to unveil remnants of ancient WGDs [45]. In total, we found support

199  for 18 hypothetical WGDs within the backbone of ferns, of which five with attained

200  collinear support (red stars, Figure 1c, see Supplemental Methods 1). Seven of the WGDs

201  were found within Polypodiales, two in Ophioglossales, and one in each of the lineages

202  of Equisetales, Psilotales, Marattiales, Salviniales and Cyatheales, together with four

203  shared by more than one order. WGDs were previously identified in other studies with

204  different phylogenetic locations [28,46]. We reevaluated possible scenarios thereof and

205  proposed WGDs with both $K_S$ and phylogenomic support after correcting rate variation

206  and taking into account the uncertainty in gene tree and gene tree-species tree

207  reconciliation (Supplemental Methods 1). Next, we tested the species richness and

208  genome size as a function of the number of ancient WGD events shared by different major

209  clades (excluding Polypodiales due to numerous nested WGD events therein), and

210  observed a significant positive correlation between the number of ancient WGD events

211  and the number of species in a lineage (Figure S3, adjusted $R^2$ = 0.41, p-value = 0.02).

212  Conversely, we did not observe any significant correlation between the number of ancient

213  WGD events and genome size (Figure S3).

214

215  **Diversity and conservation of fern gene functions and expression patterns**

216  To investigate predicted gene functions within ferns, we conducted a phylostratigraphic

217  analysis encompassing one glaucophyte, seven chlorophytes, three bryophytes, two

218  lycophytes, 26 ferns (22 fern transcriptomes and four genomes), two gymnosperms, and

219  six angiosperms (see methods). The 47 species are used to study the gain and losses of

220   gene families across Archaeplastida, by assigning the orthogroups to nodes ranging from
221   node 1 (the earliest ancestor of Archaeplastida) to node 13 (the ancestor of
222   Polypodiales)(Figure 2a). The nodes are based on the fern phylogeny tree shown in
223   Figure 1c and known relationships between Archaeplastida [47].

224   Orthogroup gains, a measure of new gene family acquisition, were highest in the
225   early stages of algae evolution (nodes 1 and 2 gained 5368 and 3306 orthogroups,
226   respectively) and when plants colonized the land (node 3, 2629 gained
227   orthogroups)(Figure 2a). However, a substantial number of orthogroup gains and losses
228   were also observed within the different fern lineages (e.g., 2408 gained and 193 lost
229   orthogroups in node 9, Figure 2a).

230   The analysis further revealed that ~50% of fern gene families are fern-specific
231   (Figure 2b nodes 6–13 in red, Table S7). For example, >50% of gene families in
232   *Stenochlaena palustris* belong to nodes 6–13 (fern-specific nodes 6–13 in red, Figure 2a,
233   Table S7), suggesting that the fern lineage has evolved genes with novel, unexplored
234   functions. We analysed sequences of 25 orthogroups comprising at least ten fern species.
235   Representatives of 17 orthogroups show no significant sequence or structural similarity
236   to non-fern species (Figure S4), with eight of them being disordered (few to no secondary
237   structures). These results show that >50% of gene families in ferns represent novel,
238   uncharacterized proteins, indicating that studying ferns will likely provide new insights into
239   plant biology and evolution.

240   To understand how gene age correlates with gene expression specificity, we
241   identified organ-specific genes with specificity measure (SPM) analysis (distributions of
242   SPM values are shown in Figure S5, expression profiles of organ-specific genes are
243   shown in Figure S6, Table S8)[47]. Genes belonging to older nodes (nodes 1–4) were less
244   organ-specific (<20%) than younger nodes (nodes 5–13, 30–60%, Figure 2c). Younger
245   genes (nodes 5–13) had specialised functions in a particular organ, with roots having the
246   highest number of specifically expressed genes (Figure 2c). Furthermore, species-
247   specific genes also tended to show a less ubiquitous, more organ-specific expression
248   (Figure S7), indicating an overall negative association between gene age and organ
249   specificity, which is in line with similar observation in land plants [47].
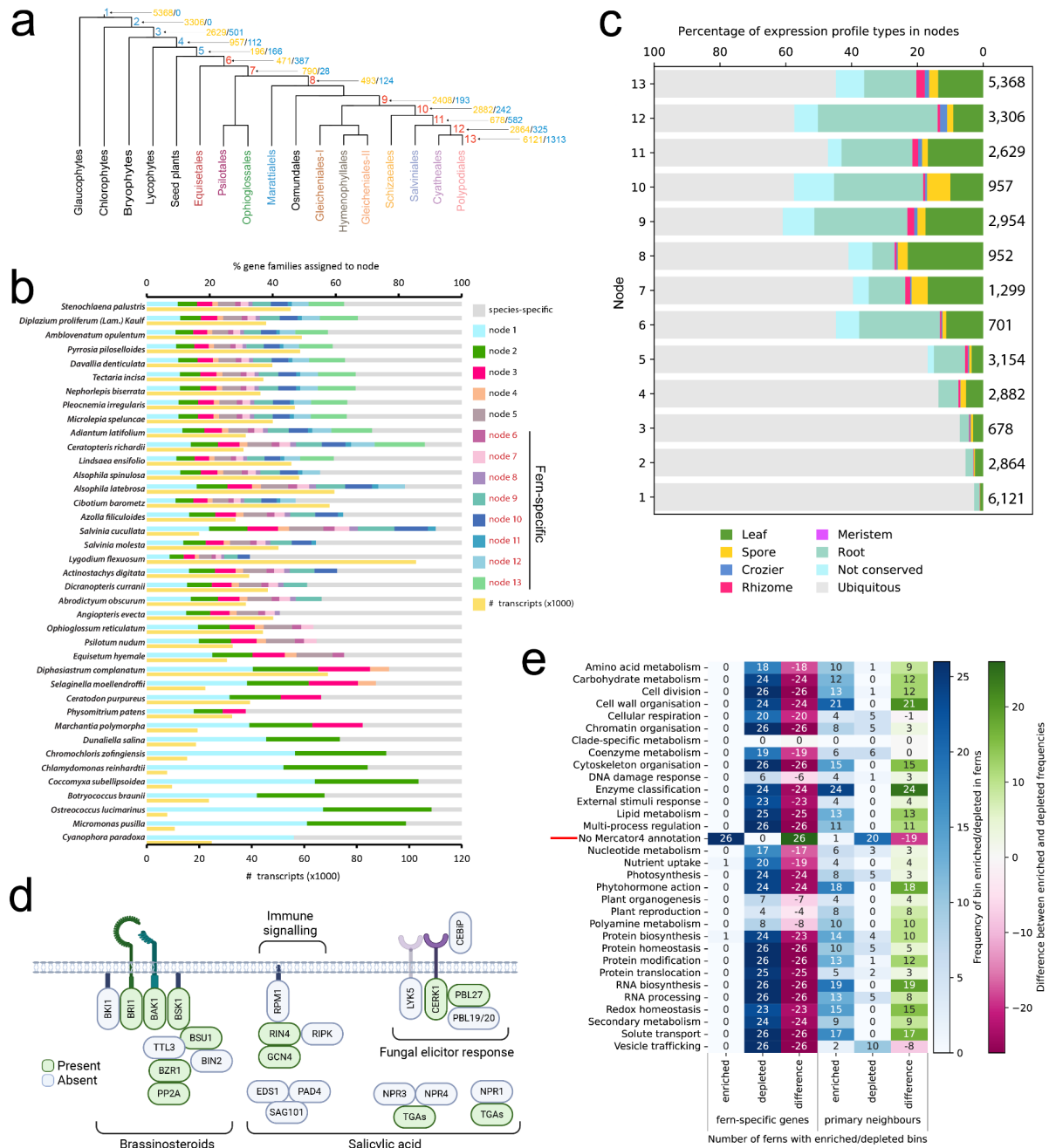
**Figure 2. Gene functions in ferns.** a) Division tree of Archaeplastida. Leaves represent orders, while node numbers correspond to the phylostrata. The orange and blue numbers indicate the number of gains and losses of orthogroups, respectively. b) Stacked bar plot showing the percentage of gene families belonging to a phylostrata (node). Species-specific gene families comprised genes from only one species and were not assigned to nodes. c) Percentage of gene families identified as organ-specific for each node. Numbers on the right side of the plot indicate the number of orthogroups per node. d) Examples of signalling components present (green shapes) or absent (grey shapes) in ferns. e) Clustered heatmap showing enrichment and depletion of biological processes of

10

fern-specific genes and primary co-expression neighbours in 26 (22 from this study and four sequenced genomes) fern species analysed. The left column indicates the biological processes defined by Mapman bins. The scores in the 'enriched' and 'depleted' column indicate in how many of the 26 species fern-specific genes are significantly (BH adj.p < 0.05), connected or disconnected to genes belonging to a specific bin, respectively. A higher value in difference (# enriched - # depleted) indicates overall enrichment, while lower values indicate overall depletion for respective bins in the 26 ferns.

Finally, we investigated whether organ-specific genes are conserved across ferns and other land plants. While many organs express significantly similar sets of genes across ferns and other land plants, we observed a clear difference between fern and seed plant transcriptomes (Figure S8). Not surprisingly, the organ-specific gene sets of seed-containing plants show higher mutual similarity to those of other seed plants (Figure S8, blue box), while those of ferns show the highest similarity to those of other ferns (Figure S8, green box).

**Ferns lack several genes essential for hormonal signalling, defence and development in angiosperms, indicating their unique developmental and environmental strategies**

Terrestrialization and the evolution of seeds and flowers required the evolution of many biological functions [48], which is readily visible when comparing gene inventories of algae, land, seed and flowering plants (Figure S9, gene function completness indicated by darker cells). We used MapMan sequence-based annotations and compared the gene function repertoire of ferns and model angiosperms and observed the absence of several components in ferns (missing functional categories indicated with red text in Table S9, Supplemental Methods 2).

Hormone signalling

Missing components include abscisic acid regulation, perception, and transport, auxin methylation-based degradation, brassinosteroid signalling (Figure 2d)[49] and degradation, cytokinin degradation and transport, and degradation of gibberellins and jasmonic acid and their transport genes.

292       For example, several components of the salicylic acid (*SAG101, EDS1, PAD4*,

293   Figure 2d)[50] and strigolactone signalling pathways were missing (Table S9), as is the

294   degradation component of the former hormone. To test this further, we investigated the

295   presence of canonical NPR domains (NPR1-like C superfamily, BTB/POZ NPR plant

296   domain or BTB/POZ superfamily, and an ANKYRIN domain) in our fern transcriptomes.

297   Of the 26 ferns we studied, 22 had at least one canonical NPR (Supplemental Methods

298   2, Figure SM7). This is consistent with previous evidence that the duplication of NPR1/3/4

299   happened sometime during angiosperm diversification, long after the split between

300   flowering plants and ferns [51]. The SAG101/PAD4/EDS1 module, on the other hand,

301   appears to be a more recent invention of flowering plants, as it is mainly absent in non-

302   seed plants (Supplemental Methods 2, Figure SM8).

303       Further, we investigated perception and downstream signalling of jasmonic acid

304   (JA), focusing on the JA receptor COI1 and the JAZ transcriptional repressors. Both COI1

305   and JAZ candidates are encoded in fern transcriptomes. While *Arabidopsis thaliana* and

306   *Marchantia polymorpha* both have only one copy of COI1, ferns show several gene

307   duplication events, some of which are species-specific, and some of which appear more

308   ancient (Supplemental Methods 2, Figure SM9). The current evolutionary model of JA

309   perception is that the COI1 ligand switched from *dn*-cis-OPDA to JA-Ile in the ancestor of

310   vascular plants [52]. The radiation of COI1 in ferns, however, suggests functional

311   divergence in jasmonate perception, possibly complicating its evolutionary history. This

312   highlights ferns as a key lineage for further functional investigation to understand the

313   evolution of plant immunity.

314

315   <u>Secondary metabolism</u>

316   Phytochemical studies on ferns have revealed that they contain a wide range of

317   secondary metabolites, many of which are function herbivore defense and show bioactive

318   properties [9]. For secondary metabolite pathways associated with biotic interactions, we

319   observed that multiple genes known to act in the flavonoid biosynthesis pathway were

320   missing in all fern species analysed (Table S9), agreeing with previous datasets on the

321   evolution of red pigmentation in land plants [53,54]. Yet, ferns are able to synthesise

322   flavonoids [55].

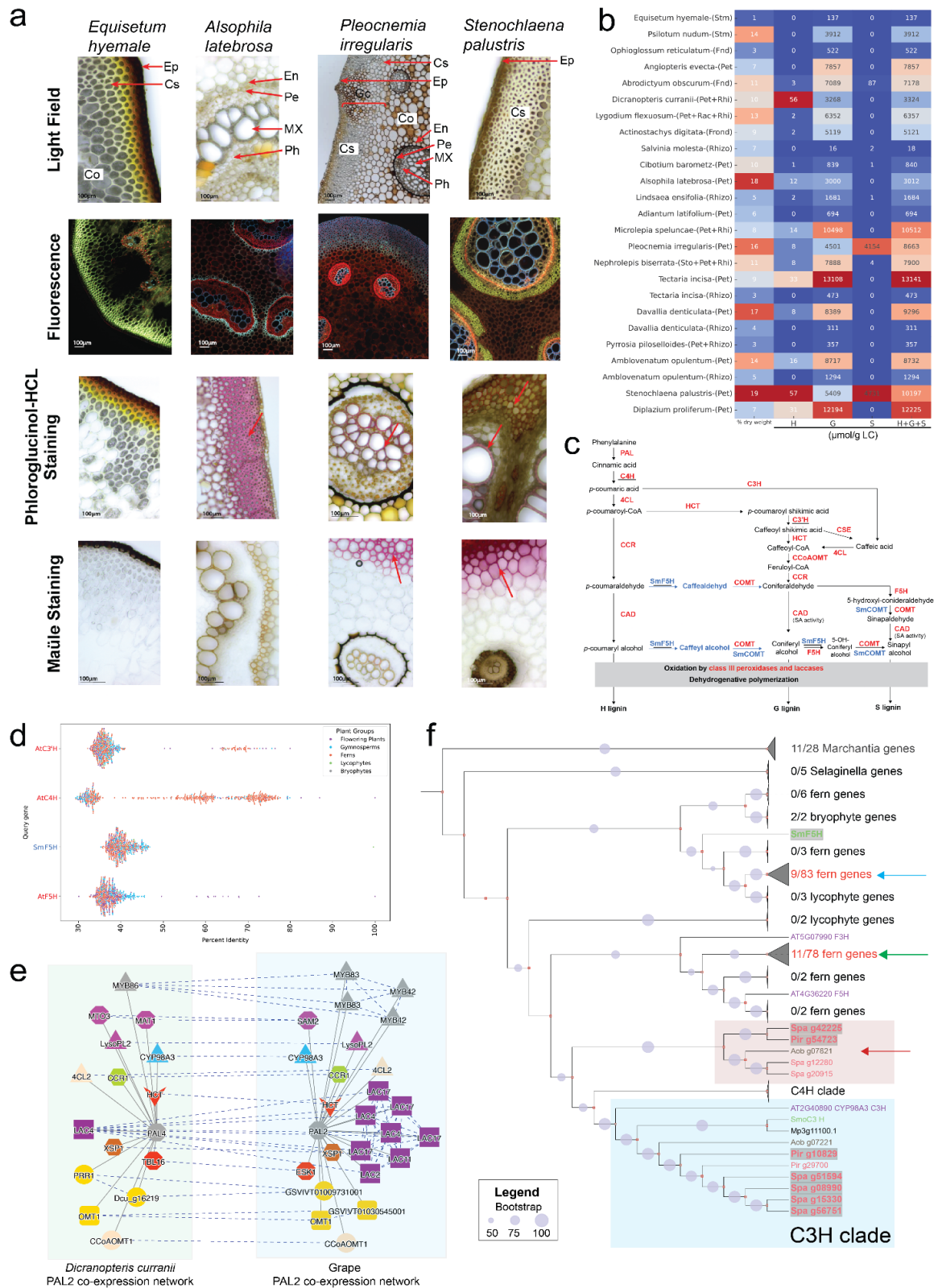**Figure 3. Lignin analysis of ferns.** a) Sections of stems (Equisetum) and petioles (Alsophila, Pleocnemia, Stenochlaena). Red arrows with labels indicate cortex (Co),

cortical sclerenchyma (Cs), endodermis (En), epidermis (Ep), metaxylem (MX), pericycle (Pe) and phloem (Ph). Red arrows without labels indicate stained cell walls. b) Percentage of lignin (1st column), H, G, S (2nd-4th) and H+G+S (5th) thioacidolysis lignin units. c) The H, G, and S lignin unit biosynthesis pathway for angiosperms (red text) and lycophytes (blue text). Intermediate metabolites are indicated by black text, while the grey box contains enzymes involved in lignin polymerization. d) BLAST scores (x-axis) of AtC3H, AtC4H, SmoF5H and AtF5H against the translated transcriptomes found in the CoNekT database. Each point represents a protein. e) Comparative co-expression network analysis of *PAL* genes from grape and fern *Dicranopteris curranii.* Nodes represent genes, solid edges connect co-expressed genes, while dashed edges connect orthologs. Coloured shapes represent different orthogroups, while the gene names are based on the best BLAST hits to *Arabidopsis thaliana*. f) Phylogenetic tree of land CYP450s of *P. irregularis* (genes starting with Pir), *S. palustris* (Spa), *A. obscurum* (Aob), *Selaginella moellendorffii* (Smo), *Marchantia polymorpha* (Mp). The Arabidopsis (AT) lignin-related C4H, C3'H and F5H and flavonoid-related F3H are included.

### External stimuli response

Plants have evolved elaborate signalling and response pathways to cope with the changing environment. For several of these pathways, we observed that orthologs of phototropin-mediated receptors, all $CO_2$ sensing and signalling components, and many gravity-sensing proteins were absent in our fern transcriptomes (Table S9). Genes known to be essential for sensing and responding to temperature in flowering plants were present in ferns, but acquired thermotolerance factors were missing. Other missing proteins include those involved in several pathogenesis-related processes, such as pattern- and effector-triggered immunity (16 out of 36 factors)(Figure 2d), WRKY33-dependent immunity, pathogen polygalacturonase inhibitors, and basic chitinases. While some components of symbiosis pathways are present in our fern transcriptomes (Table S9), many factors are absent, such as mycorrhizal response genes and transporters.

### Transcript control and modification

Several components controlling mRNA and protein levels are also absent, such as more than half of the subgroups of MYBs and most REMs. Organellar RNA processing is lacking plastidial and mitochondrial CFM-type splicing factors, a majority of mitochondrial RBA splicing components (>20), C-to-U RNA editing (>50 factors), and mRNA

361  stabilisation and deadenylation factors. For protein homeostasis, we found only class C-
362  I and C-II small HSP holdase chaperones, while the ten other classes were absent (C-III
363  to ER), together with E3 ubiquitin ligases from groups IV and V.

364

365  Reproduction and organ development

366  Not surprisingly, our transcriptomes indicate that ferns differ from flowering plants in their
367  gene inventories related to reproduction. Ferns lack genes associated with anther
368  dehiscence (*PCS1, NST1/2, MYB26*), pollen aperture formation (*INP1/2*), pollen tube
369  growth (except *GEX3*), embryo axis formation (except *ATML1*), endosperm formation
370  (exception *GLAUCE*) and seed formation and dormancy (Table S9). On the other hand,
371  ferns contain nearly all male gametogenesis (e.g., *DUO1/3, DAZ, APD*) and exine
372  (*ROCK/TEX2, DEX1, NEF1*) formation factors, stamen (*TPD1, EMS1, JAG*) and tapetum
373  (*DYT1, TDF1*) regulators and most factors important for female gametophytes (*AMP1,*
374  *CYP78a, RKD, MAA3*) but lack genes essential for central cell formation (Table S9).
375  Surprisingly, while ferns are seedless, they contain most genes important for seed
376  maturation and globulins.

377  Interestingly, most flower formation photoperiodic and autonomous promotion
378  pathway genes are present in ferns and bryophytes. However, as expected, most genes
379  important for floral transition are missing (*FRIGIDA, FRL1/2, FES1*), except *FRI-C* effector
380  complex genes, floral meristem identity (*LMI2, AP1/3, PISTILLATA, SEPALLA*), and
381  morphogenesis (*BLR, ETT*). This suggests that the flower formation pathways have other
382  roles in ferns, possibly linked to photoperiodic response, developmental timing,
383  sporulation control, or other process.

384  For organ development, ferns lack several key genes essential for leaf adaxial and
385  abaxial polarity, guard cell formation, and stomatal density (Table S9). Their root
386  developmental programs are likely also different from flowering plants, as they lack the
387  entire MYB-bHLH-WD40 transcriptional regulatory module, and genes controlling
388  columella apical meristem (*WOX5, FEZ, SMB, BRN1/2*), and endodermis meristem
389  regulation and signalling (*SHR, SCR, KOIN, IRK*). More than half of Casparian strip
390  factors are missing, and nearly all vascular system formation factors (only 2 out of 14
391  transcription factors present).

15

392       Taken together, the analyses shown in Figure 2a-e provide further support for the
393    presence of unique growth, development, and survival strategies in ferns, and suggests
394    that additional research on them is worthwhile.

395

**Co-expression-based prediction of gene functions in the fern lineage**

397    The presence of >50% of fern-specific orthogroups (Figure 2b) indicated that ferns might
398    have evolved as-yet unknown gene functions on a massive scale. To investigate whether
399    the fern-specific genes can be annotated by sequence-similarity approaches, performed
400    an enrichment analysis of their biological functions. We observed enrichment for MapMan
401    bin containing uncharacterized genes ('No Mercator4 annotation', red line, Figure 2e),
402    and a depletion of bins related to known biological processes. This indicates that
403    sequence similarity approaches cannot infer the functions of most fern-specific genes.

404       Gene co-expression networks can reveal the functions of non-annotated genes
405    based on the guilt-by-association principle [56], where genes with similar expression profiles
406    tend to be involved in the same biological process. To predict the functions of fern-specific
407    genes without relying on sequence similarity to genes with known functions, we calculated
408    the functional enrichment of their direct neighbours in the co-expression networks.
409    Interestingly, fern-specific genes are significantly co-expressed (> 17 fern species) with
410    biological processes such as 'Cell wall organisation', 'Enzyme classification',
411    'Phytohormone action' and 'RNA biosynthesis' and moderately co-expressed (> 10 fern
412    species) with 'Cell division', 'Cytoskeleton organisation', 'Lipid metabolism', 'Multi-
413    process regulation', 'Protein biosynthesis', 'Protein modification', and others (Figure S10,
414    Figure 2e). This indicates that these genes are involved in most biological processes in
415    ferns, especially cell wall, development and new metabolic pathways.

416

**S lignin has evolved independently in the fern lineage**

418    Lignin is a complex phenolic polymer that forms essential structural materials in the
419    support tissues of vascular plants. Importantly, this polymer also confers hydrophobicity
420    to xylem vessels, allowing water transport from roots to leaves and enabling plants to
421    grow out on land. Lignin is primarily composed of three monolignols: *p*-coumaryl, coniferyl
422    and sinapyl alcohols, which are named *p*-hydroxyphenyl (H), guaiacyl (G) and syringyl

423   (S) units when incorporated into the polymer [57]. Several studies, including one focused

424   on ferns [58], suggest a complex evolutionary history that may include independent

425   evolutionary paths for lignin synthesis, particularly the S units, among different plant

426   lineages [59]. However, testing this was difficult without additional genomic information. Our

427   fern transcriptome datasets provided an opportunity to explore the evolution of lignin

428   across the entire fern family.

429        We first characterised the presence and sites of deposition of the different lignin

430   units in nine ferns from three orders: Equisetales, Cyatheales and Polypodiales, using a

431   simple staining procedure. We stained cross-sections of stems and petioles with

432   Phloroglucinol-HCl, which reacts with coniferaldehyde residues of lignin to generate a red

433   condensation product [60,61]. All nine ferns showed the presence of lignin. However staining

434   was not or poorly discernable in vessels of *Equisetum hyemale* and *Adiantum latifolium*

435   (Figure 3a, Figure S11). In many species, lignin was mostly found in the subcortical

436   sclerenchyma or outer layers of the metaxylem tissues, as expected (Figure S11).

437   Interestingly, Mäule staining on stem cross-sections indicated the presence of S-units in

438   *Pleocnemia irregularis* and *Stenochlaena palustris* [62] (Figure 3a, Figure S11), indicating

439   that this subunit predominantly found in angiosperms [63], is also widespread in ferns.

440        Given these results, we further determined lignin content and structure within all

441   22 fern species, using the CASA method [64] and thioacidolysis followed by Gas

442   Chromatography-Mass Spectrometry (GC-MS) [65,66], respectively (Figure 3b). Overall, we

443   observed a large variability in total lignin content and the different units among the ferns.

444   Not surprisingly, *Equisetum hyemale* showed the lowest CASA lignin content (1 %), and

445   *Stenochlaena palustris* the highest (19 %)(Figure 3b). Most ferns contain lignin composed

446   of G units (130 - 13000 µmol/g of CASA lignin), while H units are less abundant and, in

447   some cases, not detectable (0 - 57 µmol/g). Interestingly, we observed substantial

448   differences in the lignin content of multiple organs when analyzed. For example, H units

449   were detectable in petioles but not in rhizomes of *D. denticulata* and *A. opulentum* (Figure

450   3b, Table S11). We observed S units in high quantities in *P. irregularis* and *S. palustris*

451   (>4000 µmol/g), medium quantities in *A. obscurum* (86.9 µmol/g CASA) and minute but

452   detectable quantities in *S. molesta, Cibotium barometz, Lindsaea ensifolia, Nephrolepis*

453   *biserrata* (<5.0 µmol/gCASA).

17

454        Next, we set out to identify the biosynthetic pathways of lignin, focusing on S units.

455   To analyze the pathways and make the fern gene expression data easily accessible, we

456   uploaded the expression data for the 22 ferns to our CoNeKT database

457   (https://conekt.sbs.ntu.edu.sg/)[67], upgrading the database to comprise 39 species,

458   including angiosperms, lycophytes, bryophytes and algae. S unit synthesis evolved

459   independently in the lycophyte *Selaginella* and angiosperms [68], illustrated in Figure 3c.

460   Unlike angiosperms, which require *p*-coumarate 3-hydroxylase (C3'H) and ferulate 5-

461   hydroxylase (F5H) to make S units (Figure 3c, black arrows), *Selaginella* utilises a

462   multifunctional F5H that skips several steps of the pathway to make caffealdehyde and

463   caffeyl alcohol, which can then be utilised to make G and S lignin (Figure 3c, blue text)[69].

464   Blasting *AtC3H* and *AtC4H* (Cinnamate 4-hydroxylase) against all species proteomes in

465   the CoNeKT database (https://conekt.sbs.ntu.edu.sg/blast/), showed identity scores of

466   >60% for ferns (Figure 3d, orange points), indicating that ferns likely contain C3'H and

467   C4H enzymes (Table S12). However, *AtF5H* and *SmF5H* showed only low sequence

468   identity to fern proteomes (~40%), which according previous studies indicates an absence

469   of known F5H enzymes in ferns [70].

470        To identify candidate fern F5H enzymes, we took advantage of the observation

471   that lignin biosynthetic genes tend to be tightly coexpressed and that these relationships

472   are conserved even across large evolutionary distances [71]. Indeed, comparing the co-

473   expression networks of *PAL* genes from fern *Dicranopteris* and angiosperm *Vitis vinifera*

474   (grape, Vitaceae) revealed many of the expected enzymes and a CYP98A3-like gene that

475   could likely represent *C3H* (Figure 3e, query gene *Dcu_g01768*, co-expression networks

476   of the lignin genes are in Supplemental Data 1). Furthermore, most of the fern lignin

477   biosynthetic genes are co-expressed with at least two other relevant enzymes (Figure

478   S12a), and the co-expression networks can suggest the unknown components (Figure

479   3f), including transcription factors and CYP450 enzymes (Figure 3f, Figure S12). To

480   suggest the identity of fern F5H enzymes, we first performed phylogenetic analysis of all

481   CYP450s of S lignin-producing ferns Pir, Spa, Aob, angiosperm *Arabidopsis*, lycophyte

482   *Selaginella* and included the outgroup bryophyte *Marchantia* that does not produce S

483   units (Figure 3g). We then indicated which CYPs are co-expressed with at least one lignin

484   enzyme. As expected, *C3H* genes are co-expressed with the other lignin biosynthetic

485    enzymes (co-expressed genes indicated with grey boxes, Figure 3g). However, we

486    observed several clades in the tree that likely emerged independently in ferns and

487    contained groups of co-expressed CYP450 enzymes (Figure 3g, indicated by red, green,

488    and blue arrows). These enzymes comprise prime candidates for the discovery of F5H

489    enzymes in ferns.

490

491    **Members of the Polypodiales contain a non-canonical cell wall sugar**

492    To further understand the evolution of fern cell walls, we carried out a monosaccharide

493    composition analysis using Gas Chromatography on the 22 ferns that were part of our

494    transcriptomic study (Table S13). The most abundant sugars were glucose (a building

495    block of cellulose, mixed-linkage glucans, xyloglucans), mannose (mannans) and xylose

496    (xylans)(Figure 4a). Less abundant sugars were rhamnose (pectic rhamnogalacturonan

497    I, arabinogalactan-protein), fucose (rhamnogalacturonan II, xyloglucan, arabinogalactan-

498    protein), arabinose (hemicellulose arabinoxylan, rhamnogalacturonan I and II,

499    arabinogalactan-protein) and galactose (rhamnogalacturonan I, hemicellulose

500    galactomannans, arabinogalactan-protein). The proportions of various sugars changed

501    among species and among different organs of the same species, which is in line with

502    previous observations [72]. For example, the *T. incisa* rhizome exhibited a higher proportion

503    of glucose than the petiole of the same species, and higher than rhizomes in of *Davaillia*

504    *denticulata* and *Ambloventanum opulentum* (Figure 4a). In addition to these sugars, we

505    also observed trace amounts of methylated rhamnose (3*O*-MeRha*p*), a known sugar

506    found in arabinogalactan proteins in ferns [73], in all species except for *Psilotum nudum*

507    and *Dicranopteris curranii* (Figure S13),

508        Interestingly, we detected an unknown peak from samples derived from three

509    species, *T. incisa*, *A. opulentum* and *D. proliferum* (Figure 4a, red bars). Because this

510    peak was not observed with our common standards during Gas Chromatography analysis

511    (data not shown), it likely represented a novel sugar. Since initial GC-MS analyses

512    suggested the sugar to be a methylated hexose (data not shown), we synthesised a panel

513    of methylated sugars (Figure S14). Out of the six methylated sugars, only 2-*O*-Methyl-D-

514    glucopyranose (2O-Me-Glc*p*) showed identical retention time and mass spectrum to the

515    unknown sugar (Figure 4b, Figure S15), indicating that these three species of ferns

516  produce a novel sugar. While methylated sugars such as 4-*O*-D-Methyl-Glucuronic Acid

517  are present in plant cell walls [74], this is to our knowledge the first report of 2O-Me-Glc*p*

518  and warrants further study.



519

**Figure 4. Polysaccharide analysis.** a) Percentage of total neutral sugars estimated by GC-MS. The sugars are rhamnose, fucose, arabinose, xylose, mannose, galactose and glucose. b) GC spectra of the unknown peak from c) The number of orthogroups involved in cell wall biosynthesis in land plants. Columns represent species, while rows correspond to a given gene family. The rows are further divided into different polysaccharide classes, separated by horizontal lines. Red and blue numbers indicate that a given species contains significantly more/less genes than others (adjusted p-value < 0.05). Darker colors of boxes indicate more gene copies in each row. The red arrows indicate rows which are particularly depleted in ferns. Ferns are indicated with bold names and thick black lines. d) Relative epitope abundance for five fern species quantified by Comprehensive Microarray Polymer Profiling (CoMMP). Rows indicate the different species and organs, while columns represent the obtained signal from the different antibodies. The colours of the cells correspond to the signal strength. e) Schematic drawing of the major cell wall polysaccharides. Each colour-coded shape represents a sugar or amino acid. The antibodies binding to a respective epitope are coloured with blue letters, while the black bold letters indicate the biosynthetic enzymes.

## Ferns contain most but not all cell wall polysaccharides of angiosperms

We next performed a large-scale comparative analysis of Carbohydrate-Active enZYmes (CAZymes) in land plants (Figure 4c)[75]. The CAZyme database contains genes involved in cell wall biosynthesis, allowing us to compare similarities and differences of ferns and other land plants. To do this, we calculated with species contain significantly (adjusted p-value <0.05) more (red numbers) or less (blue numbers) than the other species. Ferns contain fewer xyloglucan-related gene families involved in remodelling (*BGAL10*)[76], fucosylation (*FUT1*) and no genes involved in O-acetylation (*AXY4*)[77] (Figure 4c, red arrows)[78]. Although the number of FUT1 homologs are fewer in ferns, the evolutionary history of GT37 sequences was shown to be more complex when it comes to substrate specificity [79]. For xylans, ferns showed a near absence of genes involved in methylation of glucuronic acid in glucuronoxylan (*IRX15*, *GXMT1-3* [74]) and xylan acetylation (*ESK1*)[80]. Homogalacturonan pectins showed fewer fern genes involved in xylogalacturonan synthesis (*XGD1*)[81] and galacturonan acetylation (*PMR5*)[82]. Finally, rhamnogalacturonan I pectins showed fewer fern genes involved in remodelling (*TBG4*)[83]. A similar analysis of hydroxyproline-rich glycoproteins did not reveal significant absences of these proteins in ferns (Figure S16). Overall, ferns tend to contain a lower number of genes per family and do not have the *DUF579* (sugar methylesterification)[84] and DUF231 (sugar acetylation)[85]

21

555    gene families (Figure 4c).

556         To directly compare the polysaccharide inventories of angiosperms and ferns, we

557    performed a Comprehensive Microarray Polymer Profiling (CoMPP)[86], where we probed

558    102 cell wall extracts from nine ferns from six fern orders, including 36 organs at different

559    developmental stages, with 48 antibodies targeting different cell wall epitopes The

560    epitopes recognised by the antibodies are given in Table S14.

561         Our analysis revealed the relative abundance of cell wall polysaccharides (Table

562    S15), which showed that different organs from the same species tend to have similar cell

563    wall composition (Figure S17a) and the polysaccharide profiles within species tend to be

564    more correlated than across species (Figure S17b).

565         We found that Eusporangiate ferns were generally richer in easily extracted

566    polymers than leptosporangiate ferns, with Adiantum as a notable outlier (Figure 4d).

567    Surprisingly, while mixed-linkage $\beta$-glucan (MLG) was only reported so far in Equisetum

568    [87,88], we observed a clear MLG signal outside of Equisetidae in *O. pendulum, Angiopteris*

569    *sp1, A. capillus-veneris* and *S. palustris* (Figure 4d). For pectins, we observed a high

570    abundance of homogalacturonan at different grades of methyl esterification (antibodies

571    CCRCM38, LM18, LM19, LM20, JIM5, JIM7), but low signal from RG-I backbone (INRA-

572    RU1, INRA-RU2), galactosylated RG-I (LM16) and arabinan (LM13). For hemicelluloses,

573    we observed a strong signal for xyloglucan (CCRC-M87), both fucosylated (CCRC-M102,

574    CCRC-M1, CCRC-M39) and non-fucosylated (CCRC-M50), suggesting that xyloglucan

575    might be a quantitatively important hemicellulose, which contrasts with a previous study

576    reporting mainly mannan-rich cell walls [89]. We also observed signals from xylan (CCRC-

577    M154, LM11, CCRC-M159, LM10, LM23) and galactomannan (CCRC-M175) and various

578    mannan-containing polysaccharides (LM22, LM21), with galactomannan showing signal

579    only in *E. palustre* (CCRC-M170, -M167). The weakest signals were observed for

580    epitopes in arabinogalactan-proteins (AGPs) and extensins, as only a few antibodies

581    gave moderate signals (JIM8, MAC207, JIM13). Other antibodies showed weak signals

582    (AGPs: JIM16, LM2, LM14, extensins: JIM20, JIM11). Finally, feruloylated

583    polysaccharides that crosslink with arabinan and galactan residues of cell wall pectin via

584    ester bonds [90] showed no signal (Figure 4d). Taken together, these results indicate that

585    fern and angiosperms share most of the polysaccharides and their biosynthetic enzymes,

586    but ferns might lack certain sugar modifications and AGP structures found in flowering

587    plants.

588

589    **Evolution of the cellulose synthase superfamily in Archaeplastida**

590    In addition to lignin, cellulose is one of the major load-bearing polymers. Angiosperms

591    contain primary and secondary cell walls enriched in cellulose, which are biosynthesised

592    by CESA1,3,6 and CESA4,7,8 in *Arabidopsis thaliana* [91]. The CESA complexes are

593    arranged in hexameric complexes called rosettes in angiosperms [92,93], or as linear

594    terminal complexes in bacteria [94,95]. *Selaginella* is the latest diverging plant known to

595    possess both CesA hexameric complexes and CesAs of the type that forms linear

596    complexes in bacteria [96]. The plant kingdom has also evolved cellulose synthase-like

597    (*CSL*) genes to produce other polysaccharides [97], such as mannans (*CSLA*)[98], glucan

598    chain of xyloglucan (*CSLC*)[99], cellulose in tip growing cells (*CSLD*)[100,101] and mixed-

599    linkage glucans (*CSLF*)[102]. However, the evolution of the CESA superfamily is not well

600    understood in ferns.

601         A phylogenetic analysis of the *CESA* superfamily built from algal and land plant

602    protein sequences showed that ferns contain both linear (*CESA* linear) and hexameric

603    (*CESA1/3/10* and *CESA6*) *CESA* genes (Figure 5a, Figure S18). The *CSLA*, *CSLC* and

604    *CSLD* families were found in all land plants, including ferns (Figure 5a). The *CSLB* and

605    *CSLG* families are only found in seed plants [103], but ferns contain one clade of genes that

606    is likely ancestral to the two families (Figure 5a, green arrow). We also observed that the

607    angiosperm secondary cell wall enzymes *CESA4,7,8* and the fern *CESAs* do not form a

608    monophyletic group (Figure 5a, black arrows), indicating that ferns either lack CESA4, 7,

609    8 or have evolved versions that no longer form clear clades with them. Conversely, ferns

610    form two distinctive groups with Arabidopsis *CESA6* and *CESA1,3,10* (Figure 5a),

611    suggesting that the ancestor of ferns and seed plants contained two *CESAs* that gave

612    rise to *CESA6-like* and *CESA1-like* clades. The phylogenetic tree revealed four

613    independent duplication events of the *CESAs* within ferns (Figure 5a, light blue arrows),

614    suggesting that ferns have likely evolved cell walls with properties distinctive from

615    flowering plants.

616

23

**Fern-specific evolution of secondary cellulose synthases**

To better understand the function of the duplicated CESAs ferns, we first analyzed the gene tree of two CESA clades (Figure 5a, red and blue arrow). Both clades contain genes from *Dicranopteris curranii* (red clade: *Dcu_g31359,* blue clade: *Dcu_g12277*, black arrows), suggesting a duplication in the ancestor of Gleicheniales-II (Figure 5b). To suggest the function of the red and blue clades, we first examined the expression profiles of two representative genes (*Ceric.13G049300.1*) and (*Ceric.09G024100.1*) from *Ceratopteris richardii* (Figure 5b, blue and red solid arrows, respectively). While *Ceric.13G049300.1* showed the highest expression in vegetative fronds, *Ceric.09G024100.1*'s expression was highest in leaf and shoot tips, suggesting different biological processes for the two genes (Figure 5c).

We next compared the co-expression network of *Ceric.13G049300.1* to other ferns using the CoNekT's Expression Context Conservation (ECC) panel (https://conekt.sbs.ntu.edu.sg/sequence/view/2353618). The fern gene with the most similar expression network was *Stenochlaena palustris Spa_g26805*, which happens to be found in the same clade (Figure 5b, blue solid and double arrow). The networks of *Ceric.13G049300.1 and Spa_g26805* contain orthologs involved in lignocellulose production, such as *CESAs, 4CLs, OMT1s, CYP98A3* (lignin-related *C3'H*), *CSI1* [104] and laccases [105](Figure 5d). Thus, the genes from the blue clade are likely involved in secondary well wall biosyntesis, indicating that ferns independently evolved this module.

Conversely, *Ceric.09G024100.1* and its most similar co-expression ortholog was *Pyrrosia piloselloides Ppi_g22229* (Figure 5b, green solid and double arrow) were co-expressed with genes unrelated to lignocellulose production (e.g., genes similar to monoterpenol-associated *CYP76C2* [106])(Figure 5d). This suggests that the second *CESA* clade might be involved in another unknown biological process. Taken together, this indicated that ferns have independently evolved a secondary cell wall module, and further duplicated the CEASs to perform yet unknown functions.
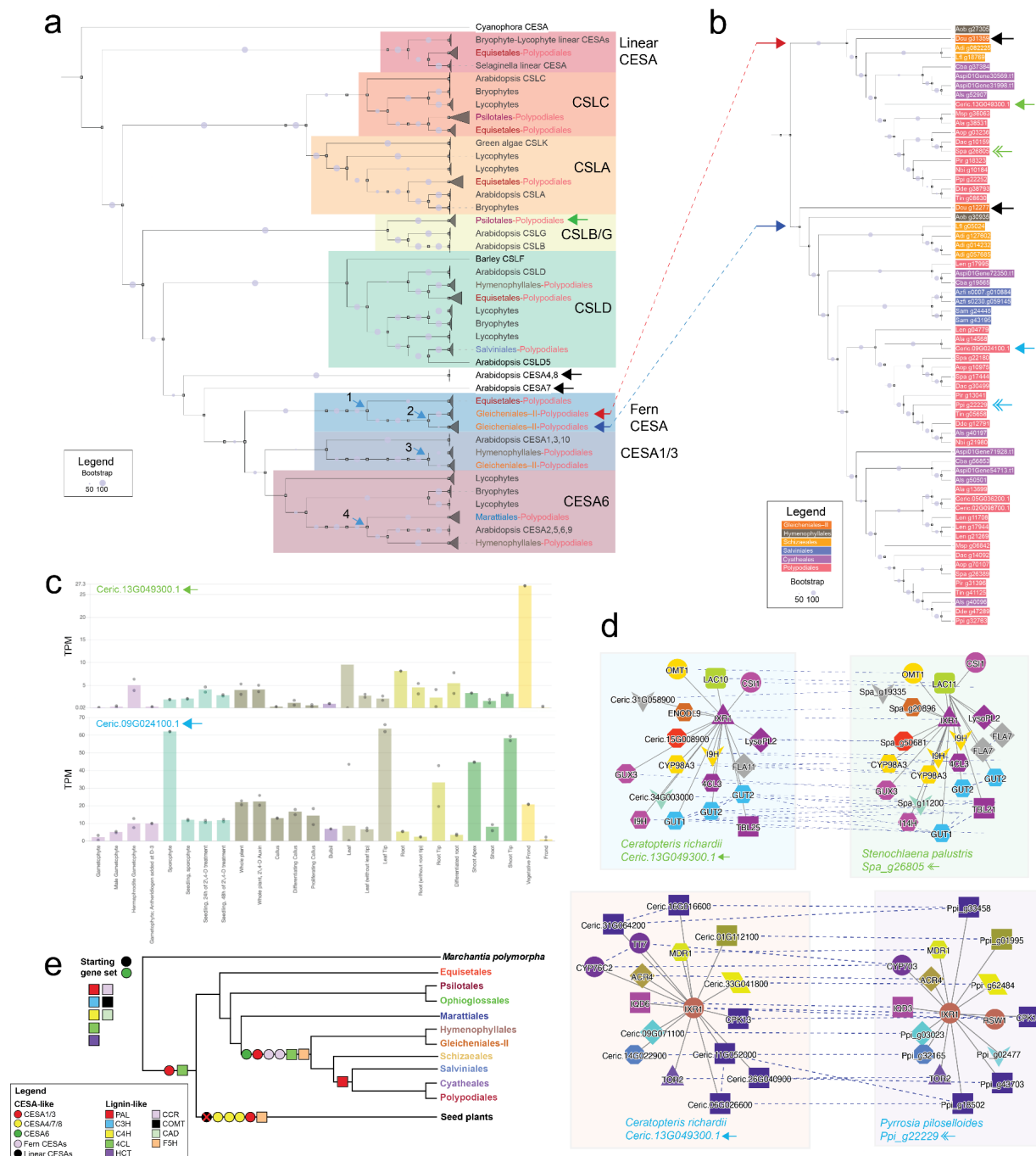
644

**Figure 5. Independent duplication of cell wall-related modules in ferns.** a) Maximum likelihood cellulose synthase superfamily gene tree. Gray triangles indicate collapsed clades. The arrows indicate the discussed clades and duplication events. The grey circles indicate the bootstrap support of the nodes, while the colored text represents the different fern clades. b) Closeup on the fern-specific clades. The fern orders are colour-coded, while the blue and green arrows indicate the discussed genes. c) CoNekT expression profiles of two CESA genes from *Ceratopteris richardii*, where the x-axis indicates organs and tissues, and the y-axis represents transcripts per million (TPM) values. d) CoNekT comparative co-expression analysis of the four CESA genes indicated in panel b. Nodes

25

654 represent genes, while solid and dashed edges connect co-expressed and orthologous
655 genes. Coloured shapes indicate the different orthogroups. e) Order tree summarising
656 the duplication events of cell wall-related genes. The tree is based on the gene tree of
657 CESAs and lignin-related genes. Coloured shapes represent the different gene classes.
658

659 **The evolution of lignocellulose-biosynthesizing genes in land plants**

660 To better understand how the genes involved in lignin and cellulose synthesis have
661 evolved in land plants, we mapped the timing of gene duplications onto the land plant
662 species tree (Figure 5e). Because both *CESA6-like* and linear *CESA* clades both contain
663 bryophytes (Figure 5a), we propose that the ancestor of land plants contained a *CESA6-*
664 *like* and a *CESA* of linear-type. The ancestor of ferns and seed plants evolved *CESA1/3-*
665 *like* genes, that further expanded in seed plants into *CESA1* and *3* and secondary
666 *CESAs4,7,8* (Figure 5a). Within ferns, we observed duplications of *CESA1/3-like* in the
667 ancestor of Gleicheniales-II (Figure 5a, duplication 3), *CESA6-like* in Hymenophyllales
668 (duplication 4), and two duplications of the fern-specific CESAs (duplications 1,2) in
669 Gleicheniales-II. We also observed a complete gene set of lignin biosynthetic genes in
670 early-diverging land plants, and evidence of duplication of *4CL* in the ancestor of ferns
671 and seed plants and the ancestor of Hymenophyllales (Figure S19).

672 Taken together, the ferns show a prolific duplication of *CESA* genes deeply within
673 the fern lineage (Figure 5e), further suggesting that ferns have evolved cell walls with yet
674 unknown features.

675

676 **Discussion**

677 Despite ferns' critical evolutionary position as the sister group, no large-scale studies that
678 investigated their phylogeny, biological pathways and cell walls had been performed. To
679 remedy this, we generated gene expression atlases for 22 ferns and covered ten out of
680 12 fern orders (Figure 1), allowing us to generate a high-quality species tree that resolved
681 the long-standing relationship between Gleicheniales and Hymenophyllales [12,107,108]. The
682 tree is supported by outgroups comprising lycophytes, horsetails or seed plants, but not
683 bryophytes. We speculate that the greater phylogenetic distance between bryophytes and
684 the other lineages, combined with reduced single-copy gene dataset obtained from
685 OrthoFinder (Supplemental Methods 1), and degenerated phylogenetic signals in amino

686   acid sequences contributed to the discordance between bryophyte-based and the other
687   outgroups.

688         The species tree of ferns allowed us to estimate the time of speciation and whole
689   genome duplication events. The WGD analysis revealed that WGD events likely
690   contributed to species diversity (Figure S3), but we observed no correlation between the
691   number of WGDs and genome size. This suggests that alternative evolutionary sources
692   contribute to the exceptional genome size of ferns, and that recent transposon activities
693   and ploidy variation might play a bigger role.

694         The stem age of early diverging ferns Equisetidae (consisting of Equisetales), with
695   95% HPD time estimates are in line with the estimate from [109] and the oldest unequivocal
696   euphyllophyte fossils [1]. Polypodiidae was originated in the time range between Lower
697   Carboniferous (323.2 ± 0.4 - 358.9 ± 0.4 mya) and Middle Devonian (382.7 ± 1.6 - 393.3
698   ± 1.2 mya), with 95% HPD time estimate as 345.09 - 389.61 mya and posterior mean
699   369.01 mya, which might have first survived the Hangenberg and Kellwasser extinction
700   events before its substantial diversification. The early diverging leptosporangiate fern
701   order Osmundales originated amid Upper to Lower Carboniferous (298.9 ± 0.15 - 358.9
702   ± 0.4 mya), consistent with [13]. The aquatic Salviniales, the only extant ferns with
703   heterospory, was originated between Upper Triassic (201.4 ± 0.2 - 237 mya) and Lower
704   Permian (273.01 ± 0.14 - 298.9 ± 0.15 mya), with 95% HPD time estimate as 211.15 -
705   273.19 mya and posterior mean 241.70 mya, which might correlate with the P-T event
706   after which a vast majority of aquatic environment became empty and the innovative
707   microspores might facilitate their spread and survival. The two most species-diverse
708   suborders, Polypodiineae (i.e., eupolypods I) and Aspleniineae (i.e., eupolypods II) of the
709   Polypodiales were originated amid the Cretaceous (66.0 - 145.0 mya), a relatively warm
710   and ice-free period, with 95% HPD time estimate as 89.75 -144.92 mya and posterior
711   mean 115.93 mya, coincident with the burst of angiosperms in the mid-Cretaceous as
712   highlighted by Darwin [110] and the decline of gymnosperms [111].

713         Our gene inventory analysis shows massive gains of genes in the fern lineage
714   (Figure 2a), resulting in ~50% being fern-specific (Figure 2b). Expression analysis
715   revealed that the fern-specific genes tend to be organ-specific, suggesting their role in
716   fern-specific adaptations. Conversely, older genes are ubiquitously expressed (Figure

27

717   2c), which aligns with our previous observation that these genes tend to have basal,
718   essential functions (e.g., photosynthesis, protein synthesis, DNA duplication)[47]. Many of
719   the genes involved in angiosperms' hormonal and developmental pathways were missing
720   (Table S9), showing that ferns have organised these pathways differently.

721   Signalling and biosynthetic pathways may significantly vary within land plants, and
722   these pathways tend to expand to support increased anatomical and lifestyle complexity
723   [112].Thus, the arguably simpler fern hormonal pathway genes might suggest that these
724   pathways can function in ferns without their angiosperm counterparts. Alternatively, ferns
725   might have evolved equally complex but alternative signalling pathway components that
726   show no sequence similarity to known angiosperm genes. This idea is exemplified by our
727   analysis of flavonoid biosynthesis genes. The lack of these enzymes in ferns–but the
728   presence of flavonoids–indicates that the 'canonical' flavonoid pathway is an angiosperm-
729   specific invention and suggests that ferns have either convergently evolved other
730   enzymes with similar functions or use a different pathway to synthesise these
731   compounds. As fern-specific genes are co-expressed with genes involved in
732   development, reproduction and various signalling pathways (Figure 2e), ferns likely have
733   independently expanded these pathways.

734   The observed high amounts of lignin S units in *P. irregularis* and *S. palustris*
735   (Figure 3ab), and the absence of angiosperm- or lycophyte-specific *F5H* enzymes
736   suggest that the S unit has independently evolved at least four times in the plant lineage:
737   angiosperms, lycophytes, gymnosperms and now ferns [113]. The re-emergence of S lignin
738   in distantly related plant lineages implies that it may have an essential role in plants'
739   environmental adaptation, such as improved mechanical properties or herbivore
740   resistance [113]. While lycophytes have evolved a C3'H-independent pathway by inventing
741   a dual meta-hydroxylase *SmF5H* (Figure 3c, blue pathway)[69], we observed the presence
742   of *C3'H* genes in ferns (Figure 3d), suggesting that ferns have independently evolved a
743   *F5H* enzyme, and likely follow the biosynthetic route of angiosperms. By combining
744   phylogenetic and co-expression analysis, we propose that the red clade shown in Figure
745   3h, which contains the highest density of *CYP450s* co-expressed with the lignin
746   biosynthetic genes, comprises the fern *F5H* enzymes. The biosynthetic activity of these
747   genes could be tested by in vitro studies, as done for Selaginella Sm*F5H* [69].

748      Our comparative analysis of cell wall-related genes indicated that ferns and

749    angiosperms contain similar gene sets but that ferns have smaller acetyl and methyl-

750    transferase gene families (Figure 4c). Cell wall composition varies considerably between

751    fern species, corroborating findings in earlier glycan array surveys of ferns [114].

752    Surprisingly, we observed a clear mixed-linkage $\beta$-glucan signal from *O. pendulum*,

753    *Angiopteris sp1*, *A. capillus-veneris* and *S. palustris* (Figure 4d), demonstrating that this

754    unusual polymer is found outside of fern Equisetum [87,88]. While the AGP epitopes showed

755    weak signals, ferns contain AGPs with special features, such as 3-*O*-methylrhamnose,

756    that are not known in angiosperms [73,115].

757      Surprisingly, we observed a wide-spread occurrence of mixed-linkage glucans

758    outside of Equisetidae (Figure 4d), and we propose two candidate enzyme families that

759    could produce this hemicellulose. First, bryophytes, ferns and Selaginella both contain

760    MLGs [96] and linear CesAs (Figure 5a). In the moss *Physcomitrium*, linear CesAs produce

761    arabinoglucan [116], and the authors point out that these CESAs are related to an

762    ascomycete MLG synthase and thus represent an early system for MLG-synthesis. A

763    second candidate could be the fern CSLB/G-related clade, as the functions of these

764    genes are currently unknown in ferns and angiosperms. The biosynthetic activity of these

765    enzymes could be tested in vivo, as done for the barley MLG synthase [117].

766      Separate sets of CesAs for primary and secondary cell wall synthesis are a shared

767    feature of spermatophytes [91]. Whereas tracheids probably evolved once uniting all

768    tracheophytes [118], vessels have evolved in angiosperms and independently in

769    *Selaginella, the Gnetales, Equisetum* and other ferns [119]. Surprisingly, *Selaginella* does

770    not have two sets of CesAs, and we did not observe CesA related to angiosperm

771    secondary cell walls in ferns (Figure 5a), suggesting that their vascular elements result

772    from convergent evolution [96]. Convergent evolution is supported by a fern-specific CesA

773    clade containing genes involved in lignocellulose biosynthesis (Figure 5). We observed

774    several duplications of the cell wall-related genes within ferns (Figure 5e), which aligns

775    with similar observations in angiosperms and bryophytes [120,121]. Combined with the

776    presence of a non-canonical sugar 2-O-Methyl-D-glucopyranose observed in

777    Polypodiales (Figure 4a), our data suggest that cell walls underwent independent

778    innovations within ferns. While it is unclear whether 2-*O*-Methyl-D-glucopyranose is

779   biosynthesized by ferns or bacterial or fungal organisms found in the environment, their

780   significant presence in the fern cell walls indicates that they might have a role in fern

781   biology.

782   We anticipate that the availability of the comprehensive coding sequence and

783   transcriptomic data from ferns - and their availability as a user-friendly CoNekT database

784   (https://conekt.sbs.ntu.edu.sg/)-will be mined to lead to vital insights into the evolution of

785   plant genes and gene families. Implementing fern data into the existing comparative

786   genomic framework will enhance our understanding of the plant tree of life.

787

788   **Methods**

789   **Sampling of ferns**

790   22 ferns from 22 families were sampled across Singapore (Table S1). Fern organs were

791   sampled as three biological replicates, where organs were selected to capture the highest

792   variance in the developmental stages and morphological characteristics (Figure S1).

793   Samples were placed into 15 ml falcon tubes and kept in liquid nitrogen and subsequently

794   at -80°C to prevent degradation of RNA.

795

796   **RNA isolation and sequencing**

797   After collection, each sample was ground in liquid nitrogen to a fine powder. RNA was

798   extracted from 100 mg of plant material using Spectrum™ Total Plant RNA Kit (Sigma)

799   Protocol A following the manufacturer's instructions. Quality control of all extracted RNA

800   was carried out by Novogene (Singapore). Each sample was evaluated for its quantity,

801   integrity and purity using agarose electrophoresis and Nanodrop. Library construction

802   was performed by Novogene, and mRNA was enriched from total RNA with oligo(-dT)

803   magnetic beads. The library was then quantified with Qubit and real-time PCR and

804   sequenced using Illumina NovaSeq 6000, with paired-end sequencing of 150 base pairs

805   (bp) per read and a sequencing depth of approximately 60-70 million reads.

806

807   **Transcriptome assembly**

808   Low-quality RNA-seq reads were removed, and the remaining reads were trimmed with

809   Fastp (v0.23.2)[122]. Reads were assembled via a curated transcriptome assembly pipeline

810    (Figure S2). Reads were assembled in three biological replicates for each organ, with all

811    organs concatenated and filtered. Each organ consisting of three reads was assembled

812    using SOAPdenovo-Trans (v1.03)[37] with 10 single *K-mer* (21-39) and Trinity (v2.8.5)[36]

813    with 25 single *K-mer*. All reads were concatenated into a Trinity-SOAPdenovo assembly

814    and    filtered    through    the    Evidential    Gene    Pipeline

815    (http://arthropods.eugenes.org/EvidentialGene/) using trformat.pl and tr2aacds.pl for

816    removal of any redundant coding sequences. The filtered transcriptome was evaluated

817    using BUSCO (v5.4.3)[123] and embryophyta as the dataset. All organ assemblies within

818    the sampled fern species were concatenated and filtered again through the

819    EvidentialGene Pipeline, with the output transcriptome being filtered via transcripts per

820    million (TPM) reads using kallisto (v0.50.1)[124] against each RNA-seq. TPM scores were

821    averaged per organ, and coding sequences with scores < 1 were removed from the

822    assembly. Assembly was then filtered using GC% content, with redundancy coding

823    sequences with less than 40% and more than 60% removed. All transcriptomes were

824    blasted against the NCBI database, and sequences with a percentage identity of more

825    than 70% and an e-value less than e$^{-10}$ were removed from the final assembly. The quality

826    of assembly was determined by BUSCO using the embryophyta dataset (Table S2).

827

### Construction of $K_S$-based age distributions

829    $K_S$-age distributions for all paralogous genes (paranome) of genomes and transcriptomes

830    were constructed by ksrates (v1.1.1)[125]. In brief, the ksrates pipeline entails firstly

831    translating the coding nucleotide sequences into peptide sequences assuming standard

832    genetic code, filtering out sequences whose sequence length is not divisible by 3,

833    containing invalid codons or in-frame stop codon, after which an all-versus-all BLASTP

834    was implemented with *E*-value set as 1 x 10$^{-10}$ in BLASTP (v2.11.0+)[126] and the resultant

835    subject-query hit table was fed into MCL (v14-137) [127] with clustering inflation factor set

836    as 3.0 to delineate paralogous gene families while filtering out gene families whose size

837    is larger than 200, secondly calling the aligner MUSCLE (v3.8.1551)[128] under default

838    parameter to obtain a multiple sequence alignment (MSA) at the protein level for each

839    paralogous gene family while filtering out sequence pairs whose gap-stripped alignment

840    length was shorter than 100, which was then back-translated into a codon alignment and

31

841  subsequently fed into the CODEML function within PAML (v4.9j) [129] to acquire the
842  maximum likelihood estimate (MLE) of $K_S$ values under non-pairwise mode using the
843  default control file defined by wgd (v1.1.1) [130] and then calling FastTree (v2.1.11)[131] upon
844  the peptide MSA under default parameter to attain a midpoint-rooted phylogenetic tree of
845  each paralogous gene family for retrieving the weight of each paralogous gene pair with
846  or without outliers, and eventually building the $K_S$-age distribution with de-redundancy
847  achieved by node-weighted method after excluding outliers. The collinear gene pairs
848  (anchor pairs) were identified by i-ADHoRe (v3.0.01)[132] under the default control file
849  defined by wgd, and the weight values for anchor pairs whose corresponding $K_S$ values
850  were between 0.05 and 20 were recalculated and reassigned while the weight of
851  remaining pairs was set as zero. For orthologous $K_S$-age distributions, the process of
852  MCL clustering was superseded as reciprocal best hits (RBH) searching to identify
853  orthologous gene pairs while the weighting process was revoked on that only one-versus-
854  one orthologues were inferred. CD-HIT (v4.8.1)[133] was applied for the de-redundancy of
855  transcriptome assemblies with the clustering threshold set as 0.99 before $K_S$-age
856  analysis.

857

858  **Correction of differences of synonymous substitution rates**
859  Synonymous substitution rates were corrected in ksrates. The principle leans on a
860  number of trios of species, including a focal species, a sister species and an outgroup
861  species. The disparate synonymous substitution rate between focal species and sister
862  species since the divergence is represented indeed by the branch-specific contribution of
863  accumulated synonymous substitutions per synonymous site in respective branches. The
864  mode of orthologous $K_S$-age distributions inferred from the kernel density estimate (KDE)
865  using Gaussian kernels within the python package scipy was designated as the proxy of
866  the peak $K_S$ value of each orthologous $K_S$-age distribution. 200 iterations of bootstrap with
867  replacements were implemented for each orthologous $K_S$-age distribution, and the mean
868  along with standard deviations (STD) of mode across the replicates was determined as
869  the final peak $K_S$ value representing divergence distance and its associated STD. The
870  original accumulated synonymous substitutions per synonymous site of focal species-
871  sister species pair consisting of the branch-specific contribution of both species since

872 diversification was transformed into two times the branch-specific contribution of focal
873 species with the prop of outgroup species to resemble the timescale of focal species. The
874 mean of rescaled peak $K_S$ values of focal species-sister species pair against various
875 outgroup species was taken as consensus-adjusted peak $K_S$ value. The maximum
876 number of trios was set as 20. The species tree inferred by ASTRAL-Pro2 using seed
877 plants as outgroup species were adopted in ksrates.

878

879 **Construction of orthologous families and single-copy gene trees**

880 Orthofinder (v2.5.4)[134] was performed upon the protein sequences of 107 ferns and
881 outgroup species with an inflation factor set as 3 to delineate the orthologous families. No
882 single-copy gene families were identified by Orthofinder, probably because of the
883 universal and unique gene duplication and loss scenario across species and gene
884 isoforms [135]. To recover reliable and adequate single-copy gene families, we constructed
885 mostly single-copy gene families [51], wherein most species were in single-copy while the
886 remaining species had no more than four copies which were assumed to be transcript
887 variants of the same gene, by retaining the longest copy, if applied, of each species. We
888 referred to the mostly single-copy gene families as single-copy gene families thereafter.
889 In total, 140, 112, 107 and 57 single-copy gene families were constructed from dataset
890 107 ferns, 107 ferns plus seed plants, 107 ferns plus seed plants and lycophytes, 107
891 ferns plus seed plants, lycophytes and bryophytes, respectively. MAFFT (v7.475)[136] was
892 performed to obtain a peptide multiple sequence alignment (MSA) for each single-copy
893 gene family with the parameter "–auto". Trimal (v1.4.1)[137] was then performed to trim the
894 MSA and back-translate it into a codon-level nucleotide MSA with parameter "-
895 automated1". IQ-TREE (v1.6.12)[138] was implemented on each codon-level nucleotide
896 MSA wherein ModelFinder [139] was called to find the best-fit codon substitution model in
897 terms of Bayesian Information Criterion (BIC) upon which a maximum likelihood (ML)
898 gene tree was inferred and assigned with bootstrap support values from 1000 ultrafast [140]
899 bootstrap replicates with parameter "-bnni" to further optimize each bootstrap iteration
900 through a hill-climbing nearest neighbor interchange (NNI) search based directly on the
901 corresponding bootstrap alignment to avoid severe model violations. The same process
902 was further applied to the 108 ferns dataset including *Marsilea vestita*, in which a total of

903    136, 108, 103 and 55 single-copy gene families were reconstructed from datasets varied

904    in outgroups with both nucleotide and peptide molecules.

905

**Species tree inference**

907    Three methods, ASTRAL-Pro2 [41], STAG [42], and a concatenated-based method were

908    implemented to infer the species tree. The acquired individual ML single-copy gene trees

909    and gene name-species name map files were imported into ASTRAL-Pro2 and STAG

910    under default parameters to estimate a consensus species tree with support values for

911    each bipartition denoting local posterior probabilities (localPP) and the proportion of

912    individual estimates of the species tree that contain that bipartition, respectively.  For

913    the concatenated-based method, the individual codon-level nucleotide MSA of single-

914    copy gene families were concatenated and then fed into IQ-TREE to infer a ML super-

915    gene tree as above. *Dicranopteris curranii*, *Dicranopteris pedata*, *Diplopterygium*

916    *laevissimum*, *Diplopterygium glaucum* and *Sticherus truncatus* in the Gleicheniaceae

917    were named as Gleicheniales−II clade while *Cheiropleuria integrifolia*, *Dipteris conjugata*

918    and *Dipteris lobbiana* in the Dipteridaceae were named as Gleicheniales−I clade p. In

919    total, (140), (112), (107) and (57) single-copy gene families were constructed from the

920    (107 ferns dataset), (107 ferns plus seed plants), (107 ferns plus seed plants and

921    Lycopods), (107 ferns plus seed plants, Lycopods and Bryophytes), respectively. The 108

922    ferns dataset including *Marsilea vestita*, were with 136, 108, 103 and 55 single-copy gene

923    families, respectively.

924

**Estimation of absolute divergence time**

926    Mcmctree (v4.9j)[129] was implemented upon the concatenated peptide MSA of single-copy

927    gene families of 108 ferns dataset with Equisetales as outgroup to infer the absolute

928    divergence time for each bipartition. The independent rates model, which assumes a log-

929    normal distribution of evolutionary rates across branches, was selected and 18 fossil

930    calibrations of soft constraint from [141] were adopted to refine the divergence time of

931    internal nodes, as summarised in Table S2. Fossils calibrating clades within

932    Gleicheniaceae or Hymenophyllaceae were avoided for their indefinite phylogenetic

34

933    location. LG amino acid substitution matrix was selected and a gamma model of rate

934    variation was assumed with alpha as 0.5 and 5 categories in discrete gamma. Parameters

935    controlling the birth-death process were set as 1, 1, 0.1 to generate uniform age priors on

936    nodes that didn't have a fossil calibration. Gamma priors for the transition/transversion

937    rate ratio and shape parameters for variable rates among sites were set as 6 2 and 1 1.

938    A Dirichlet-gamma prior was set upon the mean rate across loci and the variance in

939    logarithm as 2 20 1 and 1 10 1. The first 2000 iterations were discarded as burn-in and

940    then 20,000,000 iterations were performed with sampling per 1000 iterations. The

941    effective sample size (ESS) of all parameters was larger than 200, suggesting adequate

942    sampling and convergence.

943

944    **Phylogenomic analysis of gene tree - species tree reconciliation**

945    To estimate the retention rate and interrogate hypothetical WGDs over competing

946    scenarios in different clades, we implemented 4 categories of statistical gene tree -

947    species tree reconciliation analysis using Whale (v.2.0.3)[142], as shown in Figure SM9-12.

948    Firstly, orthogroups of each category of species were inferred by Orthofinder (v2.5.4) [134]

949    with an inflation factor as 3. Gene families were filtered to assure at least one gene from

950    each descendant present at the root and to avoid large gene family size(which contains

951    noise and causes computational downshift) via "orthofilter.py"

952    (https://github.com/arzwa/Whale.jl)1000 gene families were randomly selected as

953    subsequent inputs. PRANK (v.150803)[143] was utilised to obtain a MSA for each gene

954    family and MrBayes (v.3.2.6)[144] was then applied to infer the posterior distributions of

955    gene trees under the LG + GAMMA model, with iterations set as 110,000 and sample

956    frequency as 10 to get in total 11,000 posterior samples. ALEobserve [145] was

957    subsequently performed on the tree samples to construct the conditional clade distribution

958    with a burn-in of 1000. Two gene family evolution models, the relaxed branch-specific

959    model and the critical branch-specific model, were applied as in a previous study [45] to

960    estimate the retention rates of hypothetical WGDs for each category, as shown in Figure

961    SM9-12. Hypothetical WGDs with retention rates higher than 0.05 were regarded as

962    supported WGDs, considering the incompleteness of transcriptome assemblies and the

963    stochasticity of sampled gene families.

964

**Absolute dating of WGDs**

Phylogenetic dating of AZOL and CERA WGD proceeded as follows. Firstly, an orthogroup comprising orthologues from 8 other species and anchor pair which was assumed to be retained from the corresponding WGD was constructed per anchor pair by searching the reciprocal best hits (RBH) between the anchor pair and the transcriptomes or genomes of other species by Diamond (v2.0.5.143)[146] under default parameters (Figure SM14). $K_S$ range 0.36 - 2.00 was confined for the age of anchor pairs to be adopted in terms of densest aggregation of duplicates and avoidance towards saturation for AZOL WGD and $K_S$ range 0.41 - 2.0 was bounded for CERA WGD. Secondly, the peptide sequences of each individual orthogroup were aligned by MAFFT (v7.475)[136] under default parameters and then concatenated as a single peptide MSA. The numbers of concatenated orthogroups were 45 and 14 for AZOL and CERA WGD, respectively. The adopted fossil calibrations followed Table S2 at corresponding phylogenetic locations while the boundaries of root for AZOL WGD were set as minimum bound 168 mya based on the minimum bound of fossil calibration "Stem Lygodiaceae" and safe maximum bound 345 mya as the fossil calibration "Stem Osmundaceae", as shown in (Figure SM14). Mcmctree (v4.9j)[129] was implemented for the Bayesian molecular dating for each WGD with the parameters same as above. The ESS of all parameters was larger than 200, indicating adequate sampling and convergence. The posterior distribution of time estimate for the node joining the anchor pair was retrieved and the 95% HPD, posterior mean, median and mode were adopted to characterise the age of WGD, as shown in (Figure SM13).

**Identifying organ-specific genes**

Organ-specific genes were isolated from each transcriptome via specificity measure (SPM) values [47]. For each gene, we calculated the average TPM values in each organ. Following that, the SPM value of a gene in an organ was calculated by dividing the average TPM in the organ by the sum of the average TPM values of all organs. The SPM value ranges from 0 (gene not expressed in organ) to 1 (gene fully organ-specific). To identify organ-specific genes for each organ, we first identified an SPM value threshold

995 above which the top 5-11% of SPM values were found (Figure S5). These top values

996 varied across the species sampled, depending on the number of organ-specific genes

997 identified. If the SPM value of a gene in an organ was equal or greater than the threshold

998 value, the gene was identified as organ-specific within said organ. Organ-specific genes

999 were then plotted in a heat map to show their distributions (Figure S6).

1000

**Functional annotation of genes**

1001

1002 Assembled sequences of 22 fern species, including four ferns with genomes available

1003 online (*A. spinulosa, A. filiculoides, C. richardii* and *S. cucullata*) were annotated using

1004 the online tool Mercator4 v.2.0 [147]. We visualised the Mercator4 annotation using a

1005 heatmap, showing the distribution of Mapman Bins across sampled fern species (Figure

1006 S9).

1007

**Assignment of orthogroups to phylostrata**

1008

1009 Using the coding sequences of the transcriptomes, we constructed orthologous gene

1010 groups (Orthogroups) with Orthofinder (v.2.5.4)[134]. Respective outputs for orthogroups

1011 were used for further analysis. By utilising the theoretical evolutionary line produced by

1012 the phylogenomic analysis of gene trees, phylostratic nodes were assigned to

1013 orthogroups based on plant lineages. This analysis spanned a total of 47 species across

1014 the plant kingdom, and assigned nodes ranging from node 1 (most ancient, ancestor of

1015 Archaeplastida) to node 13 (ancestor of Polypodiales). Nodes were assigned based on

1016 the fern species tree (Figure 1c), as well as known phylogenetic analyses of early plants

1017 [148]. The nodes are: node 1 (ancestor of Archaeplastida), node 2 (ancestor of green

1018 plants), node 3 (ancestor of land plants), node 4 (ancestor of vascular plants), node 5

1019 (ancestor of ferns and seed plants), and node 6-13 (various fern orders). Specifically,

1020 Equisetales is designated as node 6, Psilotales and Ophioglossales as node 7,

1021 Marattiales as node 8, Hymenophyllales and Gleicheniales-II as node 9, Schizaeales as

1022 node 10, Salviniales as node 11, Cyatheales as node 12, and Polypodiales as node 13.

1023 Species-specific gene families were characterised by gene families consisting of only one

1024 species, and hence, not assigned to nodes. In cases where nodes encompass multiple

1025 species, such as node 4, orthogroups containing only one node assignment (e.g., those

1026 with genes solely from Ophioglossales and Psilotales) were not designated to specific
1027 nodes.

1028

1029 **Orthogroup gain loss analysis**

1030 Gain and loss of orthogroups were determined by the presence of an oldest clade
1031 member in a particular node. Potential contamination by non-fern sequences due to the
1032 nature of transcriptome assembly was filtered out at this stage by checking for the
1033 presence of at least half of the expected clades in each node. For basal nodes (nodes 1
1034 to 4), the clades used were 'Glaucophytes', 'Chlorophytes', 'Bryophytes', 'Lycophytes',
1035 'Tracheophyta' and 'Spermatophyta'). Nodes were defined as lost based on the clade that
1036 they last appeared in.

1037

1038 **Identification of orthogroup expression profiles**

1039 Analysis of the expression profiles at phylostrata level was performed as in [47], by
1040 classifying orthogroups into 'organ-specific', 'ubiquitous' or 'not conserved'. Organ-
1041 specific orthogroups are orthogroups containing organ-specific genes and were
1042 subclassified according to the organ (leaf-, meristem-, crozier-, root meristem-, male-,
1043 spore-, rhizome-, root-specific). Orthogroups that are expressed in different organs for
1044 each species - that is, that do not show an 'organ-specific' expression profile in different
1045 species - were labelled as ubiquitous. Orthogroups that had different organ-specific
1046 expression profiles in different species (orthogroups containing root-specific genes for
1047 *Alsophila latebrosa* and leaf-specific genes for *Equisetum hyemale*) were labelled as not
1048 conserved. Only orthogroups that fulfilled the following criteria were identified as organ-
1049 specific: (1) Contained at least two species with transcriptome data within each
1050 orthogroup. (2) >50% of the genes within the orthogroup supported the expression profile
1051 and (3) ≥50% of the species present in the node supported the expression profile.

1052

1053 **Structural analysis of fern-specific orthogroups**

1054 25 orthogroups containing at least 10 fern species with protein sequence representatives

1055 from sequenced genomes were used to check for sequence similarity by NCBI BLASTp

1056 restricted to Viridiplantae (E-value < 1e-10, Query cover > 50%), prediction of structure

1057 (alphafold 3 -server)[149], structure similarity search using DALI (all PDB, Z score > 8, lali >

1058 0.5 of residues in the protein)[150] and foldseek [151]. The cif outputs from alphafold 3 were

1059 converted to pdb format for input to DALI and foldseek using UCSF ChimeraX version:

1060 1.8. The sequenced representatives were selected based on the highest similarity to the

1061 consensus sequence, which was derived from the multiple sequence alignment

1062 generated using Seaview v5.0.5 (-align -align_algo 1 -output_format clustal -o ) using the

1063 muscle algorithm [152].

1064

1065 **Constructing co-expression networks and addition of the ferns to the CoNekT**

1066 **database**

1067 Co-expression networks were calculated using the CoNekT framework [67], and were also

1068 used to update the existing database, available at (https://conekt.sbs.ntu.edu.sg/).

1069

1070 **CASA lignin quantification**

1071 Methods used for solvent extraction and determination of lignin content by CASA lignin

1072 method closely followed the protocol outlined in [64]. Species organs that were sampled

1073 were ground, with solvent extraction in 80% ethanol for woody samples. For non-woody

1074 samples, extraction with water using sonication was done first to remove proteins and

1075 other water-soluble components. A cysteine stock solution (0.1 g/mL) in 72% sulfuric acid

1076 (SA) was prepared by dissolving 10 g L-cysteine in 100 ml SA. 5-10 mg of the solvent

1077 extract was placed in a glass vial, where 1.0 mL of prepared stock solution was added,

1078 sealed with a Telfon-lined screw cap and stirred at 24 °C (room temperature) via a

1079 magnetic stir bar (400 rpm) for 60 mins until the biomass was completely dissolved. The

1080 dissolving temperature was decreased to 24 °C to identify a milder condition, allowing

1081 convenient operation and minimising interference from carbohydrates. The solution was

1082 diluted with deionized water to a volume of 50 or 100mL in a volumetric flask, depending

1083 on the lignin content and biomass weight used. Absorbance of the diluted solution was

1084    measured at 283 nm ($A_{283}$) in a 1 cm quartz cell using a UV spectrophotometer against a

1085    blank solution (1 mL stock solution diluted to corresponding volume).

1086

1087    **Thioacidolysis**

1088    This method is adapted from [66]. Briefly, 10 mg of alcohol-insoluble cell wall residues were

1089    incubated in 3 mL of dioxane with ethanethiol (10%), BF3 etherate (2.5%) containing

1090    0.1% of heneicosane C21 diluted in $CH_2Cl_2$ at 100 °C during 4 hr. Three ml of $NaHCO_3$

1091    (0.2 M) were added after cooling and mixed prior to the addition of 0.1 mL of HCl (6 M).

1092    The tubes were vortexed after addition of 3 mL of dichloromethane and the lower organic

1093    phase collected in a new tube before concentration under nitrogen atmosphere to

1094    approximately 0.5 ml. Then, 10 μL of the mixture was trimethylsilylated (TMS) with 100

1095    μL of N,O- bis(trimethylsilyl) trifluoroacetamide and 10 μL of ACS- grade pyridine. The

1096    trimethylsilylated samples were injected (1 μL) onto an Agilent 5973 Gas

1097    Chromatography–Mass Spectrometry system. Specific ion chromatograms reconstructed

1098    at m/z 239, 269 and 299 were used to quantify H, G and S lignin monomers respectively

1099    and compared to the internal standard at m/z 57, 71, 85.

1100

1101    **Neutral sugar content analysis**

1102    This method is adapted from [66]. Neutral monosaccharide sugar content was determined

1103    by gas chromatography after acid hydrolysis and conversion of monomers into alditol

1104    acetates as described in Hoebler et al., 1989, and Blakeney et al., 1983. Gas

1105    chromatography was performed on a DB 225 capillary column (J&W Scientific, Folsorn,

1106    CA, USA; temperature 205 °C, carrier gas $H_2$). Calibration was made with standard sugar

1107    solution and inositol as internal standard.

1108

1109    **Synthesis of methylated sugars (Figure S14)**

1110    Please see supplemental methods. The structures of compounds were ascertained by

1111    NMR spectroscopy and were in agreement with reported data.

1112

1113    **Constructing phylogenetic trees for genes controlling primary and secondary cell**

1114    **wall formation**

1115 Genes controlling primary and secondary cell wall formation were identified from a
1116 previous study of *A. thaliana* [120]. Using the known genes as reference, we utilised
1117 OrthoFinder (v.2.5.4)[134] with 49 species within Table S9. Genes that were grouped in the
1118 same orthogroups with the reference genes were aligned via MUSCLE (v.5.1)[128] and
1119 analysed via IQTree (v.1.6.12)[138] to construct trees with bootstrap values of 100 (for gene
1120 families with more than 600 genes) and 1000 (for those less than 600 genes). Trees were
1121 then visualised using MEGAX (v.1.0)[153], with bootstrap values less than 80% being
1122 condensed. Ancestor-specific duplication events were inferred using trees generated
1123 using previous phylostratigraphic analysis, where branches containing common earliest
1124 ancestors were deemed as such.

1125

1126 **Visualisation of genes controlling primary and secondary cell wall formation**
1127 Annotations of fern genes were based on the phylogenetic trees in the methods
1128 mentioned above, with genes annotated where reference genes from *A. thaliana* were
1129 found. Coexpression coefficients of each gene within 26 fern species were calculated
1130 using Pearson's correlation coefficient (PCC) and transformed into Highest Reciprocal
1131 Rank (HRR)[154] and genes of interest (GOIs) were isolated. Coexpressed GOIs were
1132 visualised using Cytoscape (v 3.10.1) (https://cytoscape.org/).

1133

1134 **Analysis of CAZymes and HRGPs**
1135 Protein files from coding sequences of 39 plant species were submitted to the dbCAN2
1136 pipeline [155], which annotates CAZymes using three tools (HMMer against the CAZyme
1137 domain database; DIAMOND for BLASTP against the CAZyme database; dbCAN-sub for
1138 HMMER detection of putative CAZy substrates). A majority vote was used and all
1139 annotations were supported by two or more tools, which were used to filter for relevant
1140 CAZy families. CAZy families, as well as respective functionally described enzymes, were
1141 aligned using MAFFT [136] (preferably L-INS-i; if the sequence dataset was too big the
1142 automatic mode was used). Sequence alignments were submitted to FastTree in default
1143 mode [131], and homologs of functionally described enzymes were filtered using iTOL [156].
1144 For DUF families (DUF579, DUF231) and selected other enzymes (CGR2-3, BS1,
1145 DARX1, P4H, QUA2 and QUA3), BLASTP was used with the described family members

1146 against the protein files as a database. E-value of $e^{-7}$ was used, with the rest followed as

1147 described above.

1148 Annotation of hydroxyproline-rich glycoproteins (HRGPs) was performed by using

1149 the workflow described in [73]. The protein sequences were filtered first for the presence of

1150 N-terminal signal peptides and then classified into 24 classes based on the presence of

1151 distinct amino acid motifs and biases (as outlined in [157]).

1152

**Comprehensive Microarray Polymer Profiling (CoMPP)**

1154 The CoMPP analysis was performed according to the method reported by [86]. Each

1155 sample was weighed out in triplicate of 10 mg AIR. The samples were sequentially treated

1156 with 300 µL 50 mM trans-1,2-diaminocyclohexane-N,N,N′,N′-tetraacetic acid (CDTA) pH

1157 7.5, followed by extraction with 300 µL 4 M NaOH containing 0.1% (v/v) $NaBH_4$. Each

1158 extraction step was carried out for 2 h in a TissueLyser II (Qiagen AB, Sollentuna,

1159 Sweden) at 6 s-1 at room temperature. After each extraction, samples were centrifuged

1160 for 10 minutes at 4000 rpm, and the supernatant was collected. The samples were added

1161 to a 384 well plate and four dilution points were prepared for each sample, then two

1162 technical replicates printed on nitrocellulose using an ArrayJet Marathon printer (ArrayJet,

1163 Roslin, UK).

1164 Separate arrays for each probe were first blocked with 5% (w/v) low-fat milk

1165 powder solution in phosphate-buffered saline (MP/PBS), then probed with a set of specific

1166 primary monoclonal antibodies (LM, JIM and MAC 207; Plant Probes, Leeds university,

1167 CCRC; Complex Carbohydrate Research Center, University of Georgia, BS-400;

1168 Biosupplies Australia and INRA-RU donated from Marie Christine Ralet, INRA, France)

1169 (Table S14) for 2 h. After three washes with PBS, the arrays were incubated with 1:5000

1170 solutions of either anti-mouse or anti-rat secondary antibodies (depending on the source

1171 of the primary antibody) conjugated with alkaline phosphatase for another 2 h. Following

1172 three washes with PBS the array had a final wash in Milli-Q water. Arrays were developed

1173 with a 5-bromo-4-chloro-3-indolyl-phosphate (BCIP)/nitro-blue tetrazolium chloride (NBT)

1174 substrate andscanned using a flatbed scanner (CanoScan 9000 Mark II; Canon, Søborg,

1175 Denmark) at 2400 dpi converting the dots to grayscale. The calculated intensity of the

1176 signal was quantified using the microarray analysis software ProScanArray Express

1177 (PerkinElmer, Waltham, Massachusetts, USA). The relative intensity values were
1178 normalised to a scale from 0 to 100 and transformed into a heatmap.

1179

1180 **Enrichment and depletion of biological processes in fern-specific genes and**
1181 **neighbourhood**
1182 Fern-specific genes were defined as genes within orthogroups located in nodes 6-13
1183 (Figure 2a), excluding those orthogroups that contained only genes from one species.
1184 Neighbours of fern-specific genes were defined as genes that are co-expressed with a
1185 fern-specific gene and are not fern-specific genes themselves. The functional annotations
1186 of genes were retrieved (first-level Mapman bins) and subjected to enrichment and
1187 depletion analysis against a background of genes assigned to orthogroups. The analysis
1188 was performed for each fern using a hypergeometric test and adjusted for multiple testing
1189 via Benjamini-Hoechberg correction (q < 0.05)[158]. The overall trend of enrichment or
1190 depletion of biological processes across fern species was derived by subtracting the
1191 number of depleted Mapman bins across all species from the number of enriched bins.

1192

1193 **Data availability**
1194 The raw sequencing data is available at E-MTAB-13848, while the CDS and protein
1195 sequences are found at https://doi.org/10.6084/m9.figshare.26347330. The co-
1196 expression networks are available at https://conekt.sbs.ntu.edu.sg/species/.

1197

1213

**Supplementary Methods**

**Supplementary Methods 1: Inferences of species tree, whole genome duplications, analysis of salicylic acid- and jasmonic acid-mediated signalling in ferns, chemical synthesis of methylated sugars.**

**Supplementary tables:**

**Table S1. 22 species of Tracheophyta from 22 different families.** Samples were collected from various locations in Singapore, with each species having multiple organs harvested.

**Table S2. Sequencing statistics.** The columns contain descriptions of the 415 sample names, including the species, organs, organ types, and data statistics.

**Table S3. Transcriptome assembly statistics for the 22 ferns.** BUSCO value, MapMan annotation percentage, number of transcripts, GC% content, N50, BUSCO scores and percentage of genes annotated by MapMan are shown.

**Table S4 Clade, order, family, species and source of data of the 108 ferns used in this study.**

**Table S5. Fossil calibrations of soft constraint adopted in this study.**

**Table S6. The 95% HPD and posterior mean age estimates (mya) for the origin of each major clade**

**Table S7. Phylostratigraphic assignments of orthogroups to nodes.** The table shows the orthogroups,     clades which are present in the orthogroup and the node where the orthogroup appeared in.

**Table S8. Gene-Organ Specificity.** The table shows the differen species (given by mnemonic), SPM value in a given sample, the number of genes in an organ and the gene ids speciffically expressed in the ogan.

**Table S9. Missing/Present Mapman Bins across 39 species, comprising of Glaucophytes, Chlorophytes, Bryophytes, Lycophytes and Ferns.**

**Table S10. The percentage of annotated clades by MapMan bins.** The species comprising the clades are indicated in column B.

**Table S11. CASA lignin content and thioacidolysis analysis, showing content of H, G and S units.**

**Table S12, BLAST scores of Arabidopsis and Selaginella lignin-related genes against 39 species contained in conekt.sbs.ntu.edu.sg.**

**Table S13. Neutral Sugar Analysis.** The different species and their organs are shown in rows. The columns indicate the abundance of sugars and the standard deviation.

**Table S14. Antibodies, their immunogens and declared specificities and the references where the antibodies were generated/described.**

**Table S15 CoMPP profiles of the 102 cell wall extracts probed with 48 antibodies.** The used solvent and antibodies are shown in columns, the species, organs are shown in rows.


**Supplementary Figures:**

**Figure S1. Pictures of the 22 ferns and their sampled organs.**

**Figure S2. Transcriptome assembly (blue boxes) and subsequent analyses (red boxes).**

**Figure S3.** Genomic properties of ferns in relation to whole genome duplication events. The plots show the correlation between WGD events (x-axis) and species richness (the number of species within a lineage), holoploid genome size (total DNA content), monoploid genome size (DNA content of a single set of chromosomes) and others.

**Figure S4. AlphaFold3-derived structures of the 17 fern-specific proteins.** The colors indicate confidence scores of the structures.

**Figure S5. Number of genes (y-axis) with a given SPM value (x-axis).** The SPM value cutoff is indicated by the red line.

**Figure S6. Gene expression profiles of organ-specific genes.** Each gene's expression has been scaled to range from 0 to 1.

**Figure S7. Expression profiles for species-specific genes.**

**Figure S8. Transcriptome similarity comparison of Archaeplastida.** The heatmap shows the conservation of organ-specific orthogroups across the species. The Jaccard index of across species similarities are indicated by red shades, and for within species similarly with blue shades.

**Figure S9. Gene functions found in Archaeplastida.** Mapman bins (rows) are found in the different species (columns). The colours indicate the fraction of found bins in a given species, where 1 indicates that all genes in a given bin are present, while 0 indicates complete absence.

**Figure S10. Enrichment and depletion of biological processes in neighbours of fern-specific genes.** Cluster map showing significantly enriched and depleted primary Mapman bins (y-axis) in neighbours of fern-specific genes across ferns. The colour map indicates the significance of the biological processes, with yellow representing a p-value of 0.05 and blue representing a p-value of 0.00. P-values above 0.05 are masked.

**Figure S11. Light field, fluorescence, phloroglucinol and Maule staining of the selected ferns.**

**Figure S12. The number (x-axis) of the transcription factors, CYP450s enzyme families and lignin-related enzymes co-expressed with at least two lignin biosynthetic enzymes in the analyzed ferns.**

**Figure S13. GC-MS analysis of 3O-MeRhap.** a) MS profile of 3O-MeRhap. b) Tectaria incisa contains 3O-MeRhap (read arrow) and rhamnose. c) Psilotum nudum contains rhamnose but no 3O-MeRhap.

**Figure S14. The protocol of chemical synthesis of 2-O-methyl- and 3-O-methyl-α,β-D-galactopyranose and 2-O-methyl-D-glucopyranose.**

**Figure S15. GC-MS analysis of cell wall sugars.** a) Profile of Tectaria incisa, b) Tectaria incisa and 2O-MeGlcp standard and c) 2O-MeGlcp standard. d) GC-MS spectra of the unknown peak and the methylated sugar standards.

**Figure S16. Gene copy number analysis of hydroxyproline-rich glycoproteins (HRGPs).** Columns represent species, while rows correspond to a given class of HRGP.

1299 Red and blue numbers indicate that a given species contains significantly more/less

1300 genes than others.

1301 **Figure S17. CoMPP analysis of ferns.** a) The clustermap shows fern samples (rows)

1302 and antibodies (columns). The cells indicate the signal scaled from 0 (dark blue) to 1

1303 (bright yellow). b) Pearson Correlation Coefficient (PCC) distribution of CoMPP profiles

1304 within (blue) and across (brown) species.

1305 **Figure S18. Phylogenetic analysis of CESA genes.** The blue circles represent

1306 bootstrap values (value <50 are not indicated by a circle). The leaf colors represent the

1307 different species and orders.

1308 **Figure S19.** Phylogenetic analysis of lignin biosynthetic genes. Any inferred duplication

1309 events are indicated.

1310

1311 **Supplementary Data 1. Co-expression networks of lignin-related genes, cellulose**

1312 **synthases and cell wall-related transcription factors.** The file should be opened in

1313 cytoscape.

1314

1315 **References**

1316 1. Kenrick, P. & Crane, P. R. *The Origin and Early Evolution of Plants on Land. Nature* vol.

1317 389 (1997).

1318 2. LLOYD, R. M. Mating systems and genetic load in pioneer and non-pioneer Hawaiian

1319 Pteridophyta*. *Botanical Journal of the Linnean Society* **69**, 23–35 (1974).

1320 3. Ellwood, M. D. F. & Foster, W. A. Doubling the estimate of invertebrate biomass in a

1321 rainforest canopy. *Nature* **429**, 549–551 (2004).

1322 4. González de León, S., Briones, O., Aguirre, A., Mehltreter, K. & Pérez-García, B.

1323 Germination of an invasive fern responds better than native ferns to water and light stress

1324 in a Mexican cloud forest. *Biol Invasions* **23**, 3187–3199 (2021).

1325 5. Pryer, K. M. *et al.* Phylogeny and evolution of ferns (monilophytes) with a focus on the

1326 early leptosporangiate divergences. *American Journal of Botany* **91**, 1582–1598 (2004).

1327 6. I, P. A community-derived classification for extant lycophytes and ferns. *Journal of*

1328    *Systematics and Evolution* **54**, 563–603 (2016).

1329    7.    *Fern Ecology*. (Cambridge University Press, Cambridge, 2010).

1330    doi:10.1017/CBO9780511844898.

1331    8.    Page, C. N. Ecological strategies in fern evolution: a neopteridological overview. *Review of*

1332    *Palaeobotany and Palynology* **119**, 1–33 (2002).

1333    9.    Cao, H. *et al.* Phytochemicals from fern species: potential for medicine applications.

1334    *Phytochem Rev* **16**, 379–440 (2017).

1335    10.   Goswami, H. K., Sen, K. & Mukhopadhyay, R. Pteridophytes: evolutionary boon as

1336    medicinal plants. *Plant Genetic Resources* **14**, 328–355 (2016).

1337    11.   Shukla, A. K. *et al.* Expression of an insecticidal fern protein in cotton protects against

1338    whitefly. *Nat Biotechnol* **34**, 1046–1051 (2016).

1339    12.   Shen, H. *et al.* Large-scale phylogenomic analysis resolves a backbone phylogeny in

1340    ferns. *GigaScience* **7**, gix116 (2018).

1341    13.   Qi, X. *et al.* A well-resolved fern nuclear phylogeny reveals the evolution history of

1342    numerous transcription factor families. *Mol Phylogenet Evol* **127**, 961–977 (2018).

1343    14.   Nitta, J. H., Schuettpelz, E., Ramírez-Barahona, S. & Iwasaki, W. An open and

1344    continuously updated fern tree of life. *Front. Plant Sci.* **13**, (2022).

1345    15.   Rensing, S. A. Why we need more non-seed plant models. *New Phytologist* **216**, 355–360

1346    (2017).

1347    16.   Beerling, D. J., Osborne, C. P. & Chaloner, W. G. Evolution of leaf-form in land plants

1348    linked to atmospheric CO2 decline in the Late Palaeozoic era. *Nature* **410**, 352–354

1349    (2001).

1350    17.   Cronk, Q. C. B. Plant evolution and development in a post-genomic context. *Nat Rev*

1351    *Genet* **2**, 607–619 (2001).

1352    18.   Fernández, P. *et al.* A 160 Gbp fork fern genome shatters size record for eukaryotes.

1353    *iScience* **0**, (2024).

1354    19.    KHANDELWAL, S. Chromosome evolution in the genus Ophioglossum L. *Botanical*

1355    *Journal of the Linnean Society* **102**, 205–217 (1990).

1356    20.    Marchant, D. B. *et al.* Dynamic genome evolution in a model fern. *Nat Plants* **8**, 1038–1051

1357    (2022).

1358    21.    Klekowski, E. J. & Baker, H. G. Evolutionary significance of polyploidy in the pteridophyta.

1359    *Science* **153**, 305–307 (1966).

1360    22.    Haufler, C. H. Ever since Klekowski: testing a set of radical hypotheses revives the

1361    genetics of ferns and lycophytes. *Am J Bot* **101**, 2036–2042 (2014).

1362    23.    Clark, J. *et al.* Genome evolution of ferns: evidence for relative stasis of genome size

1363    across the fern phylogeny. *New Phytologist* **210**, 1072–1082 (2016).

1364    24.    Wood, T. E. *et al.* The frequency of polyploid speciation in vascular plants. *Proceedings of*

1365    *the National Academy of Sciences* **106**, 13875–13879 (2009).

1366    25.    Nakazato, T., Barker, M. S., Rieseberg, L. H. & Gastony, G. J. Evolution of the nuclear

1367    genome of ferns and lycophytes. in *Biology and Evolution of Ferns and Lycophytes* (eds.

1368    Haufler, C. H. & Ranker, T. A.) 175–198 (Cambridge University Press, Cambridge, 2008).

1369    doi:10.1017/CBO9780511541827.008.

1370    26.    Barker, M. S. Karyotype and Genome Evolution in Pteridophytes. in *Plant Genome*

1371    *Diversity Volume 2: Physical Structure, Behaviour and Evolution of Plant Genomes* (eds.

1372    Greilhuber, J., Dolezel, J. & Wendel, J. F.) 245–253 (Springer, Vienna, 2013).

1373    doi:10.1007/978-3-7091-1160-4_15.

1374    27.    Baniaga, A. E. & Barker, M. S. Nuclear Genome Size is Positively Correlated with Median

1375    LTR-RT Insertion Time in Fern and Lycophyte Genomes. *amfj* **109**, 248–266 (2019).

1376    28.    Huang, C.-H., Qi, X., Chen, D., Qi, J. & Ma, H. Recurrent genome duplication events likely

1377    contributed to both the ancient and recent rise of ferns. *J Integr Plant Biol* **62**, 433–455

1378    (2020).

1379    29.    Soltis, D. E. & Soltis, P. S. Polyploidy and Breeding Systems in Homosporous

1380        Pteridophyta: A Reevaluation. *The American Naturalist* **130**, 219–232 (1987).

1381   30.  Nakazato, T., Jung, M.-K., Housworth, E. A., Rieseberg, L. H. & Gastony, G. J. Genetic

1382        Map-Based Analysis of Genome Structure in the Homosporous Fern Ceratopteris richardii.

1383        *Genetics* **173**, 1585–1597 (2006).

1384   31.  Marchant, D. B. *et al.* The C-Fern (Ceratopteris richardii) genome: insights into plant

1385        genome evolution with the first partial homosporous fern genome assembly. *Sci Rep* **9**,

1386        18181 (2019).

1387   32.  Huang, X. *et al.* The flying spider-monkey tree fern genome provides insights into fern

1388        evolution and arborescence. *Nat. Plants* **8**, 500–512 (2022).

1389   33.  Li, F.-W. *et al.* Fern genomes elucidate land plant evolution and cyanobacterial symbioses.

1390        *Nature Plants* **4**, 460–472 (2018).

1391   34.  Fang, Y. *et al.* The genome of homosporous maidenhair fern sheds light on the

1392        euphyllophyte evolution and defences. *Nat. Plants* **8**, 1024–1037 (2022).

1393   35.  Zhong, Y. *et al.* Genomic Insights into Genetic Diploidization in the Homosporous Fern

1394        Adiantum nelumboides. *Genome Biology and Evolution* **14**, evac127 (2022).

1395   36.  Grabherr, M. G. *et al.* Trinity: reconstructing a full-length transcriptome without a genome

1396        from RNA-Seq data. *Nat Biotechnol* **29**, 644–652 (2011).

1397   37.  Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-read de novo

1398        assembler. *Gigascience* **1**, 18 (2012).

1399   38.  Leebens-Mack, J. H. *et al.* One thousand plant transcriptomes and the phylogenomics of

1400        green plants. *Nature* **574**, 679–685 (2019).

1401   39.  Rahmatpour, N. *et al.* Analyses of Marsilea vestita genome and transcriptomes do not

1402        support widespread intron retention during spermatogenesis. *New Phytol* **237**, 1490–1494

1403        (2023).

1404   40.  Vanneste, K., Sterck, L., Myburg, A. A., Van de Peer, Y. & Mizrachi, E. Horsetails are

1405        ancient polyploids: Evidence from Equisetum giganteum. *Plant Cell* **27**, 1567–1578 (2015).

1406    41.    Zhang, C. & Mirarab, S. ASTRAL-Pro 2: ultrafast species tree reconstruction from multi-

1407           copy gene family trees. *Bioinformatics* **38**, 4949–4950 (2022).

1408    42.    Emms, D. M. & Kelly, S. STAG: Species Tree Inference from All Genes. Preprint at

1409           https://doi.org/10.1101/267914 (2018).

1410    43.    Le, S. Q. & Gascuel, O. An improved general amino acid replacement matrix. *Mol Biol Evol*

1411           **25**, 1307–1320 (2008).

1412    44.    Wang, F.-G. *et al.* Genome size evolution of the extant lycophytes and ferns. *Plant Divers*

1413           **44**, 141–152 (2022).

1414    45.    Chen, H. *et al.* Revisiting ancient polyploidy in leptosporangiate ferns. *New Phytologist*

1415           **237**, 1405–1417 (2023).

1416    46.    Pelosi, J. A., Kim, E. H., Barbazuk, W. B. & Sessa, E. B. Phylotranscriptomics Illuminates

1417           the Placement of Whole Genome Duplications and Gene Retention in Ferns. *Front Plant*

1418           *Sci* **13**, 882441 (2022).

1419    47.    Julca, I. *et al.* Comparative transcriptomic analysis reveals conserved programmes

1420           underpinning organogenesis and reproduction in land plants. *Nat Plants* **7**, 1143–1159

1421           (2021).

1422    48.    Feng, X. *et al.* Genomes of multicellular algal sisters to land plants illuminate signaling

1423           network evolution. *Nat Genet* **56**, 1018–1031 (2024).

1424    49.    Amorim-Silva, V. *et al.* TTL Proteins Scaffold Brassinosteroid Signaling Components at the

1425           Plasma Membrane to Optimize Signal Transduction in Arabidopsis. *The Plant Cell* **31**,

1426           1807–1828 (2019).

1427    50.    Kazan, K. A new twist in SA signalling. *Nature Plants* **4**, 327–328 (2018).

1428    51.    Li, F.-W. *et al.* Anthoceros genomes illuminate the origin of land plants and the unique

1429           biology of hornworts. *Nat. Plants* **6**, 259–272 (2020).

1430    52.    Monte, I. *et al.* Ligand-receptor co-evolution shaped the jasmonate pathway in land plants.

1431           *Nat Chem Biol* **14**, 480–488 (2018).

53. Piatkowski, B. T. *et al.* Phylogenomics reveals convergent evolution of red-violet coloration in land plants and the origins of the anthocyanin biosynthetic pathway. *Mol Phylogenet Evol* **151**, 106904 (2020).

54. Davies, K. M. *et al.* Evolution and function of red pigmentation in land plants. *Annals of Botany* **130**, 613–636 (2022).

55. Güngör, E. *et al.* Azolla ferns testify: seed plants and ferns share a common ancestor for leucoanthocyanidin reductase enzymes. *New Phytologist* **229**, 1118–1132 (2021).

56. Usadel, B. *et al.* Co-expression tools for plant biology: Opportunities for hypothesis generation and caveats. *Plant, Cell and Environment* **32**, 1633–1651 (2009).

57. Weng, J.-K. & Chapple, C. The origin and evolution of lignin biosynthesis. *New Phytologist* **187**, 273–285 (2010).

58. Rencoret, J. *et al.* New Insights on Structures Forming the Lignin-Like Fractions of Ancestral Plants. *Front. Plant Sci.* **12**, (2021).

59. Renault, H., Werck-Reichhart, D. & Weng, J.-K. Harnessing lignin evolution for biotechnological applications. *Current Opinion in Biotechnology* **56**, 105–111 (2019).

60. Pomar, F., Merino, F. & Barceló, A. R. O-4-Linked coniferyl and sinapyl aldehydes in lignifying cell walls are the main targets of the Wiesner (phloroglucinol-HCl) reaction. *Protoplasma* **220**, 17–28 (2002).

61. Blaschek, L. *et al.* Cellular and Genetic Regulation of Coniferaldehyde Incorporation in Lignin of Herbaceous and Woody Plants by Quantitative Wiesner Staining. *Front. Plant Sci.* **11**, (2020).

62. Yamashita, D., Kimura, S., Wada, M. & Takabe, K. Improved Mäule color reaction provides more detailed information on syringyl lignin distribution in hardwood. *J Wood Sci* **62**, 131–137 (2016).

63. Logan, K. J. & Thomas, B. A. Distribution of Lignin Derivatives in Plants. *New Phytologist* **99**, 571–585 (1985).

1458    64.  Lu, F., Wang, C., Chen, M., Yue, F. & Ralph, J. A facile spectroscopic method for

1459         measuring lignin content in lignocellulosic biomass. *Green Chem.* **23**, 5106–5112 (2021).

1460    65.  Lapierre, C., Pollet, B. & Rolando, C. New insights into the molecular architecture of

1461         hardwood lignins by chemical degradative methods. *Res. Chem. Intermed.* **21**, 397–412

1462         (1995).

1463    66.  Kairouani, A. *et al.* Cell-type-specific control of secondary cell wall formation by Musashi-

1464         type translational regulators in Arabidopsis. *eLife* **12**, RP88207 (2023).

1465    67.  Proost, S. & Mutwil, M. CoNekT: An open-source framework for comparative genomic and

1466         transcriptomic network analyses. *Nucleic Acids Research* **46**, W133–W140 (2018).

1467    68.  Weng, J. K., Li, X., Stout, J. & Chapple, C. Independent origins of syringyl lignin in vascular

1468         plants. *Proceedings of the National Academy of Sciences of the United States of America*

1469         **105**, 7887–7892 (2008).

1470    69.  Weng, J. K. *et al.* Convergent evolution of syringyl lignin biosynthesis via distinct pathways

1471         in the lycophyte Selaginella and flowering plants. *Plant Cell* **22**, 1033–1045 (2010).

1472    70.  Werck-Reichhart, D. & Feyereisen, R. Cytochromes P450: a success story. *Genome*

1473         *Biology* **1**, reviews3003.1 (2000).

1474    71.  Ferrari, C. *et al.* Expression Atlas of Selaginella moellendorffii Provides Insights into the

1475         Evolution of Vasculature, Secondary Metabolism, and Roots. *The Plant Cell*

1476         tpc.00780.2019 (2020) doi:10.1105/tpc.19.00780.

1477    72.  Biswal, A. K. *et al.* Comparison of four glycosyl residue composition methods for

1478         effectiveness in detecting sugars from cell walls of dicot and grass tissues. *Biotechnology*

1479         *for Biofuels* **10**, 182 (2017).

1480    73.  Mueller, K.-K. *et al.* Fern cell walls and the evolution of arabinogalactan proteins in

1481         streptophytes. *The Plant Journal* **114**, 875–894 (2023).

1482    74.  Urbanowicz, B. R. *et al.* 4-O-methylation of glucuronic acid in Arabidopsis glucuronoxylan

1483         is catalyzed by a domain of unknown function family 579 protein. *Proc Natl Acad Sci U S A*

**109**, 14253–14258 (2012).

75. Cantarel, B. L. *et al.* The Carbohydrate-Active EnZymes database (CAZy): an expert resource for Glycogenomics. *Nucleic Acids Res* **37**, D233–D238 (2009).

76. Sampedro, J. *et al.* AtBGAL10 Is the Main Xyloglucan β-Galactosidase in Arabidopsis, and Its Absence Results in Unusual Xyloglucan Subunits and Growth Defects1[W][OA]. *Plant Physiol* **158**, 1146–1157 (2012).

77. Gille, S. *et al.* O-acetylation of Arabidopsis hemicellulose xyloglucan requires AXY4 or AXY4L, proteins with a TBL and DUF231 domain. *Plant Cell* **23**, 4041–4053 (2011).

78. Rocha, J. *et al.* Structure of Arabidopsis thaliana FUT1 Reveals a Variant of the GT-B Class Fold and Provides Insight into Xyloglucan Fucosylation. *Plant Cell* **28**, 2352–2364 (2016).

79. Mikkelsen, M. D. *et al.* Ancient origin of fucosylated xyloglucan in charophycean green algae. *Commun Biol* **4**, 1–12 (2021).

80. Yuan, Y. *et al.* Mutations of Arabidopsis TBL32 and TBL33 affect xylan acetylation and secondary wall deposition. *PLoS ONE* **11**, e0146460 (2016).

81. Jensen, J. K. *et al.* Identification of a xylogalacturonan xylosyltransferase involved in pectin biosynthesis in Arabidopsis. *Plant Cell* **20**, 1289–1302 (2008).

82. Chiniquy, D. *et al.* PMR5, an acetylation protein at the intersection of pectin biosynthesis and defense against fungal pathogens. *Plant J* **100**, 1022–1035 (2019).

83. Wang, D. *et al.* Characterization of CRISPR Mutants Targeting Genes Modulating Pectin Degradation in Ripening Tomato. *Plant Physiology* **179**, 544–557 (2019).

84. Temple, H. *et al.* Two members of the DUF579 family are responsible for arabinogalactan methylation in Arabidopsis. *Plant Direct* **3**, e00117 (2019).

85. Zhong, R., Cui, D. & Ye, Z.-H. Members of the DUF231 Family are O-Acetyltransferases Catalyzing 2-O- and 3-O-Acetylation of Mannan. *Plant and Cell Physiology* **59**, 2339–2349 (2018).

86. Moller, I. *et al.* High-throughput mapping of cell-wall polymers within and between plants using novel microarrays. *Plant J* **50**, 1118–1128 (2007).

87. Fry, S. C., Nesselrode, B. H. W. A., Miller, J. G. & Mewburn, B. R. Mixed-linkage (1-->3,1-->4)-beta-D-glucan is a major hemicellulose of Equisetum (horsetail) cell walls. *New Phytol* **179**, 104–115 (2008).

88. Sørensen, I. *et al.* Mixed-linkage (1-->3),(1-->4)-beta-D-glucan is not unique to the Poales and is an abundant component of Equisetum arvense cell walls. *Plant J* **54**, 510–521 (2008).

89. Silva, G. B. *et al.* Cell wall polysaccharides from fern leaves: evidence for a mannan-rich Type III cell wall in Adiantum raddianum. *Phytochemistry* **72**, 2352–2360 (2011).

90. Fry, S. C. Feruloylated pectins from the primary cell wall: their structure and possible functions. *Planta* **157**, 111–123 (1983).

91. Lampugnani, E. R. *et al.* Cellulose Synthesis – Central Components and Their Evolutionary Relationships. *Trends in Plant Science* **24**, 402–412 (2019).

92. Tsekos, I. The Sites of Cellulose Synthesis in Algae: Diversity and Evolution of Cellulose-Synthesizing Enzyme Complexes. *Journal of Phycology* **35**, 635–655 (1999).

93. Brown Jr, R. M. The Biosynthesis of Cellulose. *Journal of Macromolecular Science, Part A* **33**, 1345–1373 (1996).

94. Seifriz, W. The origin, composition, and structure of cellulose in the living plant. *Protoplasma* **21**, 129–159 (1934).

95. Pear, J. R., Kawagoe, Y., Schreckengost, W. E., Delmer, D. P. & Stalker, D. M. Higher plants contain homologs of the bacterial celA genes encoding the catalytic subunit of cellulose synthase. *Proc Natl Acad Sci U S A* **93**, 12637–12642 (1996).

96. Harholt, J. *et al.* The glycosyltransferase repertoire of the spikemoss selaginella moellendorffii and a comparative study of its cell wall. *PLoS ONE* **7**, (2012).

97. Yin, Y., Huang, J. & Xu, Y. The cellulose synthase superfamily in fully sequenced plants

1536      and algae. *BMC Plant Biology* **9**, 99 (2009).

1537  98. Goubet, F. *et al.* Cell wall glucomannan in Arabidopsis is synthesised by CSLA

1538      glycosyltransferases, and influences the progression of embryogenesis. *Plant J* **60**, 527–

1539      538 (2009).

1540  99. Cocuron, J.-C. *et al.* A gene from the cellulose synthase-like C family encodes a β-1,4

1541      glucan synthase. *Proceedings of the National Academy of Sciences* **104**, 8550–8555

1542      (2007).

1543  100. Wang, W. *et al.* Arabidopsis CSLD1 and CSLD4 are required for cellulose deposition and

1544      normal growth of pollen tubes. *J Exp Bot* **62**, 5161–5177 (2011).

1545  101. Bernal, A. J. *et al.* Functional analysis of the cellulose synthase-like genes CSLD1,

1546      CSLD2, and CSLD4 in tip-growing arabidopsis cells. *Plant Physiology* **148**, 1238–1253

1547      (2008).

1548  102. Burton, R. A. *et al.* Cellulose synthase-like CslF genes mediate the synthesis of cell wall

1549      (1,3;1,4)-beta-D-glucans. *Science* **311**, 1940–1942 (2006).

1550  103. Liu, X. *et al.* Genome-wide bioinformatics analysis of Cellulose Synthase gene family in

1551      common bean (Phaseolus vulgaris L.) and the expression in the pod development. *BMC*

1552      *Genomic Data* **23**, 9 (2022).

1553  104. Bringmann, M. *et al.* POM-POM2/CELLULOSE SYNTHASE INTERACTING1 is essential

1554      for the functional association of cellulose synthase and microtubules in Arabidopsis. *Plant*

1555      *Cell* **24**, 163–177 (2012).

1556  105. Wang, Y. *et al.* LACCASE5 is required for lignification of the brachypodium distachyon

1557      culm. *Plant Physiology* **168**, 192–204 (2015).

1558  106. Höfer, R. *et al.* Dual Function of the Cytochrome P450 CYP76 Family from Arabidopsis

1559      thaliana in the Metabolism of Monoterpenols and Phenylurea Herbicides. *Plant Physiology*

1560      **166**, 1149–1161 (2014).

1561  107. Christenhusz, M. J. M. & Chase, M. W. Trends and concepts in fern classification. *Annals*

1562    *of Botany* **113**, 571–594 (2014).

1563    108. Rothfels, C. J. *et al.* The evolutionary history of ferns inferred from 25 low-copy nuclear

1564        genes. *Am J Bot* **102**, 1089–1107 (2015).

1565    109. Morris, J. L. *et al.* The timescale of early land plant evolution. *Proceedings of the National*

1566        *Academy of Sciences* **115**, E2274–E2283 (2018).

1567    110. Herendeen, P. S., Friis, E. M., Pedersen, K. R. & Crane, P. R. Palaeobotanical redux:

1568        revisiting the age of the angiosperms. *Nature Plants* **3**, 1–8 (2017).

1569    111. Condamine, F. L., Silvestro, D., Koppelhus, E. B. & Antonelli, A. The rise of angiosperms

1570        pushed conifers to decline during global cooling. *Proc Natl Acad Sci U S A* **117**, 28867–

1571        28875 (2020).

1572    112. Bowman, J. L. *et al.* Insights into Land Plant Evolution Garnered from the Marchantia

1573        polymorpha Genome. *Cell* **171**, 287-304.e15 (2017).

1574    113. Weng, J.-K., Banks, J. A. & Chapple, C. Parallels in lignin biosynthesis: A study in

1575        Selaginella moellendorffii reveals convergence across 400 million years of evolution.

1576        *Commun Integr Biol* **1**, 20–22 (2008).

1577    114. Leroux, O. *et al.* Antibody-based screening of cell wall matrix glycans in ferns reveals

1578        taxon, tissue and cell-type specific distribution patterns. *BMC Plant Biol* **15**, 56 (2015).

1579    115. Bartels, D. & Classen, B. Structural investigations on arabinogalactan-proteins from a

1580        lycophyte and different monilophytes (ferns) in the evolutionary context. *Carbohydrate*

1581        *Polymers* **172**, 342–351 (2017).

1582    116. Roberts, A. W. *et al.* Functional Characterization of a Glycosyltransferase from the Moss

1583        Physcomitrella patens Involved in the Biosynthesis of a Novel Cell Wall Arabinoglucan.

1584        *Plant Cell* **30**, 1293–1308 (2018).

1585    117. Taketa, S. *et al.* Functional characterization of barley betaglucanless mutants

1586        demonstrates a unique role for CslF6 in (1,3;1,4)-β-D-glucan biosynthesis. *Journal of*

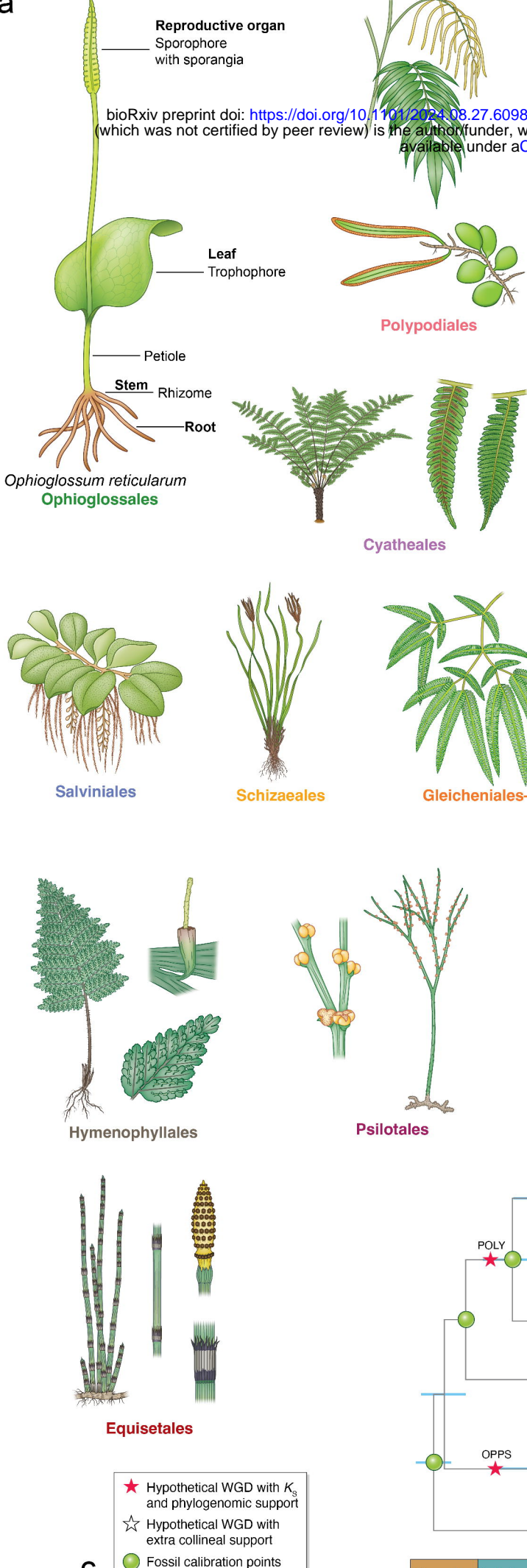1587        *Experimental Botany* **63**, 381–392 (2012).

1588    118. Sperry, J. S. Evolution of Water Transport and Xylem Structure. *International Journal of*

1589         *Plant Sciences* **164**, S115–S127 (2003).

1590    119. Baas, P. & Wheeler, E. A. Parallelism and Reversibility in Xylem Evolution a Review.

1591         (1996) doi:10.1163/22941932-90000633.

1592    120. Ruprecht, C. *et al.* Famnet: A framework to identify multiplied modules driving pathway

1593         expansion in plants. *Plant Physiology* **170**, 1878–1894 (2016).

1594    121. Ruprecht, C. *et al.* Phylogenomic analysis of gene co-expression networks reveals the

1595         evolution of functional modules. *Plant Journal* **90**, 447–465 (2017).

1596    122. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor.

1597         *Bioinformatics* **34**, i884–i890 (2018).

1598    123. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M.

1599         BUSCO: Assessing genome assembly and annotation completeness with single-copy

1600         orthologs. *Bioinformatics* **31**, 3210–3212 (2015).

1601    124. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq

1602         quantification. *Nature Biotechnology* **34**, 525–527 (2016).

1603    125. Sensalari, C., Maere, S. & Lohaus, R. ksrates: positioning whole-genome duplications

1604         relative to speciation events in KS distributions. *Bioinformatics* **38**, 530–532 (2022).

1605    126. Camacho, C. *et al.* BLAST+: Architecture and applications. *BMC Bioinformatics* **10**, 421

1606         (2009).

1607    127. Dongen, S. M. van. Graph clustering by flow simulation.

1608         https://dspace.library.uu.nl/handle/1874/848 (2000).

1609    128. Edgar, R. C. MUSCLE: a multiple sequence alignment method with reduced time and

1610         space complexity. *BMC Bioinformatics* **5**, 113 (2004).

1611    129. Yang, Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Molecular Biology and*

1612         *Evolution* **24**, 1586–1591 (2007).

1613    130. Zwaenepoel, A. & Van de Peer, Y. wgd—simple command line tools for the analysis of

1614      ancient whole-genome duplications. *Bioinformatics* **35**, 2153–2155 (2019).

1615    131. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 – Approximately Maximum-Likelihood

1616      Trees for Large Alignments. *PLoS One* **5**, e9490 (2010).

1617    132. Proost, S. *et al.* i-ADHoRe 3.0—fast and sensitive detection of genomic homology in

1618      extremely large data sets. *Nucleic Acids Research* **40**, e11 (2012).

1619    133. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-

1620      generation sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).

1621    134. Emms, D. M. & Kelly, S. OrthoFinder: solving fundamental biases in whole genome

1622      comparisons dramatically improves orthogroup inference accuracy. *Genome Biology* **16**,

1623      (2015).

1624    135. Steenwyk, J. L. *et al.* OrthoSNAP: A tree splitting and pruning algorithm for retrieving

1625      single-copy orthologs from gene family trees. *PLOS Biology* **20**, e3001827 (2022).

1626    136. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:

1627      improvements in performance and usability. *Mol Biol Evol* **30**, 772–780 (2013).

1628    137. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: A tool for automated

1629      alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973

1630      (2009).

1631    138. Nguyen, L. T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: A fast and

1632      effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular*

1633      *Biology and Evolution* **32**, 268–274 (2015).

1634    139. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermiin, L. S.

1635      ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat Methods* **14**,

1636      587–589 (2017).

1637    140. Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2:

1638      Improving the Ultrafast Bootstrap Approximation. *Molecular Biology and Evolution* **35**, 518–

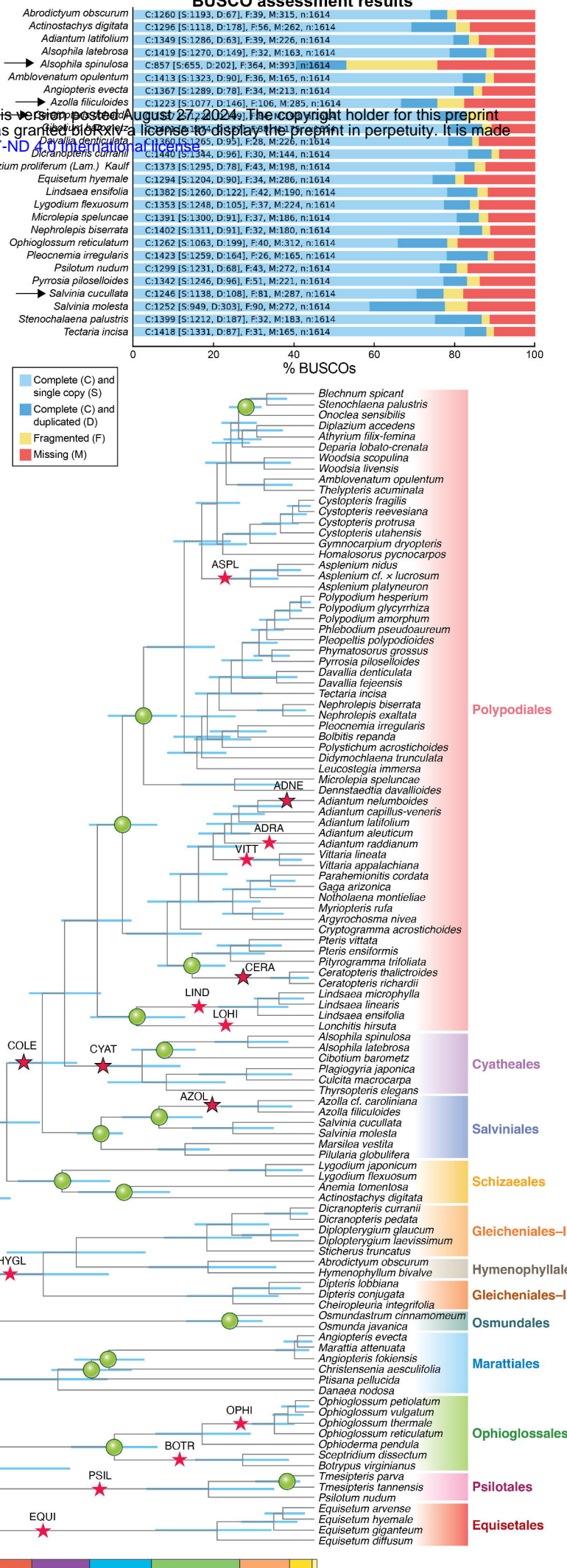1639      522 (2018).

1640    141. Lehtonen, S. *et al.* Environmentally driven extinction and opportunistic origination explain

1641        fern diversification patterns. *Sci Rep* **7**, 4831 (2017).

1642    142. Zwaenepoel, A. & Van de Peer, Y. Inference of Ancient Whole-Genome Duplications and

1643        the Evolution of Gene Duplication and Loss Rates. *Mol Biol Evol* **36**, 1384–1404 (2019).

1644    143. Löytynoja, A. Phylogeny-aware alignment with PRANK. *Methods Mol Biol* **1079**, 155–170

1645        (2014).

1646    144. Ronquist, F. *et al.* Mrbayes 3.2: Efficient bayesian phylogenetic inference and model

1647        choice across a large model space. *Systematic Biology* **61**, 539–542 (2012).

1648    145. Szöllősi, G. J., Rosikiewicz, W., Boussau, B., Tannier, E. & Daubin, V. Efficient Exploration

1649        of the Space of Reconciled Gene Trees. *Systematic Biology* **62**, 901–912 (2013).

1650    146. Buchfink, B., Reuter, K. & Drost, H.-G. Sensitive protein alignments at tree-of-life scale

1651        using DIAMOND. *Nat Methods* **18**, 366–368 (2021).

1652    147. Schwacke, R. *et al.* MapMan4: A Refined Protein Classification and Annotation Framework

1653        Applicable to Multi-Omics Data Analysis. *Molecular Plant* **12**, 879–892 (2019).

1654    148. Smith, A. R. *et al.* A classification for extant ferns. *TAXON* **55**, 705–731 (2006).

1655    149. Abramson, J. *et al.* Accurate structure prediction of biomolecular interactions with

1656        AlphaFold 3. *Nature* **630**, 493–500 (2024).

1657    150. Holm, L. Dali server: structural unification of protein families. *Nucleic Acids Research* **50**,

1658        W210–W215 (2022).

1659    151. van Kempen, M. *et al.* Fast and accurate protein structure search with Foldseek. *Nat

1660        Biotechnol* **42**, 243–246 (2024).

1661    152. Gouy, M., Tannier, E., Comte, N. & Parsons, D. P. Seaview Version 5: A Multiplatform

1662        Software for Multiple Sequence Alignment, Molecular Phylogenetic Analyses, and Tree

1663        Reconciliation. *Methods Mol Biol* **2231**, 241–260 (2021).

1664    153. Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular Evolutionary

1665        Genetics Analysis across Computing Platforms. *Mol Biol Evol* **35**, 1547–1549 (2018).

1666    154. Mutwil, M. *et al.* Assembly of an interactive correlation network for the Arabidopsis genome

1667       using a novel Heuristic Clustering Algorithm. *Plant Physiology* **152**, 29–43 (2010).

1668    155. Zhang, H. *et al.* dbCAN2: a meta server for automated carbohydrate-active enzyme

1669       annotation. *Nucleic Acids Research* **46**, W95–W101 (2018).

1670    156. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree

1671       display and annotation. *Nucleic Acids Research* **49**, W293–W296 (2021).

1672    157. Johnson, K. L. *et al.* Pipeline to Identify Hydroxyproline-Rich Glycoproteins. *Plant*

1673       *Physiology* **174**, 886–903 (2017).

1674    158. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and

1675       Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society: Series B*

1676       *(Methodological)* **57**, 289–300 (1995).

**a**

Reproductive organ
Sporophore
with sporangia

Leaf
Trophophore

Petiole

**Stem**  Rhizome

**Root**

*Ophioglossum reticularum*
**Ophioglossales**

**Cyatheales**

**Polypodiales**

**Salviniales**  **Schizaeales**  **Gleicheniales–II**

**Hymenophyllales**  **Psilotales**

**Equisetales**

★ Hypothetical WGD with *K*s and phylogenomic support

☆ Hypothetical WGD with extra collineal support

● Fossil calibration points

**c**

**b**  BUSCO assessment results

| | |
|---|---|
| *Abrodictyum obscurum* | C:1260 [S:1193, D:67], F:39, M:315, n:1614 |
| *Actinostachys digitata* | C:1296 [S:1118, D:178], F:56, M:262, n:1614 |
| *Adiantum latifolium* | C:1349 [S:1286, D:63], F:39, M:226, n:1614 |
| *Alsophila latebrosa* | C:1419 [S:1270, D:149], F:32, M:163, n:1614 |
| *Alsophila spinulosa* | C:857 [S:655, D:202], F:364, M:393, n:1614 |
| *Amblovenatum opulentum* | C:1413 [S:1323, D:90], F:36, M:165, n:1614 |
| *Angiopteris evecta* | C:1367 [S:1289, D:78], F:34, M:213, n:1614 |
| *Azolla filiculoides* | C:1223 [S:1077, D:146], F:106, M:285, n:1614 |
| | C:1260 [S:1165, D:95], F:28, M:226, n:1614 |
| *Davallia denticulata* | C:1260 [S:1165, D:95], F:28, M:226, n:1614 |
| *Dicranopteris curranii* | C:1440 [S:1344, D:96], F:30, M:144, n:1614 |
| *Diplazium proliferum (Lam.) Kaulf* | C:1373 [S:1295, D:78], F:43, M:198, n:1614 |
| *Equisetum hyemale* | C:1294 [S:1204, D:90], F:34, M:286, n:1614 |
| *Lindsaea ensifolia* | C:1382 [S:1260, D:122], F:42, M:190, n:1614 |
| *Lygodium flexuosum* | C:1353 [S:1248, D:105], F:37, M:224, n:1614 |
| *Microlepia speluncae* | C:1391 [S:1300, D:91], F:37, M:186, n:1614 |
| *Nephrolepis biserrata* | C:1402 [S:1311, D:91], F:32, M:180, n:1614 |
| *Ophioglossum reticulatum* | C:1262 [S:1063, D:199], F:40, M:312, n:1614 |
| *Pleocnemia irregularis* | C:1423 [S:1259, D:164], F:26, M:165, n:1614 |
| *Psilotum nudum* | C:1299 [S:1231, D:68], F:43, M:272, n:1614 |
| *Pyrrosia piloselloides* | C:1367 [S:1246, D:96], F:51, M:221, n:1614 |
| *Salvinia cucullata* | C:1246 [S:1138, D:108], F:81, M:287, n:1614 |
| *Salvinia molesta* | C:1252 [S:949, D:303], F:90, M:272, n:1614 |
| *Stenochalaena palustris* | C:1399 [S:1212, D:187], F:32, M:183, n:1614 |
| *Tectaria incisa* | C:1418 [S:1331, D:87], F:31, M:165, n:1614 |

% BUSCOs

Complete (C) and single copy (S)

Complete (C) and duplicated (D)

Fragmented (F)

Missing (M)

Figure S1. Pictures of the 22 ferns and their sampled organs.

**Figure S2. Transcriptome assembly (blue boxes) and subsequent analyses (red boxes).**

**Figure S3. Genomic properties of ferns in relation to whole genome duplication events.** The plots show the correlation between WGD events (x-axis) and species richness (the number of species within a lineage), holoploid genome size (total DNA content), monoploid genome size (DNA content of a single set of chromosomes) and others.
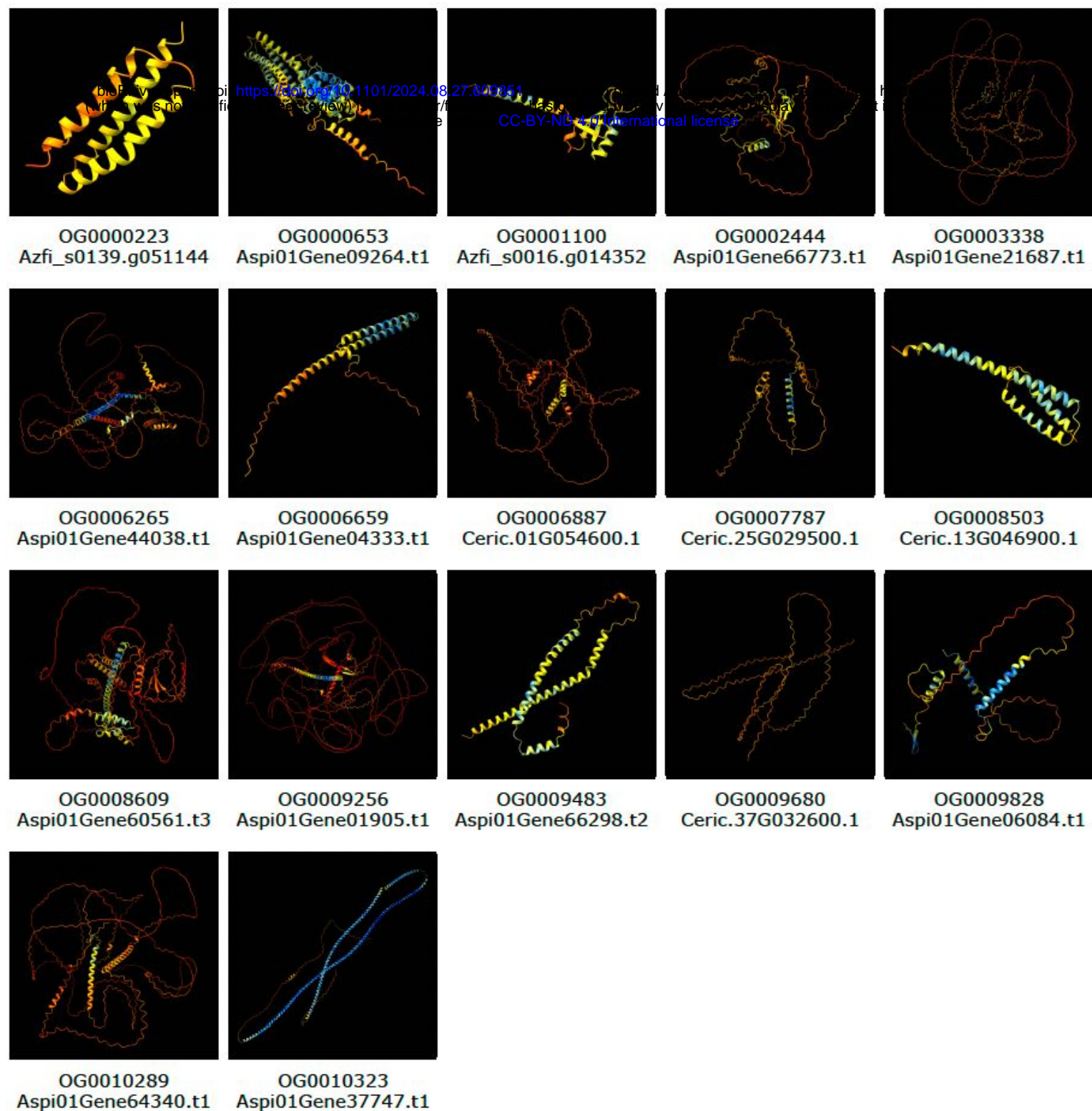
**Figure S4. AlphaFold3-derived structures of the 17 fern-specific proteins.** The colors indicate confidence scores of the structures.

**Figure S5. Number of genes (y-axis) with a given SPM value (x-axis).** The SPM value cutoff is indicated by the red line.

**Figure S6 Gene expression profiles of organ-specific genes.** Each gene's expression has been scaled to range from 0 to 1.

**Figure S7. Expression profiles for species-specific genes.**

**Figure S8. Transcriptome similarity comparison of Archaeplastida.** The heatmap shows the conservation of organ-specific orthogroups across the species. The jaccard index of across species similarities are indicated by red shades, and for within species similarly with blue shades.

**Figure S9. Gene functions found in Archaeplastida.** Mapman bins (rows) found in the different species (columns). The colours indicate the fraction of found bins in a given species, where 1 indicates that all genes in a given bin are present, while 0 indicates complete absence.

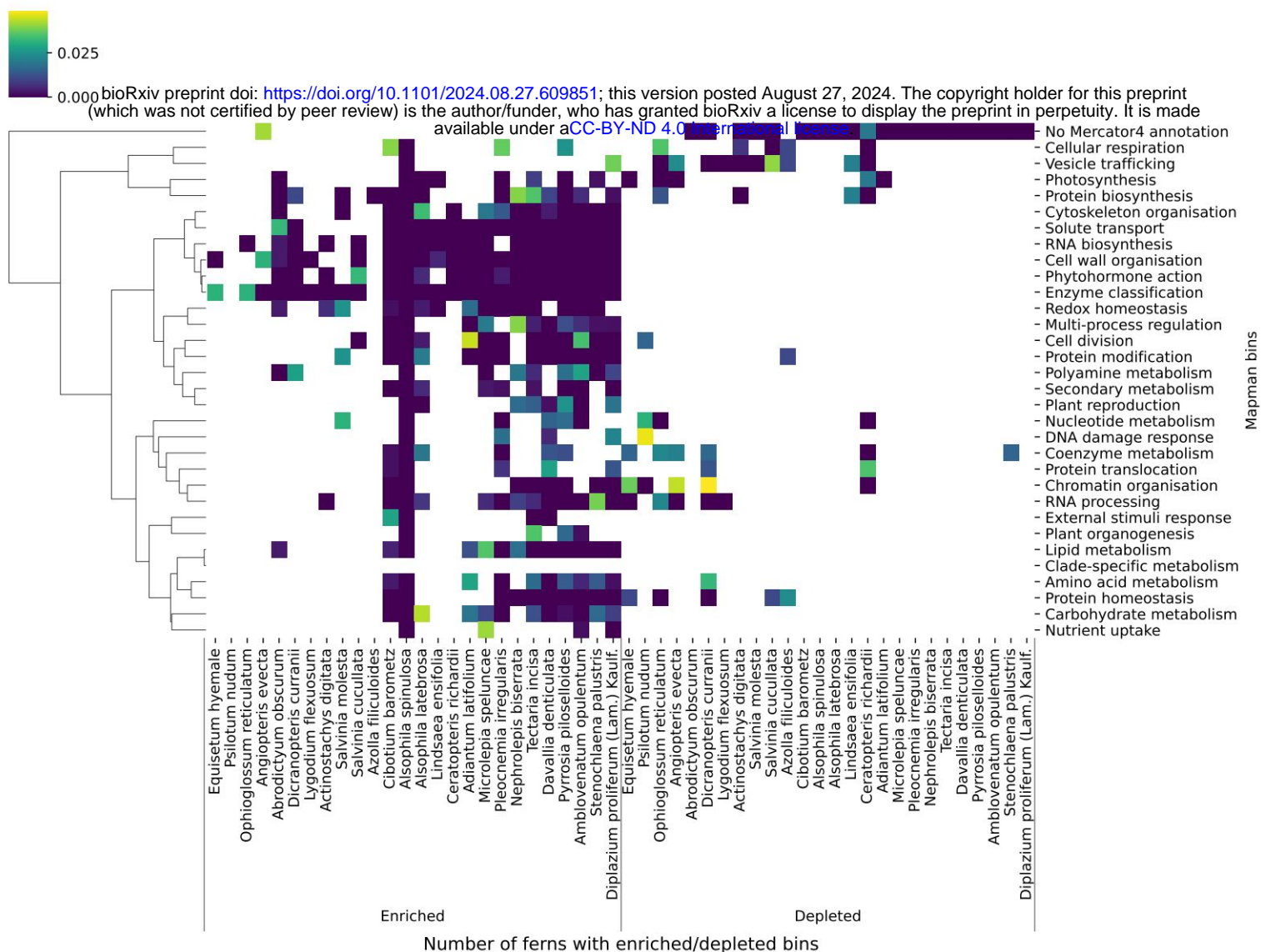**Figure S10. Enrichment and depletion of biological processes in neighbours of fern-specific genes.** Clustermap showing significantly enriched and depleted primary Mapman bins (y-axis) in neighbours of fern-specific genes across ferns. The colour map indicates the significance of the biological processes, with yellow representing a p-value of 0.05 and blue representing a p-value of 0.00. P-values above 0.05 are masked.

Ep: Epidermis, Cs: Cortical Sclerenchyma, Co: Cortex, En: Endodermis,
Pe: Pericycle, Ph: Phloem, MX: Metaxylem

**Figure S11. Light field, fluorescence, phloroglucinol and Maule staining of the selected ferns.**

**Figure S12. The number (x-axis) of the transcription factors, CYP450s enzyme families and lignin-related enzymes co-expressed with at least two lignin biosynthetic enzymes in the analyzed ferns.** Co-expression network of PAL4 (blue circle) from Dicranopteris curranii. CYP96A3 and MYB86 are connected to seven and five enzymes involved in lignin biosynthesis.

a

b



c



**Figure S13. GC-MS analysis of 3*O*-MeRha*p*.** a) MS profile of 3*O*-MeRha*p*. b) T*ectaria incisa* contains 3*O*-MeRha*p* (read arrow) and rhamnose. c) Psilotum nudum contains rhamnose but no 3*O*-MeRha*p*.

**Chart 1**: The protocol of chemical synthesis of 2-*O*-methyl- and 3-*O*-methyl-α,β-D-galactopyranose **5a** and **5b** and 2-*O*-methyl-D-glucopyranose **10**

Methyl β-D-galactopyranoside (**1**)

Methyl 4,6-*O*-benzylididene-β-D-galactopyranoside (**2**)

Methyl α-D-glucopyranoside (**6**)

Methyl 4,6-*O*-benzylidene-2-*O*-(4-methoxybenzyl)-β-D-galactopyranoside (**3a**)

Methyl 4,6-*O*-benzylidene-3-*O*-(4-methoxybenzyl)-β-D-galactopyranoside (**3b**)

Methyl 4,6-*O*-benzylididene-α-D-glucopyranoside (**7**)

Methyl 4,6-*O*-benzylidene-2-*O*-(4-methoxybenzyl)-3-*O*-methyl-β-D-galactopyranoside (**4a**)

Methyl 4,6-*O*-benzylidene-3-*O*-(4-methoxybenzyl)-2-*O*-methyl-β-D-galactopyranoside (**4b**)

Methyl 4,6-*O*-benzylidene-3-*O*-(4-methoxybenzyl)-α-D-glucopyranoside (**8**)

Methyl 4,6-*O*-benzylidene-3-*O*-(4-methoxybenzyl)-2-O-methyl-α-D-glucopyranoside (**9**)

3-*O*-Methyl-D-galactose (**5a**)

2-*O*-Methyl-D-galactose (**5b**)

2-*O*-Methyl-D-glucose (**10**)

**Figure S14. The protocol of chemical synthesis of 2-O-methyl- and 3-O-methyl-α, β-D-galactopyranose 5a and 5b and 2-O-methyl-D-glucopyranose 10.**

**Figure S15. GC-MS analysis of cell wall sugars.** a) Profile of T*ectaria incisa*, b) T*ectaria incisa* and 2*O*-MeGlc*p* standard and c) 2*O*-MeGlc*p* standard. d) GC-MS spectra of the unknown peak and the methylated sugar standards.
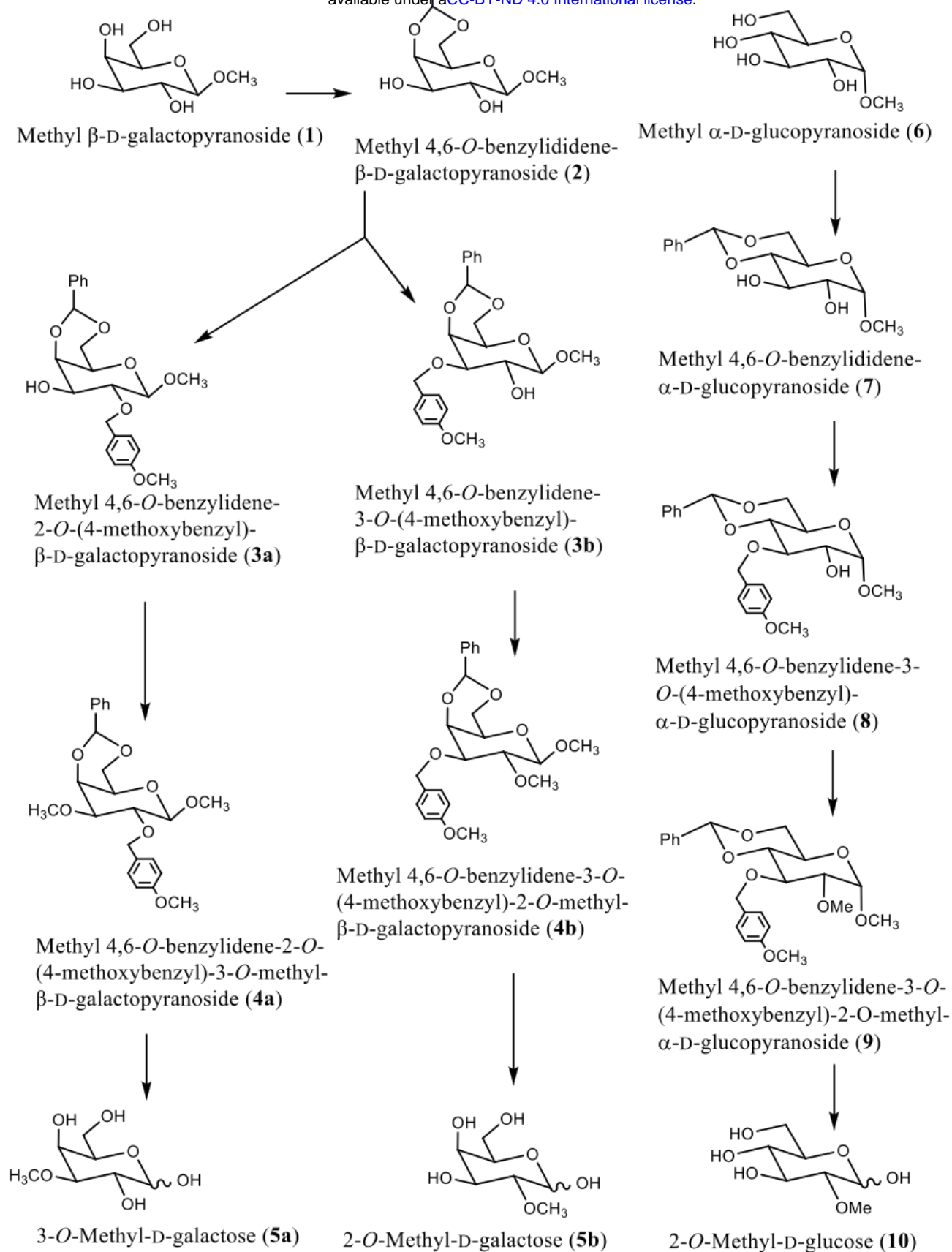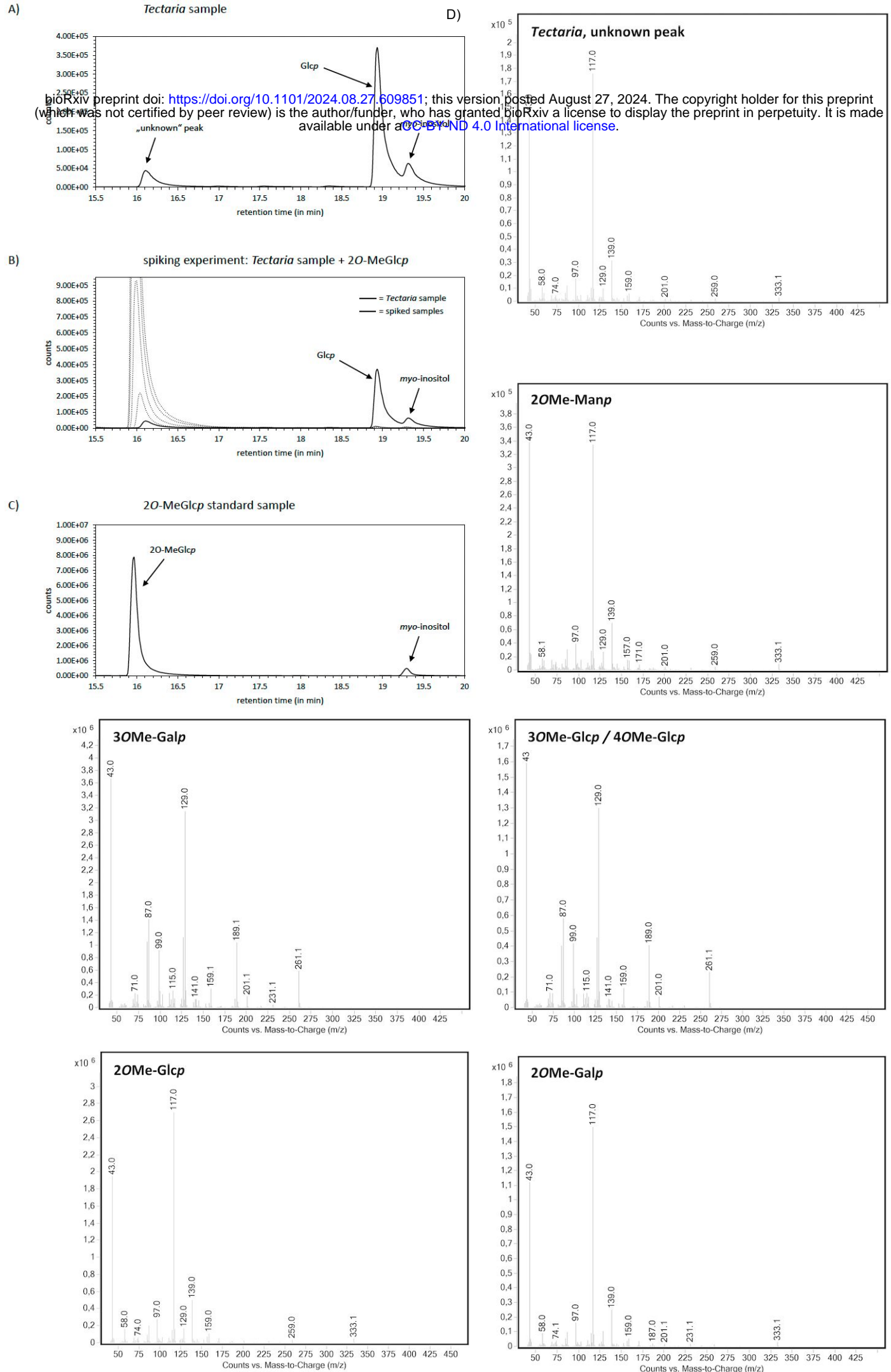
# HRGP Counts by Species

| 1KP Groups | M. polymorpha | P. patens | C. purpureus | S. moellendorffii | D. complanatum | E. hyemale | P. nudum | O. reticulatum | A. evecta | A. obscurum | D. curassavica | L. flexuosum | A. digitatum | S. molesta | S. cucumerina | A. filiculoides | C. baatz | A. spinosa | A. latiossa | L. ensifolia | C. richardii | A. lateritium | M. spicata | P. irregularis | N. bisflora | T. incisum | D. dentata | P. pilosoides | A. opulentum | S. palustris | D. prolificum | G. biloba | P. abies | A. trichopoda | O. sativa | Z. mays | S. lycopersicum | V. vinifera | A. thaliana | class |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GPI-AGPs | 20 | 16 | 6 | 3 | 2 | 1 | 2 | 7 | 2 | 11 | 6 | 1 | 14 | 4 | 3 | 0 | 22 | 3 | 3 | 9 | 9 | 3 | 10 | 13 | 10 | 0 | 1 | 1 | 8 | 6 | 4 | 5 | 9 | 3 | 12 | 20 | 11 | 0 | 17 | class 1 |
| CL-extensins | 1 | 0 | 1 | 6 | 5 | 2 | 0 | 2 | 1 | 4 | 2 | 2 | 2 | 0 | 1 | 0 | 2 | 1 | 3 | 4 | 5 | 0 | 3 | 4 | 1 | 2 | 3 | 1 | 6 | 7 | 3 | 8 | 3 | 4 | 0 | 0 | 5 | 0 | 18 | class 2 |
| PRPs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | class 3 |
| non-GPI-AGPs | 9 | 18 | 28 | 8 | 13 | 9 | 7 | 11 | 5 | 11 | 10 | 9 | 38 | 13 | 4 | 6 | 18 | 15 | 11 | 13 | 29 | 4 | 9 | 8 | 8 | 23 | 15 | 28 | 26 | 16 | 9 | 8 | 3 | 5 | 21 | 25 | 14 | 1 | 7 | class 4 |
| Hybrid AGPs | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 2 | 3 | 3 | 0 | 1 | class 5 |
| Hybrid AGPs | 1 | 1 | 0 | 0 | 2 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 3 | 0 | 0 | 2 | 1 | 1 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 2 | 0 | 1 | class 6 |
| Hybrid AGPs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | class 7 |
| Hybrid AGPs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | class 8 |
| GPI-extensins | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | class 9 |
| Hybrid extensins | 0 | 0 | 4 | 1 | 0 | 1 | 1 | 0 | 1 | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 4 | 6 | 5 | 0 | 1 | 5 | 3 | 2 | 1 | 0 | 10 | 3 | 8 | 3 | 0 | 4 | 0 | 0 | 0 | 10 | 0 | 1 | class 10 |
| Hybrid extensins | 2 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 2 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | class 11 |
| Hybrid extensins | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 3 | 0 | 0 | 0 | 1 | 2 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 10 | 0 | 0 | class 12 |
| Hybrid extensins | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | class 13 |
| Hybrid PRPs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | class 14 |
| Hybrid PRPs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 2 | 5 | 2 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 4 | 0 | 0 | 2 | 0 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | class 15 |
| Hybrid PRPs | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | class 16 |
| Hybrid PRPs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | class 17 |
| Hybrid PRPs | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 4 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | class 18 |
| Shared | 1 | 2 | 2 | 1 | 0 | 0 | 1 | 0 | 1 | 2 | 2 | 0 | 6 | 0 | 0 | 0 | 1 | 1 | 0 | 3 | 0 | 0 | 2 | 0 | 0 | 2 | 2 | 3 | 3 | 1 | 0 | 1 | 3 | 1 | 0 | 1 | 2 | 0 | 1 | class 19 |
| Shared | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 2 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 2 | 7 | 2 | 0 | 0 | 1 | 0 | 6 | class 20 |
| Shared | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | class 21 |
| Shared | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 0 | 0 | 0 | 1 | 0 | 1 | 2 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | class 22 |
| Shared | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | class 23 |
| < 15% motif | 13 | 1 | 6 | 2 | 3 | 10 | 2 | 7 | 9 | 15 | 9 | 3 | 19 | 10 | 2 | 1 | 13 | 0 | 6 | 9 | 12 | 3 | 8 | 13 | 4 | 1 | 5 | 18 | 18 | 11 | 7 | 8 | 3 | 3 | 25 | 16 | 6 | 0 | 7 | class 24 |

**Figure S16. Gene copy number analysis of hydroxyproline-rich glycoproteins (HRGPs).** Columns represent species, while rows correspond to a given class of HRGP. Red and blue numbers indicate that a given species contains significantly more/less genes than others.
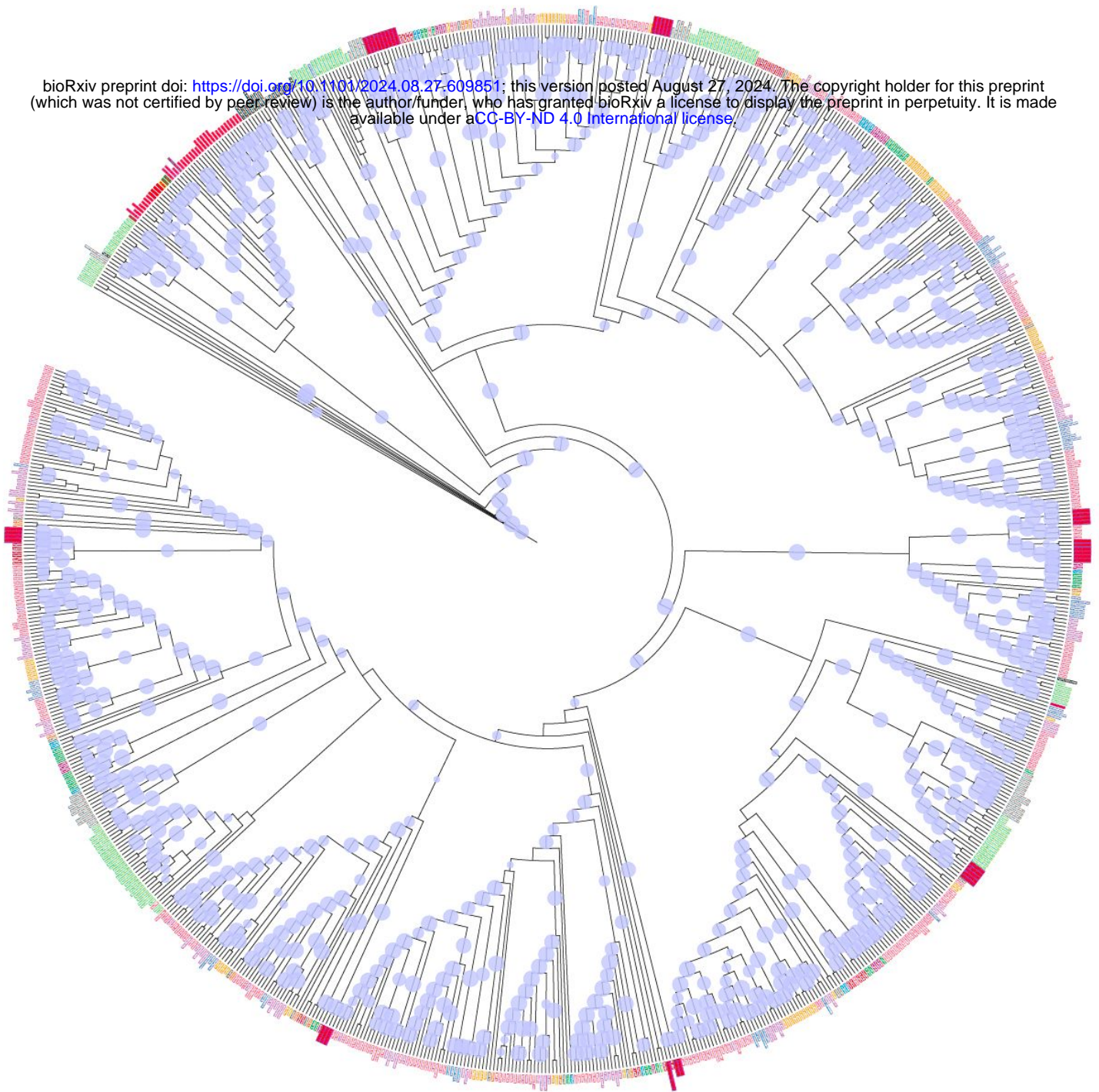
**Figure S17. CoMPP analysis of ferns.** a) The clustermap shows fern samples (rows) and antibodies (columns). The cells indicate the signal scaled from 0 (dark blue) to 1 (bright yellow). b) Pearson Correlation Coefficient (PCC) distribution of CoMPP profiles within (blue) and across (brown) species.

**Figure S18. Phylogenetic analysis of CESA genes.** The blue circles represent bootstrap values (value <50 are not indicated by a circle). The leaf colors represent the different species and orders.
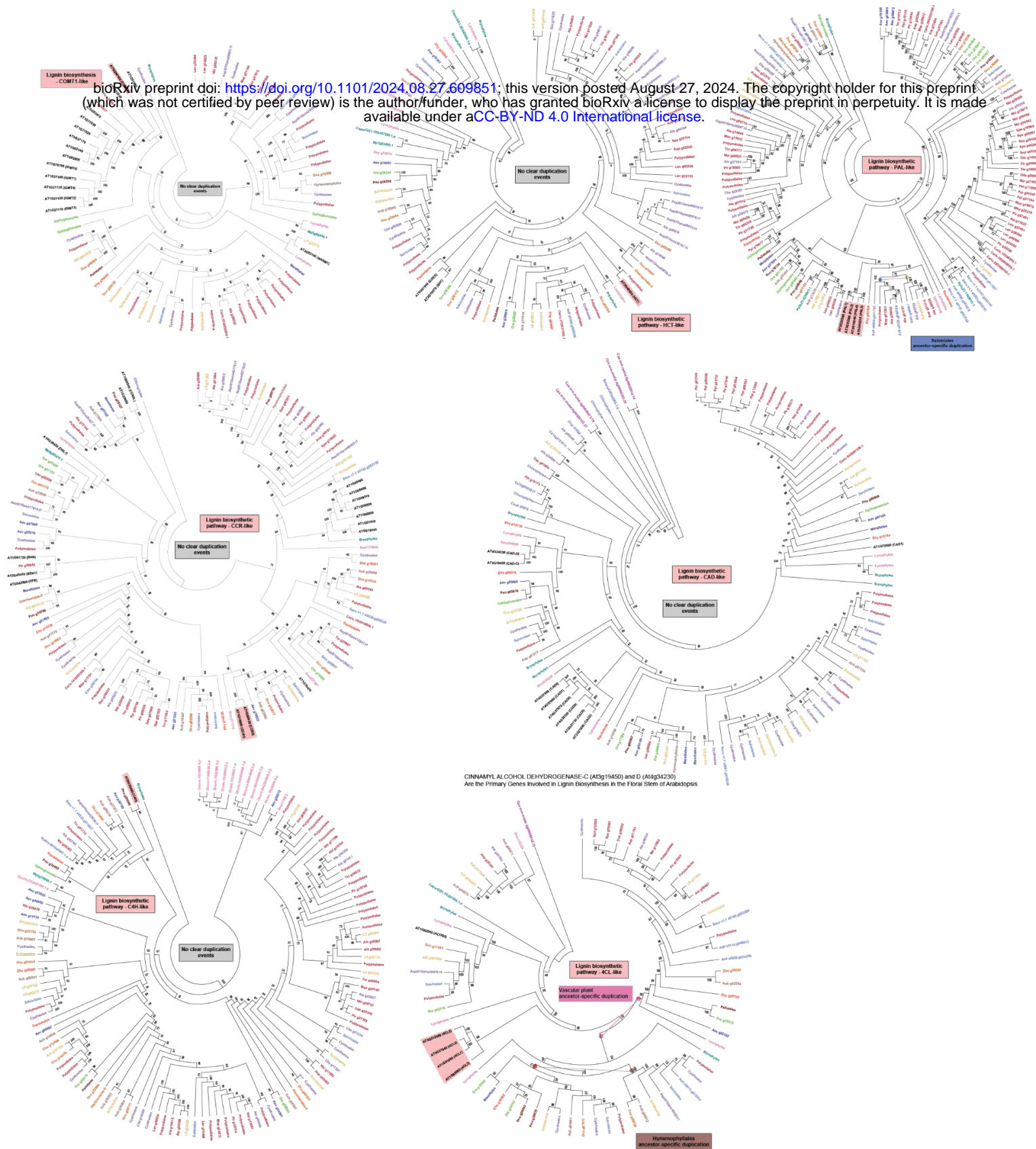
**Figure S19. Phylogenetic analysis of lignin biosynthetic genes.** Any inferred duplication events are indicated.