

1 **Comparative transcriptomic analysis reveals conserved transcriptional programs underpinning** 2 **organogenesis and reproduction in land plants**

3 **Authors:** Irene Julca¹, Camilla Ferrari², María Flores-Tornero³, Sebastian Proost^{2,4,5}, Ann-Cathrin
4 Lindner⁶, Dieter Hackenberg^{7,8}, Lenka Steinbachová⁹, Christos Michaelidis⁹, Sónia Gomes Pereira⁶,
5 Chandra Shekhar Misra^{6,13}, Tomokazu Kawashima^{10,11}, Michael Borg¹⁰, Frédéric Berger¹⁰, Jacob
6 Goldberg¹², Mark Johnson¹², David Honys⁸, David Twell⁷, Stefanie Sprunck³, Thomas Dresselhaus³, Jörg
7 D. Becker^{6,13*}, Marek Mutwil^{1*}

8

9 1) School of Biological Sciences, Nanyang Technological University, 60 Nanyang Drive, Singapore,
10 637551, Singapore

11 2) Max Planck Institute for Molecular Plant Physiology, Am Muehlenberg 1, 14476 Potsdam-Golm,
12 Germany

13 3) Cell Biology and Plant Biochemistry, University of Regensburg, Universitätsstraße 31, 93053
14 Regensburg, Germany

15 4) Laboratory of Molecular Bacteriology, Department of Microbiology and Immunology, Rega Institute,
16 KU Leuven, Herestraat 49, 3000 Leuven, Belgium

17 5) VIB, Center for Microbiology, Kasteelpark Arenberg 31, 3000 Leuven, Belgium

18 6) Instituto Gulbenkian de Ciência, Rua da Quinta Grande 6, 2780-156 Oeiras, Portugal

19 7) Department of Genetics and Genome Biology, University of Leicester, University Road, Leicester,
20 LE1 7RH, UK.

21 8) School of Life Sciences, Gibbet Hill Campus, The University of Warwick, Coventry, CV4 7AL, UK

22 9) Laboratory of Pollen Biology, Institute of Experimental Botany of the Czech Academy of Sciences,
23 Rozvojová 263, 165 02, Prague, Czech Republic

24 10) Gregor Mendel Institute (GMI), Austrian Academy of Sciences, Vienna, BioCenter (VBC), Dr. Bohr-
25 Gasse 3, 1030 Vienna, Austria

26 11) Dept. of Plant and Soil Sciences, University of Kentucky, 321 Plant Science Building, 1405 Veterans
27 Dr., Lexington, KY 40546-0312

28 12) Department of Molecular Biology, Cell Biology, and Biochemistry, Brown University, Providence,
29 RI, 02912, USA

30 13) Instituto de Tecnologia Química e Biológica, Universidade Nova de Lisboa, Av. da República, 2780-
31 157 Oeiras, Portugal

32

33 *Corresponding authors:

34 Marek Mutwil (mutwil@ntu.edu.sg)

35 Jörg D. Becker (jbecker@igc.gulbenkian.pt)

36 **Abstract**

37 The evolution of plant organs, including leaves, stems, roots, and flowers, mediated the explosive
38 radiation of land plants, which shaped the biosphere and allowed the establishment of terrestrial animal
39 life. Furthermore, the fertilization products of angiosperms, seeds serve as the basis for most of our food.
40 The evolution of organs and immobile gametes required the coordinated acquisition of novel gene
41 functions, the co-option of existing genes, and the development of novel regulatory programs. However,
42 our knowledge of these events is limited, as no large-scale analyses of genomic and transcriptomic data
43 have been performed for land plants. To remedy this, we have generated gene expression atlases for
44 various organs and gametes of 10 plant species comprising bryophytes, vascular plants, gymnosperms,
45 and flowering plants. Comparative analysis of the atlases identified hundreds of organ- and gamete-
46 specific gene families and revealed that most of the specific transcriptomes are significantly conserved.
47 Interestingly, the appearance of organ-specific gene families does not coincide with the corresponding
48 organ's appearance, suggesting that co-option of existing genes is the main mechanism for evolving new
49 organs. In contrast to female gametes, male gametes showed a high number and conservation of specific
50 genes, suggesting that male reproduction is highly specialized. The expression atlas capturing pollen
51 development revealed numerous transcription factors and kinases essential for pollen biogenesis and
52 function. To provide easy access to the expression atlases and these comparative analyses, we provide an
53 online database, www.evorepro.plant.tools, that allows the exploration of expression profiles, organ-
54 specific genes, phylogenetic trees, co-expression networks, and others.

55

56 **Introduction**

57 The evolution of land plants has completely changed the appearance of our planet. In contrast to their
58 algal relatives, land plants are characterized by three-dimensional growth and the development of
59 complex and specialized organs. They possess a host of biochemical adaptations, including those

60 necessary for tolerating desiccation and UV stress encountered on land, allowing them to colonize most
61 terrestrial surfaces. The earliest land plants which arose ~470 million years ago ¹, were speculatively
62 similar to extant bryophytes, possessing tiny fertile axes or an axis terminated by a sporangium²⁻⁴. The
63 innovation of shoots and leaves mediated the 10-fold expansion of vascular plants ^{5,6} and an 8–20-fold
64 atmospheric CO₂ drawdown ⁷, which significantly shaped the Earth's geosphere and biosphere ⁸. To
65 enable soil attachment and nutrient uptake, the first land plants only had rhizoids, filamentous structures
66 homologous to root hairs ⁹. Roots evolved to provide increased anchorage (and thus increased height) and
67 enable survival in more arid environments. Parallel with innovations of vegetative cell types, land plants
68 evolved new reproductive structures such as spores, pollen, embryo sacs, and seeds together with the
69 gradual reduction of the haploid phase. In contrast to algae, mosses, and ferns that require moist habitats,
70 the male and female gametophytes of gymnosperms and angiosperms are strongly reduced, consisting of
71 only a few cells, including the gametes ^{10,11}. Moreover, sperm cells have lost their mobility and use pollen
72 grains as a protective vehicle for long-distance transport and a pollen tube for their delivery deep into
73 maternal reproductive tissues ^{12,13}. The precise interaction of plant male and female gametes, leading to
74 cell fusion, karyogamy, and development of both the embryo and endosperm after double fertilization has
75 just begun to be deciphered at the molecular level ^{14,15}. These anatomical innovations are mediated by
76 coordinated changes in gene expression and the appearance of novel genes and/or repurposing of existing
77 genetic material. Genes that are specifically expressed in these organs often play a major role in their
78 establishment and function ^{16,17}, but the identity and conservation of these specifically-expressed genes
79 have not been extensively studied.

80 Nowadays, flowering plants comprise 90% of all land plants and serve as the basis for the terrestrial food
81 chain, either directly or indirectly. The use of model plants like *Arabidopsis thaliana* and maize and
82 technical advances allowing live-cell imaging of double fertilization have been instrumental for several
83 major discoveries ^{18,19}. When assessing current knowledge of male and female gamete development in
84 plants, it is evident that the male germline has been studied to a greater extent ^{11,20}. This is mainly due to

85 its accessibility and the development of methods to separate the sperm cells from the surrounding
86 vegetative cell of pollen, e.g. by FACS²¹. Analysis of male germline differentiation, for example, has led
87 to the identification of Arabidopsis *DUO POLLEN 1 (DUO1)* and the network of genes it controls, which
88 include the fertilization factors, *HAP2/GCSI* and *GEX2*²². However, as novel genes are still being
89 discovered that control the development of male and female gametes^{10,11} or their functions^{23,24}, it is clear
90 that our knowledge of the molecular basis of gamete formation and function is far from complete.

91 Current approaches to study evolution and gene function mainly use genomic data to reveal which gene
92 families are gained, expanded, contracted, or lost. While invaluable, genomic approaches alone might not
93 reveal the function of genes that show no sequence similarity to known genes²⁵. To remedy this, we
94 combined comparative genomic approaches with newly established, comprehensive gene expression
95 atlases of two bryophytes (*Marchantia polymorpha*, *Physcomitrium patens*), a lycophyte (*Selaginella*
96 *moellendorffii*), gymnosperms (*Ginkgo biloba*, *Picea abies*), a basal angiosperm (*Amborella trichopoda*),
97 eudicots (*Arabidopsis thaliana*, *Solanum lycopersicum*) and monocots (*Oryza sativa*, *Zea mays*). We then
98 compared these organ-, tissue- and cell-specific genes to identify novel and missing components involved
99 in organogenesis and gamete development.

100 We show that transcriptomes of most organs are conserved across land plants and report the identity of
101 hundreds of organ-specific gene families. We demonstrate that the age of gene families is positively
102 correlated with organ-specific expression and the appearance of organ-specific gene families does not
103 coincide with the appearance of the corresponding organ. We observed a high number of male-specific
104 gene families and strong conservation of male-specific transcriptomes, while female-specific
105 transcriptomes showed fewer specific gene families with less conservation. Our detailed analysis of gene
106 expression data capturing the development of pollen revealed numerous transcription factors and kinases
107 potentially important for pollen biogenesis and function. Finally, we present a user-friendly, online
108 database www.evorepro.plant.tools, which allows the browsing and comparative analysis of the genomic

109 and transcriptomic data derived from sporophytic and gametophytic samples across 13 members of the
110 plant kingdom.

111

112 **Results**

113 **Constructing gene expression atlases and identifying organ-specific genes**

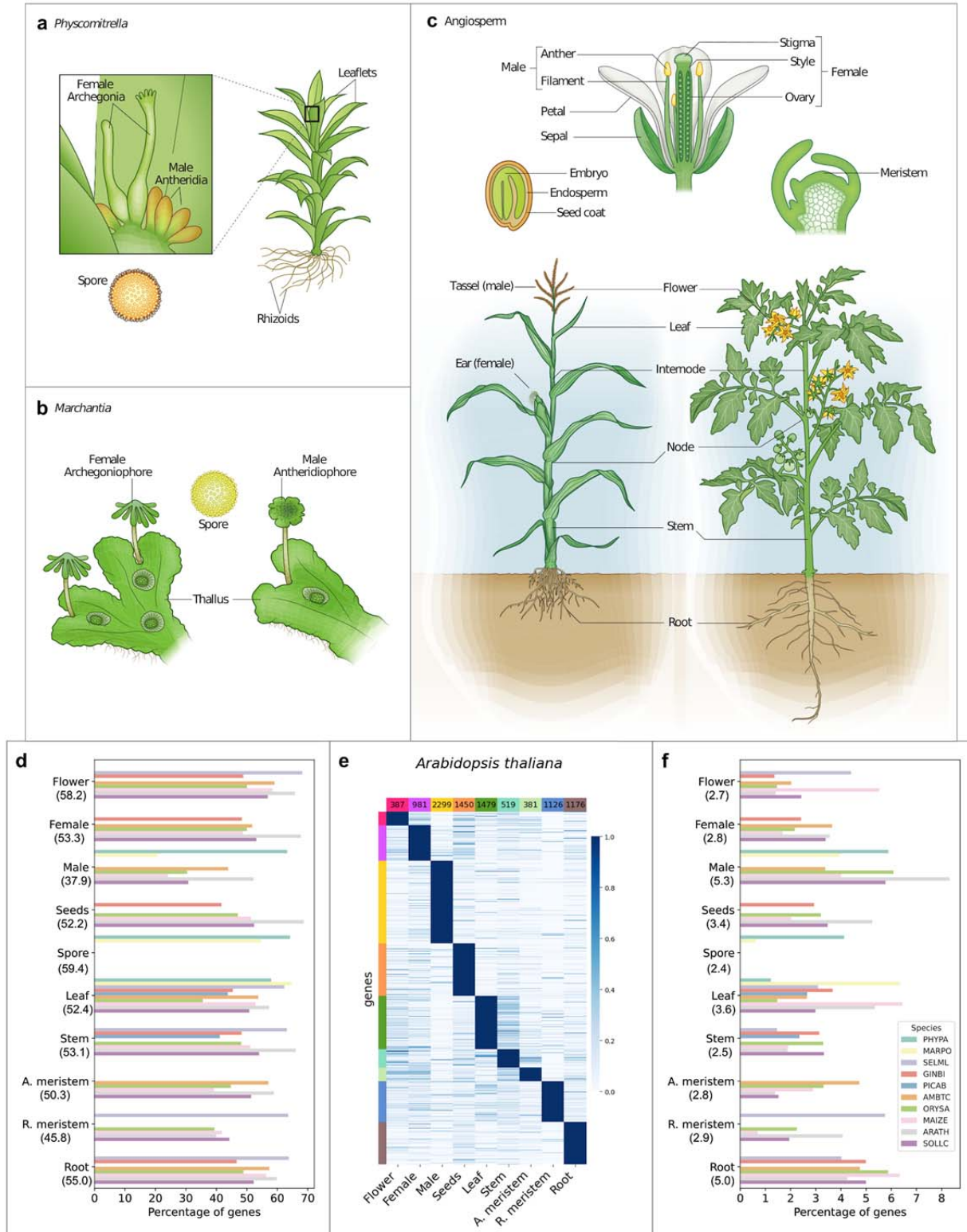
114 We constructed gene expression atlases for ten phylogenetically important species (Table 1). These
115 include the bryophytes *Physcomitrium patens* (*Physcomitrella*) (Fig. 1a) and *Marchantia polymorpha*
116 (Fig. 1b), the lycophyte *Selaginella moellendorffii*, the gymnosperms *Ginkgo biloba* and *Picea abies*, the
117 basal angiosperm *Amborella trichopoda*, the monocots *Oryza sativa* and *Zea mays*, and the eudicots
118 *Arabidopsis thaliana* and *Solanum lycopersicum* (Fig. 1c). The atlases were constructed by combining
119 publicly available RNA sequencing (RNA-seq) data with 134 fastq files generated by the EVOREPRO
120 consortium (see Supplementary Table 1). For each species, we generated an expression matrix that
121 contains transcript-level abundances captured by transcript per million (TPM) values²⁶. The expression
122 matrices capture gene expression values from the main anatomical sample types, which we grouped into
123 ten classes: flower, female, male, seeds, spore, leaf, stem, apical meristem, root meristem, and root (Fig.
124 1a-c). Furthermore, the expression data was used to construct co-expression networks and to create an
125 online EVOREPRO database allowing further analysis of the data (www.evorepro.plant.tools).

126 To identify genes expressed in the different samples, we included only those with an average TPM >2
127 (see methods). For all ten species, approximately 71% of their genes were expressed in at least one
128 structure (Supplementary Table 2). Interestingly, the male sample has a lower percentage (38%) followed
129 by root meristems (46%), while the other samples have between 50-60% expressed genes (Fig. 1d).

130 Organ- and cell-specific genes can often play a major role in the establishment and function of the organ
131 and cell type^{16,17}. To identify such genes, we calculated the specificity measure (SPM) of each gene,

132 which ranges from 0 (not expressed in a sample) to 1 (expressed only in the sample). A threshold
133 capturing top 5% of the SPM values was used to identify the sample-specific genes for all species
134 (Supplementary Fig. 1, Supplementary Table 3). To examine the sample-specific genes' expression
135 profiles, we plotted the scaled TPM values of these genes for *A. thaliana*. Visual inspection shows that
136 the TPM values of the sample-specific genes are in all cases highest in the samples that the genes are
137 specific to (Fig. 1d, Supplementary Fig. 2). For the ten species, an average of 21% of the genes were
138 identified as sample-specific (Supplementary Table 2). The lowest percentage was found in *P. abies*
139 (5%), followed by *M. polymorpha* (11%) and *P. patens* (11%), while the highest percentage was found in
140 *A. thaliana*, where 35% of the transcripts showed sample-specific expression (Supplementary Table 2).
141 These low and high percentages observed can be partially explained by the number of organs and cell
142 types that we analyzed (Supplementary Table 1).

143 Interestingly, we observed that the male (5.3%) and root (5.0%) samples typically contained the highest
144 percentage of specific genes (Fig. 1f, Supplementary Table 2). In *A. thaliana*, the higher percentage of
145 male-specific genes was in agreement with previous studies that showed a high specialization of the male
146 transcriptome^{27,28}. Conversely, stem, spore, apical meristem, root meristem, flower, and female show
147 values lower than 3% (Fig. 1f, Supplementary Table 2). Previous studies also showed the low number of
148 genes mainly expressed in the female gametophyte^{29,30}.



149

150 **Fig. 1: Expression atlases for seven land plant species.** Depiction of the different organs, tissues, and cells
 151 collected for (a) *P. patens* (b) *Marchantia polymorpha*, and (c) angiosperms. d, The percentage of genes (x-axis)

152 found to be expressed (defined as TPM>2) in organs (y-axis) of the different species (indicated by colored bars as in
 153 (f)). The numbers beneath the organs (y-axis) indicate the average percentage of genes for all species. **e**, Expression
 154 profiles of organ-specific genes from *Arabidopsis thaliana*. Genes are in rows, organs in columns and the genes are
 155 sorted according to the expression profiles (e.g., flower, female). The numbers at the top of each column indicate the
 156 total number of genes per organ. Gene expression is scaled to range from 0-1. Bars on the left of each heatmap show
 157 the sample-specific genes and correspond to the samples on the bottom: pink - Flower, purple - Female, yellow -
 158 Male, orange - Seeds/Spore, dark-green - Leaf, medium-green - Stem, light-green - Apical meristem, blue - Root
 159 meristem, brown - Root. **f**, The percentage of genes with specific expression in the ten species.

160

161 **Table 1. Organs, tissues, and cell types used in the expression atlases analyzed.** The different species
 162 are shown in columns, while the rows organize the organs, tissues and cell types into rows.

163

Organ/tissue/cell type	Marchantia	Physcomitrella	Selaginella	Ginkgo	Spruce	Amborella	Arabidopsis	Tomato	Rice	Maize
Flower	N/A	N/A	Strobili	Strobili, microstrobilus	N/A	Tepals, buds, opened flowers	Buds, stamen filaments, carpels, petals, stigmatic tissue, sepals	Buds, opened flowers	Buds, panicles	Tassels, ear, silk
Apical meristem	-	-	-	-	-	Apical meristem	Apical meristem	Apical meristem	Apical meristem	Apical meristem
Male	Sperm	Sperm	-	-	-	Pollen (mature, tube), sperm, microspores, generative cell	Pollen (mature, tube, bicellular, tricellular), microspore, sperm	Pollen (mature, tube), microspore, sperm cell, generative cell	Pollen (bi-, tricellular), microspore, sperm	Pollen (mature, tube), microspore, sperm
Female	-	-	-	Ovules	-	Ovary, egg apparatus cell	Egg cell, ovule	Ovaries, ovary walls, ovules	Ovary, ovule, egg cell	Ovary, ovule, nucellus, egg cell, embryo sac
Root	N/A	N/A	Roots, rhizophores	Roots	-	Roots	Root (apex, tip, differentiation zone, stele, elongation zone)	Root (differentiation zone, elongation zone)	Root (differentiation zone, elongation zone)	Root (tip, secondary, stele, elongation zone, maturation zone)
Root meristem	N/A	N/A	Meristematic zone	-	-	-	Meristematic and QC zone	Meristematic zone	Meristematic zone	Root (meristematic zone)
Leaf	Thallus	Leaflets	Microphyll	Leafs	Needles	Leaves	Leaves	Leaves	Leaves	Leaves
Stem	N/A	N/A	Top stem, bottom stem	Stem, xylem, cambium	Xylem, phloem, cambium	-	Stems	Stems	Stems	Stems
Seed	N/A	N/A	N/A	Kernels	-	-	Seeds (young, germinating),	Seeds (5-30 DPA)	Seeds	Seeds (mature, germinating),

							endosperm			endosperm, pericarp and aleurone
Spore	Sporeling	Spore capsule, germinating spores	-	N/A	N/A	N/A	N/A	N/A	N/A	N/A

164

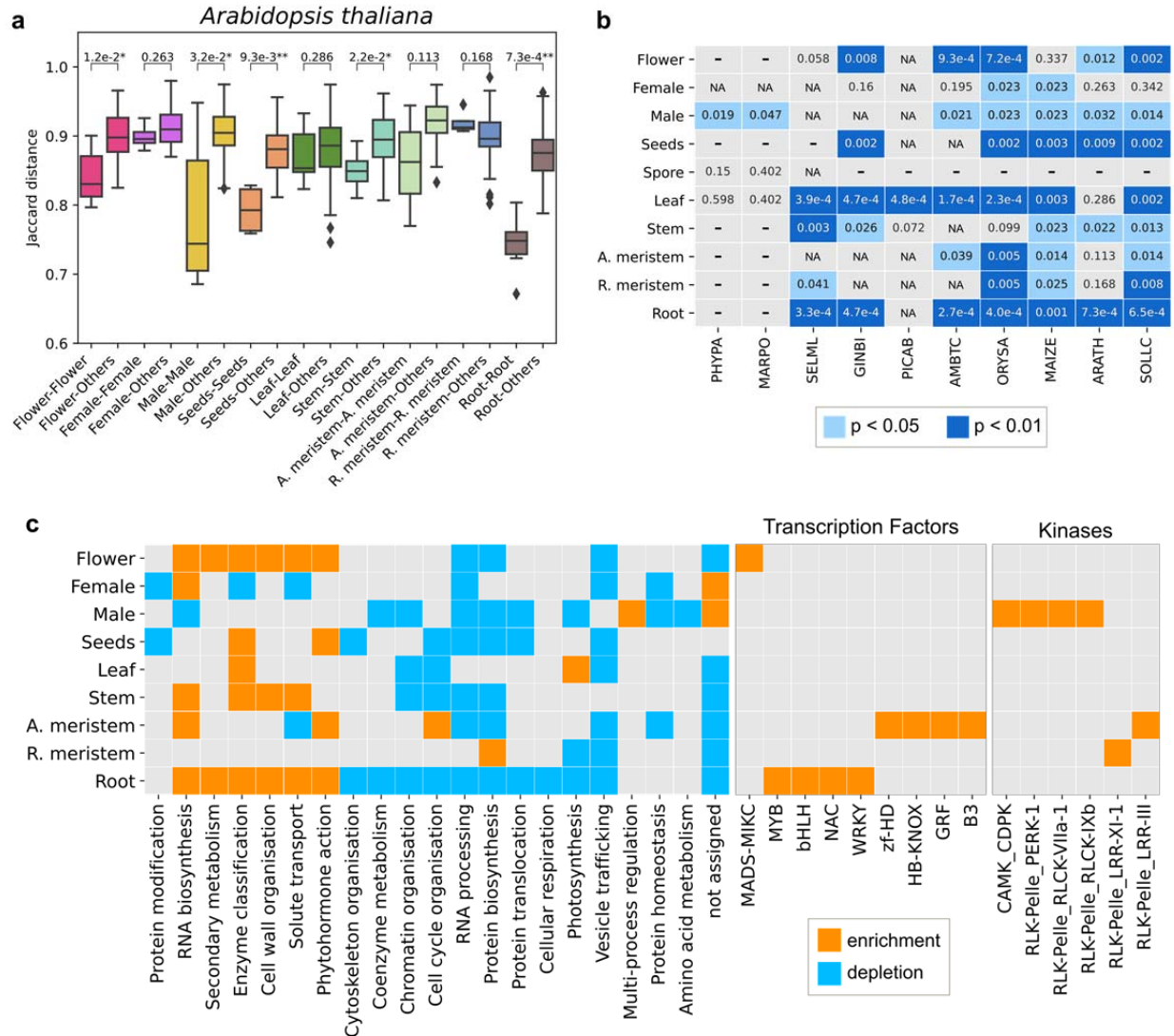
165

166 To summarize, these results show that organ-specific genes represent a significant part of the
167 transcriptome, with male and root samples possessing the most specialized transcriptomes.

168

169 **Are the transcriptomes of organs conserved across species?**

170 Our above analysis suggests that sample-specific gene expression is widespread, and we set out to
171 investigate whether these patterns are conserved across species. To this end, we investigated which
172 samples specifically expressed similar sets of gene families (represented by orthogroups) by employing a
173 Jaccard distance that ranges from 0 (two samples express an identical set of sample-specific gene
174 families) to 1 (none of the sample-specific gene families are the same in the two samples). We expected
175 that if, e.g., the root-specific transcriptome is conserved across angiosperms, then Jaccard distance of root
176 vs. root transcriptomes (e.g., *Arabidopsis* root vs. rice root) should be lower than when comparing root vs.
177 non-root transcriptomes (e.g., *Arabidopsis* root vs. rice leaf).



178

179 **Fig. 2: Comparison of sample-specific transcriptomes.** **a**, Bar plot showing the Jaccard distances (y-axis) when
 180 comparing the same samples (x-axis, e.g., male-male) and one sample versus the others (e.g., male-others) for
 181 *Arabidopsis thaliana*. Lower values indicate a higher similarity of the transcriptomes. **b**, Significantly similar
 182 transcriptomes are indicated by blue cells (light blue p<0.05 and dark blue p<0.01). Species are indicated by the
 183 mnemonic: PHYPA - *P. patens*, MARPO - *Marchantia polymorpha*, SELML - *Selaginella moellendorffii*, GINBI -
 184 *Ginkgo biloba*, PICAB - *Picea abies*, AMBTC - *Amborella trichopoda*, ORYSA - *Oryza sativa*, MAIZE - *Zea mays*,
 185 ARATH - *Arabidopsis thaliana*, SOLLC - *Solanum lycopersicum*. **c**, Heatmap showing the significant (p-value <
 186 0.05) functional enrichment (orange cell) or depletion (blue cell) in the ten sample classes (y-axis) in at least 50%
 187 species. The heatmap indicates Mapman bins (photosynthesis-not assigned), transcription factors, and kinases.

188

189 The analysis revealed that *Arabidopsis* flower-, male-, seeds-, stem- and root-specific transcriptomes were
190 significantly more similar to the corresponding sample in the other species (p-value < 0.05, Fig. 2a).
191 When performing the analysis for all ten species, we observed that root, male, and seeds expressed
192 specifically similar gene families in all species with the samples (7 species for root, 7 for male, and 5 for
193 seeds) and for other organs, some species show significance, flowers (5 out of 7 species with flower
194 samples), female (2 out of 6), leaf (7 out of 10), stem (5 out of 7), apical meristem (4 out of 5), root
195 meristem (4 out of 5) (Fig. 2b, Supplementary Fig. 3). Conversely, spore (0 out of 2) samples did not
196 show similar transcriptomes across *Marchantia* and *Physcomitrella* (Fig. 2b, Supplementary Fig. 3). We
197 also performed clustering analysis between all pairs of sample-specific genes in the ten species and
198 observed root-, seed-, flower, leaf-, meristem- and male-specific clusters (Supplementary Fig. 4).
199 Interestingly, the male samples in *Physcomitrella* and *Marchantia* formed a distinctive cluster
200 (Supplementary Fig. 4), suggesting that flagellated sperm of bryophytes employ a unique male
201 transcriptional program compared with non-motile sperm of angiosperms.

202 To reveal which biological processes are preferentially expressed in the different samples across the ten
203 species, we performed a functional enrichment analysis of Mapman bins, transcription factors, and
204 kinases (Fig. 2c, Supplementary Fig. 5). The analysis revealed that many functions were depleted in male
205 and root samples in at least 50% of the species, indicating that most male and roots' cellular processes
206 were significantly repressed (p-value < 0.05, Fig. 2c, Supplementary Fig. 5). As expected, genes
207 associated with photosynthesis were enriched in leaves but depleted in roots, root meristems, and male
208 samples. Genes expressed in roots were enriched in solute transport functions, enzyme classification
209 (enzymes not associated with other processes), RNA biosynthesis, secondary metabolism, phytohormone
210 action, and cell wall organization (Fig. 2c). Interestingly, female and male reproductive cells were
211 enriched for 'not assigned' bin, indicating that these organs are enriched for genes with unknown
212 functions.

213 Since the sample-specific genes (Supplementary Table 3) are likely important for the formation and
214 function of the organ, we investigated sample-specific transcription factors (Supplementary Table 4) and
215 receptor kinases (Supplementary Table 5). An enrichment analysis of transcription factors (69 families)
216 and kinases (142 families) showed that apical meristem and root samples were highly enriched in
217 transcription factors, while male and apical meristem were enriched for kinases (Fig. 2c). In apical
218 meristems, some of the enriched transcription factor families (C2C2-YABBY, GRF) were associated with
219 the regulation, development, and differentiation of meristem^{31,32}. In roots, the enriched transcription
220 factors (MYB, bHLH, WRKY, NAC) are related to biotic and abiotic stress response and root
221 development³³⁻³⁷. These sample-specific genes are thus prime candidates for further functional analysis
222 (Supplementary Table 5).

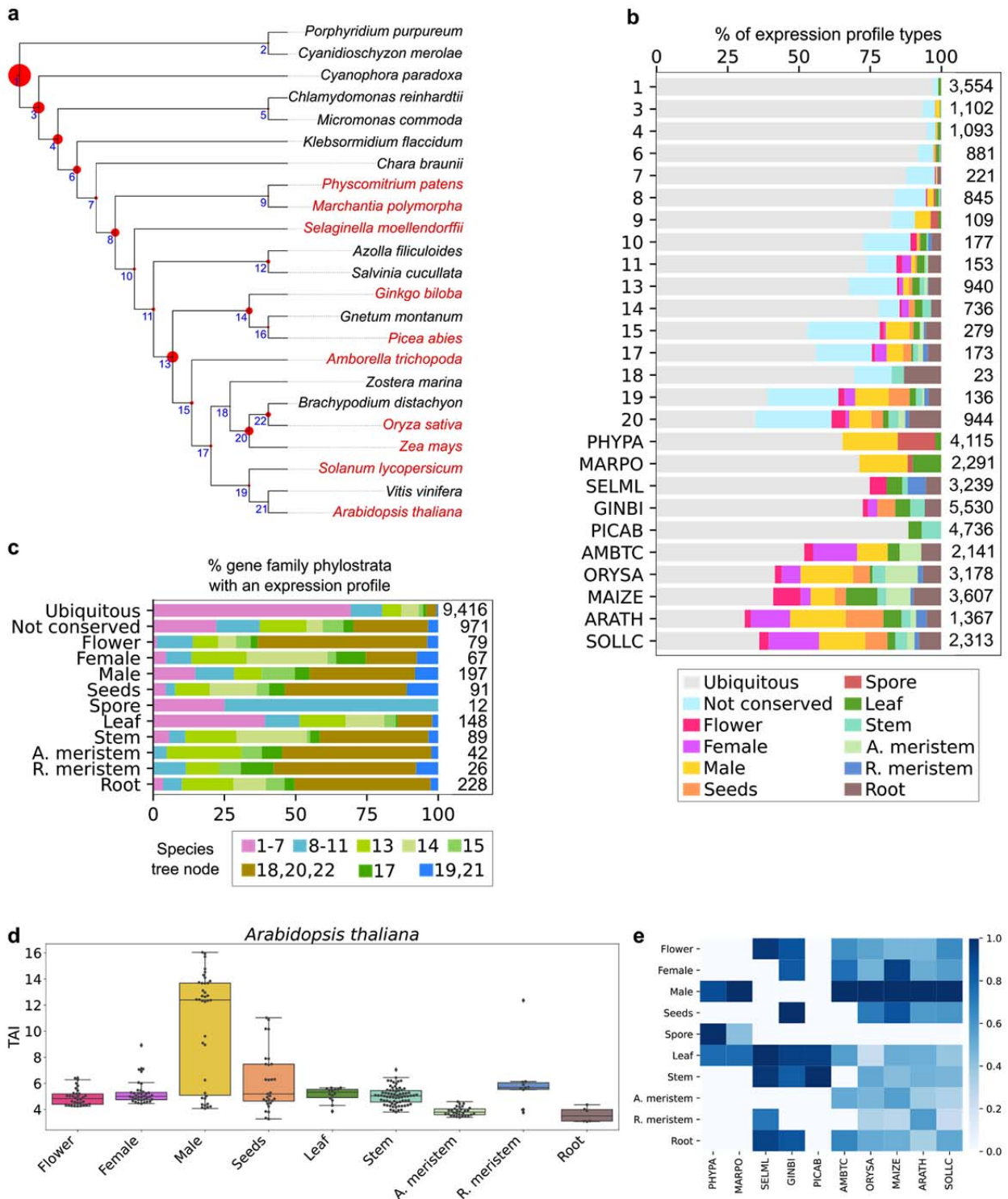
223

224 **Phylostratigraphic analysis of sample-specific gene families**

225 Organs, such as seeds and flowers, appeared at a specific time in plant evolution. To investigate whether
226 there is a link between gene families' appearance and their expression patterns, we used the proteomes of
227 23 phylogenetically important species and a derived species tree based on One Thousand Plant
228 Transcriptomes Initiative (2019). Each orthogroup was placed to one node (phylostrata) of the species
229 tree, where node 1 indicated the oldest phylostratum, and node 23 indicated the youngest, species-specific
230 phylostratum (Supplementary Table 6). A total of 131,623 orthogroups were identified in the 23
231 Archaeplastida, of which 113,315 (86%) were species-specific, and the remaining 18,308 (14%) were
232 assigned to internal nodes. Of these internal node orthogroups, most were ancestral (24% - node 1, 10% -
233 node 3), belonged to streptophytes (7%, node 6), land plants (7%, node 8), seed plants (10%, node 13),
234 monocots (0.3%, node 18), or eudicots (1%, node 19) (Fig. 3a). Analysis of phylostrata in each species
235 revealed a similar distribution of the orthogroups, with most of them belonging to node 1 (~34%) or being
236 species-specific (~31%, Supplementary Fig. 6).

237 To investigate whether the different phylostrata show different expression trends, we surveyed
238 orthogroups that contain at least two species with RNA-seq data, which resulted in 37,887 (29% of the
239 total number of orthogroups) meeting this criterion. Then, each orthogroup was assigned to different
240 expression profiles: ubiquitous (not specific in any organ), not conserved (e.g., root-specific in one
241 species, flower-specific in others), or organ-specific (for details see material and methods, Supplementary
242 Table 6). The majority of the orthogroups in internal nodes (not species-specific) of the phylogenetic tree
243 were assigned as ubiquitous (9,416), which corresponded to orthogroups that showed broad and not
244 organ-specific expression (Fig. 3b). Interestingly, we observed a clear pattern of gene families becoming
245 increasingly organ-specific as phylostratigraphic age decreased (<5% specific genes in node 1, vs. ~25%
246 in node 13), indicating that younger gene families are recruited to specific organs (Fig. 3b).

247



248

249 **Fig. 3: Genomic analysis of sample-specificity of gene families.** **a**, Species tree of the 23 species for which we
 250 have inferred orthogroups. Species in red are the ones with transcriptomic data available. Blue numbers in the nodes
 251 indicate the node number (e.g., 1: node 1). The tree's red circles show the percentage of orthogroups found at each

252 node (largest: node 1 - 24% of all orthogroups, smallest: node 21 - 0.1%). **b**, Percentage of expression profile types
253 of orthogroups per node. The expression profile types are: ubiquitous (light gray, orthogroup is not organ-specific),
254 not conserved (light blue, organ-specificity not conserved in different species), or sample-specific (e.g., brown: root-
255 specific). **c**, Percentage of phylostrata (nodes) within the different expression profile types. **d**, Transcriptome age
256 index (TAI) of the different sample-specific genes in *Arabidopsis thaliana*. The boxplots show the TAI values (y-
257 axis) in the different organs (x-axis), where a high TAI value indicates that the sample expresses a high number of
258 younger genes. **e**, Summary of the average TAI value in the ten species. The organs are shown in rows, while the
259 species are shown in columns. The TAI values were scaled to 1 for each species by dividing values in a column with
260 the highest column value.

261

262 Next, we identified sample-specific gene families and investigated when they appeared during plant
263 evolution. The number of gene families in internal nodes per sample varied from 12 (spore) to 228 (root),
264 and we observed trends of samples across the internal nodes. In general, many organ-specific orthogroups
265 were present in nodes corresponding to monocots (Node 18, 20, 22). Expectedly, the 9,416 ubiquitous
266 orthogroups were mostly of ancient (node 1-7) origin, suggesting that these old gene families tend to
267 show a broader expression. The nonconserved groups had both old and more recent gene families. From
268 the organ-specific families, leaves and spores were the groups containing more ancient families, while
269 meristems had younger families. Flower, root, seeds, stem had few older families. Interestingly, when we
270 compared male and female groups, we observed that the male-specific orthogroups had older gene
271 families than the female-specific orthogroups (Fig. 3c).

272 Several studies revealed that new genes in animals tend to be preferentially expressed in male
273 reproductive tissues, such as testis³⁸⁻⁴⁰. Similar observations have been made in *Arabidopsis*, rice, and
274 soybean⁴¹, where new genes were predominantly expressed in male reproductive cells⁴², suggesting that
275 these cells may act as an “innovation incubator” for the birth of *de novo* genes. Our gene expression data
276 also revealed that male samples possess the youngest transcriptome in *Arabidopsis* (Fig. 3d, yellow bar),
277 and in the male samples of *M. polymorpha*, *A. trichopoda*, *Z. mays*, *O. sativa*, *S. lycopersicum*, but not in
278 *P. patens* (Fig. 3e, dark-blue cells for male, Supplementary Fig. 7). With the unclear exception in
279 *Physcomitrella*, we conclude that the observation that male samples express young genes is robust in the

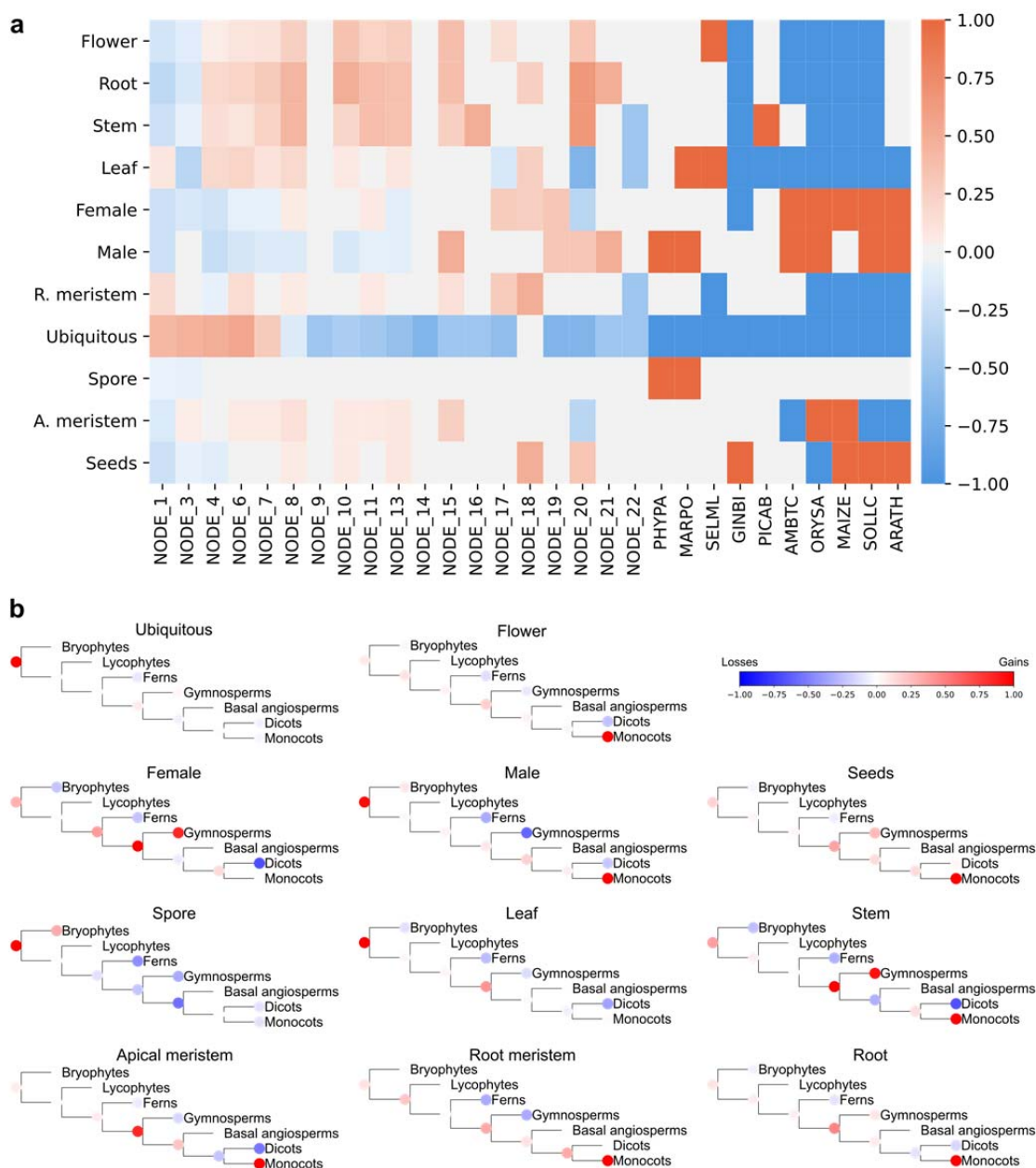
280 plant kingdom. However, pollen also expresses a substantial portion of old genes (species nodes 1-7 in
281 Figure 3c), probably representing an old transcription program present in gametes in Archaeplastida.

282

283 **Phylostratigraphic and gene expression analysis reveals that co-option drives the evolution of**
284 **organs**

285 The evolution of land plants involved many major innovations mediated by gains and losses of gene
286 families and co-option of existing gene functions. Most of the changes are related to land adaptations
287 comprising requirements for structural support, uptake of water, prevention of desiccation and gas
288 exchange⁴³. To better understand this complex process, we first analyzed the enrichment/depletion of
289 organ-specific and ubiquitous genes in each node of the species tree (Supplementary Table 7). In line with
290 previous results (Figure 3b), ubiquitous genes were enriched for genes that appeared before the
291 divergence of land plants and depleted for genes that appeared when plants colonized land (node 8, Fig.
292 4a). In line with the basal function (photosynthesis) of leaves, leaf-specific genes were enriched in
293 ancestral nodes and the species-specific nodes of *M. polymorpha* (thallus samples) and *S. moellendorffii*
294 (microphyll), and depleted in species-specific nodes of the seed plants (Fig. 4a).

295 Leaf-specific gene families were acquired mainly in two ancestral nodes, before the divergence of land
296 plants and before the divergence of seed plants (Fig. 4b). Most of the gene families were gained in node 1
297 (34 families, Supplementary Table 8). Leaves have multiple origins in land plants^{44,45}, however, the
298 programs for oxygenic photosynthesis originated in ancient organisms⁴⁶. In agreement, before the
299 divergence of land plants, we observed enrichment for functions related to photosynthesis (<N8), and
300 after the divergence of land plants, we detected enrichment for additional functions such as external
301 stimuli response, cytoskeleton organization, phytohormone action, and protein modification
302 (Supplementary Table 9).



303

304 **Fig. 4: Evolutionary analysis of organs.** **a**, Enrichment and depletion of organ-specific genes per node in the
 305 species tree (nodes are the same as in Fig. 3a). The colors correspond with the number of species showing
 306 enrichment in each case (dark red: all species show enrichment, dark blue: all species show depletion). **b**,
 307 Cladograms of the main lineages showing gain (in red) and loss (blue) of gene families with ubiquitous and sample-
 308 specific expression profiles.

309 Interestingly, stem-, root-, and flower-specific genes shared a similar pattern and appeared to be enriched
310 in nodes 4-8, 10-13, 15, and 20, and depleted in the species-specific nodes of vascular plants, except for
311 *P. abies* for stems and *S. moellendorffii* for flowers. Although the origins(s) of roots, stems, and flowers
312 are associated with vascular plants⁴⁷⁻⁴⁹, we observed gene family expansions before the divergence of
313 land plants (Fig. 4b) and in nodes as old as node 3 (2 orthogroups) for stems, node 1 (1 orthogroup) for
314 roots, and node 3 (1 orthogroup) for flowers (Supplementary Table 8). Previous studies suggested that the
315 evolution of novel morphologies was mainly driven by the reassembly and reuse of pre-existing genetic
316 mechanisms^{45,50}. It was indicated that primitive root programs may have been present before the
317 divergence of lycophytes and euphyllophytes⁵¹. Also, before the divergence of charophytes from land
318 plants, an ancestral origin was proposed for the SVP subfamily, which plays a crucial role in the control
319 of flower development^{52,53}. A recent study has shown that a moss (*Polytrichum commune*) possesses a
320 vascular system functionally comparable to that of vascular plants⁵⁴. These results support the idea that
321 primitive stem-, root-, and flower-specific gene families existed prior to vascular plants' divergence. After
322 the divergence of land plants, we can observe that there is incremental gene family gain in monocots for
323 all three organs (roots, stems, flowers, Fig. 4b, indicated by red nodes), and also to a lesser extent in the
324 ancestral node of seed plants. Specifically, for stem, we observed more gains in gymnosperms and more
325 losses in eudicots. Functional enrichment analysis supports only enrichment in nodes corresponding to
326 land plants (>N8) and not in older nodes (Supplementary Table 9).

327 Apical and root meristem-specific genes appeared enriched in ancestral nodes and depleted in species-
328 specific nodes, with the exception of apical meristem in monocots that are enriched (Fig. 4a). The
329 analysis of gain/loss of gene families showed that many apical meristem-specific orthogroups were
330 gained in seed plants and monocots and lost in eudicots. For root meristem-specific gene families we
331 observed that many orthogroups were gained in monocots (Fig. 4b). Functional enrichment analysis for
332 apical meristem-specific gene families shows enrichment of unknown functions in nodes N19 and N20,

333 and for root meristem-specific gene families shows enrichment for phytohormone action in N8 and
334 protein modification in N15 (Supplementary Table 9).

335 Seed-specific genes were enriched only in nodes of land plants. The nodes that showed enrichment were
336 N10 (vascular plants), N18 and N20 (monocots), and species-specific nodes with the exception of *O.*
337 *sativa*, which showed depletion of this set of genes. Some seed-specific families were gained before the
338 divergence of land plants, but interestingly the higher number of gains was observed in N20 (monocots -
339 39 gene families), followed by N14 (gymnosperms - 15), N13 (seed plants - 11), and N19 (eudicots - 10)
340 (see Fig. 4b, Supplementary Table 8). Enrichment of functions related to solute transport was observed
341 only in eudicots (N19, Supplementary Table 9).

342 Spore-specific genes were enriched only in the species-specific nodes of bryophytes (Fig. 4a). However,
343 gene family gains were observed in ancestral nodes (N4, N6, N8, N9, see Supplementary Table 8) and
344 lipid metabolism enrichment only in the node ancestral to bryophytes (N9, Supplementary Table 9).

345 Male-specific genes were enriched in angiosperms (N15), monocots (N20), eudicots (N19, N21), and
346 species-specific nodes, while female-specific genes were enriched only in monocots (N18, N22), eudicots
347 (N19), and species-specific nodes (Fig. 4a). Additional male-specific families were gained in older nodes
348 than female-specific families (intensity of the red color in the ancestral node of land plants, Fig. 4b). For
349 male gene families, we observed six waves of gains (>15 gene families) in nodes N3, N8 (land plants),
350 N13 (seed plants), N15 (angiosperms), N19 (eudicots), N20 (monocots). From these nodes, parallel to
351 gains, we also observed many losses (>=10 gene families) in three nodes N13 (seed plants), N15
352 (angiosperms), and N19 (eudicots) (Supplementary Table 8). For female-specific families, we observed
353 three main waves of gains (>10 gene families) in nodes N13 (seed plants), N14 (gymnosperms), N20
354 (monocots), and different waves of losses (Supplementary Table 8). Male gene families showed
355 enrichment for protein modification, enzyme classification, RNA biosynthesis, cell cycle organization,
356 phytohormone action, and female gene families showed enrichment only for RNA biosynthesis

357 (Supplementary Table 9). Considering gains and losses of gene families, male-specific families were
358 gained mainly in the node ancestral to land plants, and in monocots, and for female-specific families in
359 seed plants and gymnosperms (Fig. 4b).

360 In summary, the genetic programs for organ-specific genes are present in older nodes, before the
361 divergence of land plants. Monocots seem to be the group with more gene family gains, which is in
362 agreement with previous studies⁵⁵.

363

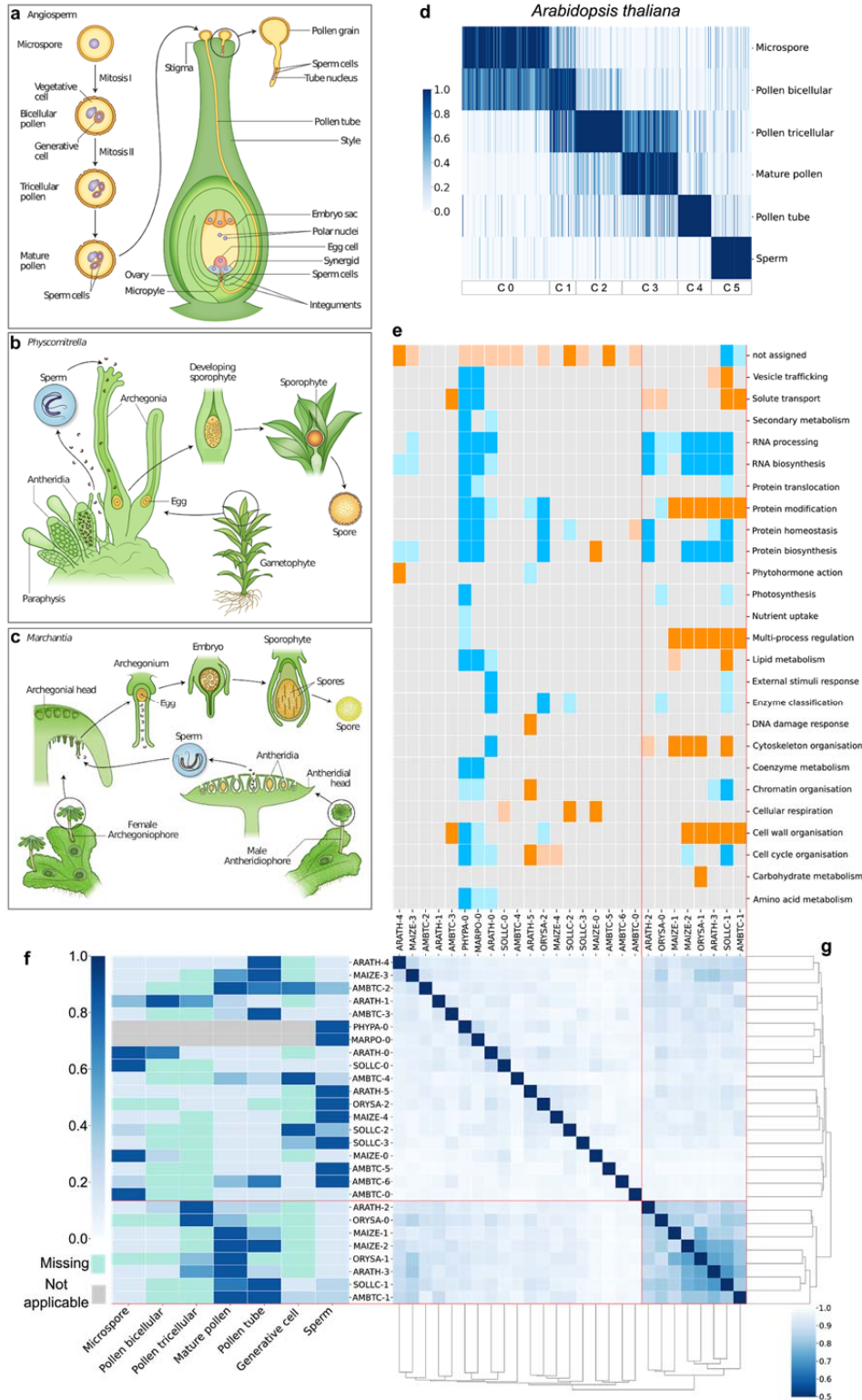
364 **Comparisons of transcriptional programs of gametes**

365 Sexual reproduction is a complex process. In diploid flowering plants involves the production of haploid
366 male and female gametes and fertilization of the female ovule by male gametes mediated by pollination
367 (Fig. 5a). The pollen delivers the sperm cell(s) to the ovary by a pollen tube, and the fertilized ovules
368 grow into seeds within a fruit (Fig. 5a). The two haploid bryophytes in our study differ in their sexual
369 reproduction. *Physcomitrella* is monoicous and bears both sperm and eggs on one individual (Fig. 5b),
370 and *Marchantia* is dioicous and bears only egg or sperm, but never both (Fig. 5c). However, both species
371 produce motile sperm that require water droplets to fertilize the egg, generating diploid zygotes. The
372 zygotes divide by mitosis and grow into a diploid sporophyte. The sporophyte eventually produces
373 specialized cells that undergo meiosis and produce haploid spores, which are released and germinate to
374 produce haploid gametophytes (Fig. 5b,c).

375 To further study whether the transcriptional programs of sexual reproduction are conserved in land plants,
376 we applied k-means clustering on the male- and female-specific genes over the RNA-seq samples
377 representing different samples of male and female organs (Supplementary Table 1). For male-specific
378 genes, the analysis assigned each sample to one or more clusters (Fig. 5d exemplifies male samples in
379 *Arabidopsis* (for other species, see Supplementary Fig. 8), with a variable number of genes assigned to
380 each cluster (Supplementary Table 10). We then inferred which biological processes were enriched in the

381 clusters (Fig. 5e), plotted an average expression profile of the genes in each cluster (Fig. 5f), and used
382 Jaccard distance to identify similar clusters across species (Fig. 5g). Interestingly, three clusters showed
383 strong similarity and were specific to pollen tricellular, mature pollen, and pollen tube for Angiosperms
384 (Fig. 5g, indicated by red lines). Functional enrichment analysis revealed that pollen tricellular, mature
385 pollen, and pollen tube samples were mainly enriched for cell wall organization, cytoskeletal
386 organization, multi-process regulation, and protein modification (supported by five species, Fig. 5e).
387 Conversely, other clusters showed enrichment for genes without assigned functions, and depletion for
388 many biological processes (Fig. 5e).

389



390

391 **Fig. 5: Comparison of male development across species.** Overview of sexual reproduction in (a) Angiosperms,
 392 (b) *Physcomitrella*, and (c) *Marchantia*. **d**, Heatmaps showing the expression of male samples genes for *Arabidopsis*
 393 *thaliana*. Genes are in columns, sample names in rows. Gene expression is scaled to range between 0-1. Darker

394 color corresponds to stronger gene expression. Bars to the bottom indicate the k-means clusters. **e**, Heatmap showing
395 enrichment (orange) and depletion (blue) of functions in the found clusters. Light colors: $p < 0.05$, dark colors: $p <$
396 0.01 . **f**, Heatmap showing the average normalized TPM value per cluster for all the species. **g**, Clustermap is
397 showing the Jaccard distance between pairs of clusters of all the species.

398 Female samples included were less diverse than male samples. In all species, each sample was assigned to
399 a cluster with exception of *O. sativa*, where ovule is divided into two clusters (Supplementary Fig. 9,
400 Supplementary Table 11). Interestingly, when we measured the Jaccard distance among all clusters
401 (including the species with one female sample), we observed no grouping of similar clusters, indicating
402 that the female gamete transcriptomes were poorly conserved (Supplementary Fig. 9). Functional
403 enrichment analysis showed enrichment mainly for not assigned functions and RNA processing, and
404 depletion for many biological processes (Supplementary Fig. 9). The *G. biloba* ovule cluster (GINBI-0,
405 ovule) showed enrichment for many functions, but ovule samples of other species did not support this
406 observation. Despite the small number of samples included these results provide evidence that female
407 gamete transcriptomes are poorly conserved across the different species analyzed.

408

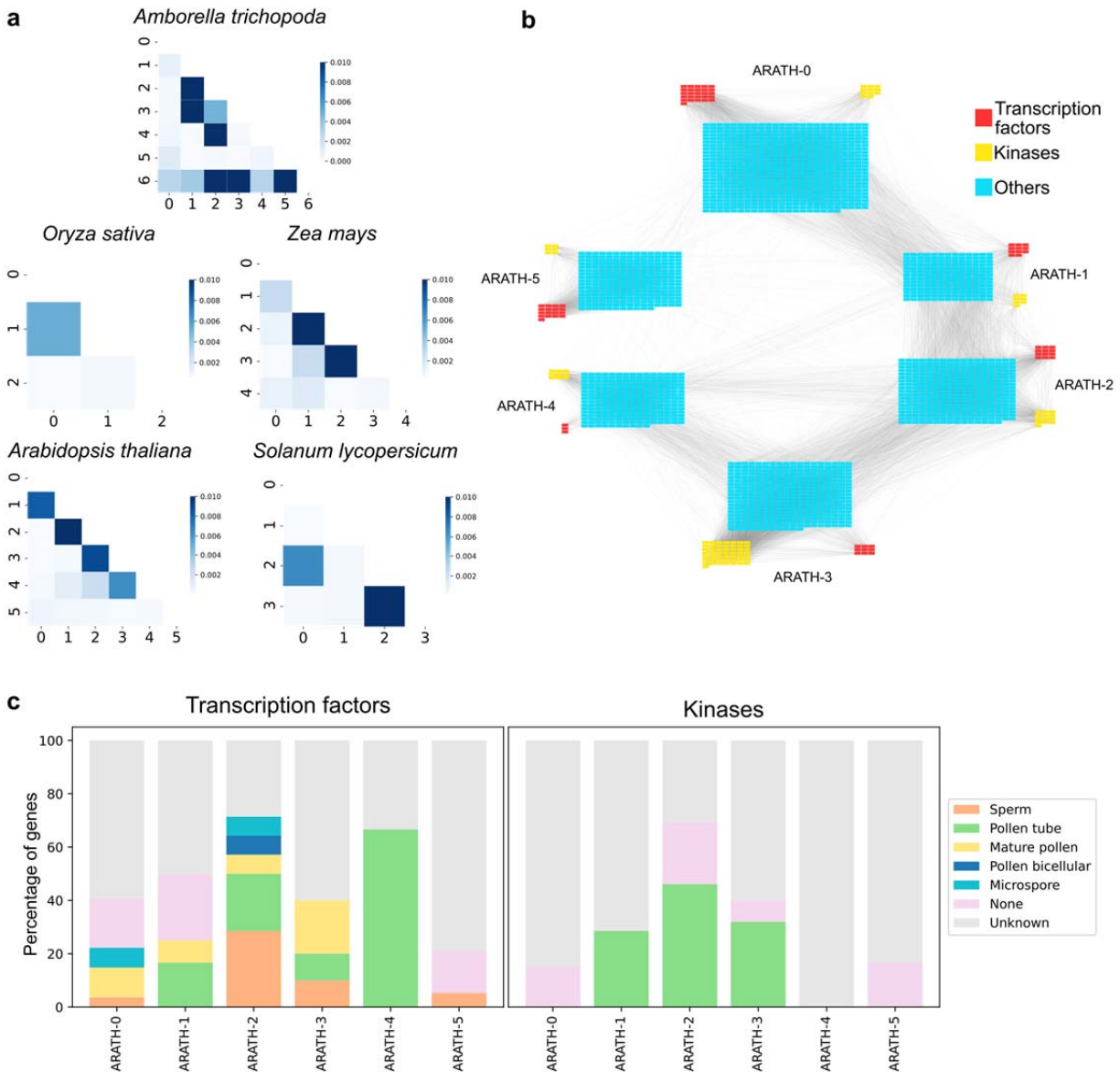
409 **Analysis of signaling networks underpinning male gametophyte development and function**

410 Gene co-expression networks help to identify sets of genes involved in related biological processes and
411 highlight regulatory relationships^{56,57}. Since we identified different gene clusters for male sub-samples
412 (see above), we decided to test whether the genes assigned to different clusters are co-expressed. For this
413 purpose, we reconstructed the co-expression networks of the ten species and analyzed whether the
414 number of observed connections was similar to the number of expected connections (see material and
415 methods). Interestingly, the clusters with expression profiles related to sperm had the least number of
416 connections with other clusters for *O. sativa*, *Z. mays*, *A. trichopoda*, and *A. thaliana* (Fig. 6a). However,
417 this pattern was not clear in *S. lycopersicum*, where the sperm cluster had connections with the cluster of
418 generative cells. Specifically, for *A. thaliana* the co-expression network revealed that cluster C5 (sperm)

419 is not well connected with other clusters (Fig. 6b), suggesting that the sperm cell transcriptome is
420 distinctive, confirming earlier observations⁵⁸⁻⁶¹. The connections between clusters followed a pattern
421 from cluster C0 to C4, which highlighted the interaction of genes among the different developmental
422 stages of male gametogenesis. The number of transcription factors and kinases present in the co-
423 expression network changed among the different clusters, where transcription factors seemed to be more
424 abundant in cluster C0 (microspore), while kinases were more abundant in cluster C3 (mature pollen)
425 (Fig. 6b, indicated by the sizes of rectangles, Supplementary Table 12).

426 Transcription factors and kinases are regulatory proteins essential for plant growth and development. To
427 uncover the regulatory mechanism underlying male gametogenesis, we analyzed all the predicted
428 transcription factors and kinases in all the male clusters of *A. thaliana*. First, we searched for all the
429 transcription factors and kinases present in the five clusters that have been characterized using
430 experimental studies with mutants (Supplementary Table 13). Then we classified the effect of each
431 mutant gene as follows: no effect related to male gametogenesis (none), no experimentally described
432 function (unknown), and important for microspore, bicellular, mature pollen, pollen tube, and sperm
433 function. Interestingly, most of the genes are described as unknown (Fig. 6c), indicating no experiments
434 associated with those genes. It is important to note that the genes classified as ‘none’ have been found to
435 have an effect in other organs, but since pollen phenotype can be easily missed, this does not rule out the
436 possibility of these genes being associated with male development. Also, many of those genes show
437 effects in roots, and it has been shown that some genes are active during tip growth of root hairs and
438 pollen tubes⁶². We observed that the transcription factors were important at different stages of male
439 development, with main phenotypes affecting pollen tube and sperm function. Conversely, kinases only
440 showed an effect on pollen tubes, which is in line with their intercellular communication involvement.
441 Interestingly, we observed that genes present in the pollen tube cluster (ARATH-4) only affected pollen
442 tube function, but pollen tube function can also be affected by genes from earlier stages of pollen
443 development (ARATH1-3). In the case of sperm function, transcription factors expressed in tricellular

444 pollen have the greatest effect, but we also observed the involvement of genes expressed in microspore,
 445 mature pollen and sperm (Fig. 6c).



446

447 **Fig. 6: A network analysis of male clusters.** **a**, Heatmaps show the number of observed connections divided by the
 448 number of expected connections. Darker colors indicate more connections between clusters. **b**, *A. thaliana* co-
 449 expression network clusters showing the edges between the different clusters (indicated as ARATH-0-5). The size of
 450 the panels indicate the number of genes in each cluster. Transcription factors, kinases, and other genes are shown in
 451 red, yellow, and blue, respectively. **c**, Percentage of genes of each *A. thaliana* male cluster. The colors indicate the
 452 different stages of male development that a given gene is known to be involved in. For example, the majority of

453 transcription factors in cluster ARATH-4 (highest expression in the pollen tube, Fig. 5f) are important for pollen
454 tube growth (green bars).

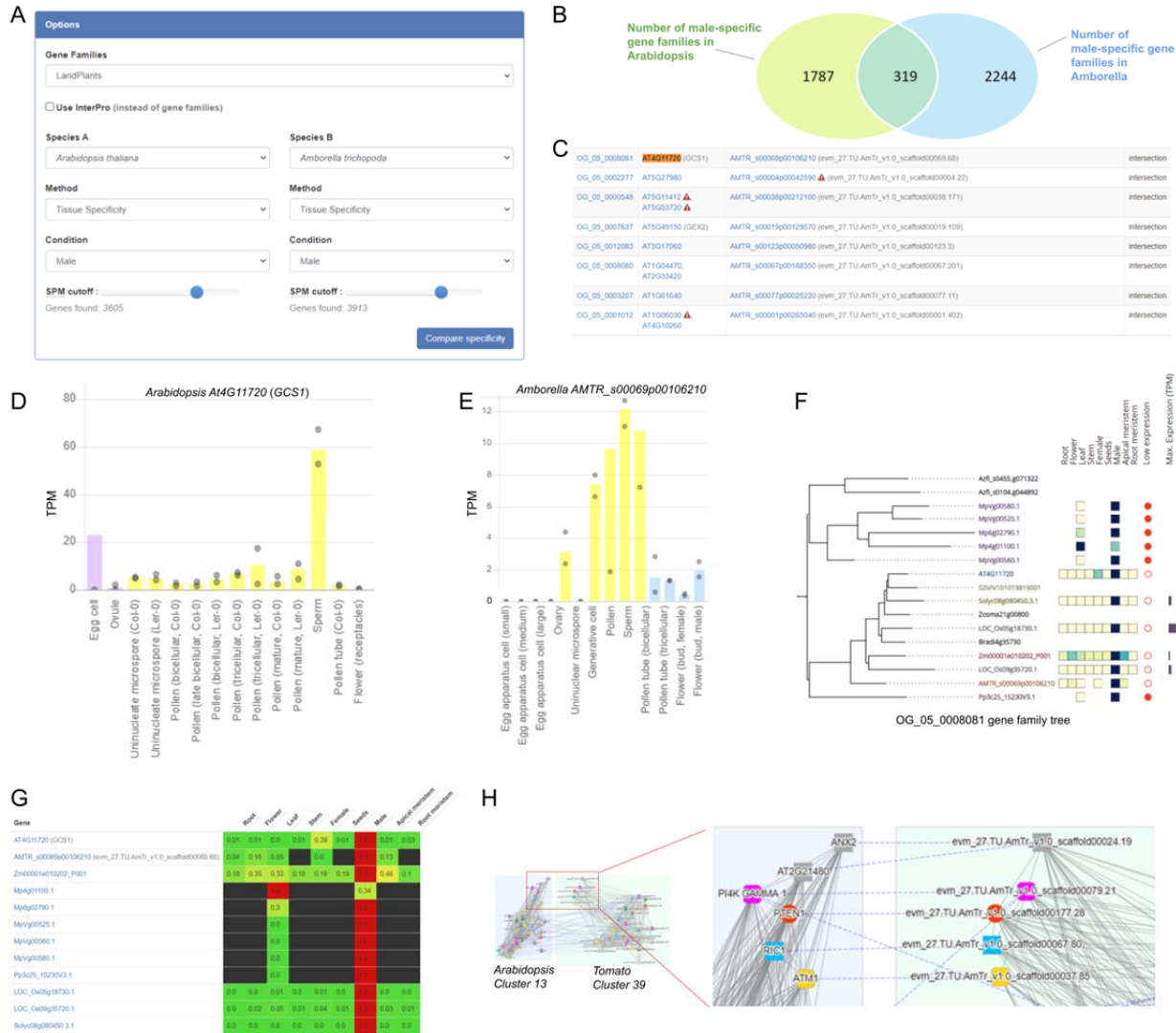
455

456 **Comparative gene expression analyses with the EVOREPRO database**

457 To provide easy access to the data and analyses generated by our consortium, we have constructed an
458 online database available at www.evorepro.plant.tools. The database is preloaded with the expression data
459 used in this study and also includes *Vitis vinifera* (eudicot, grapevine), *Chlamydomonas reinhardtii*
460 (chlorophyte), and *Cyanophora paradoxa* (glaucohyte), bringing the total number of species to 13. The
461 database can be queried with gene identifiers and sequences but also allows sophisticated, comparative
462 analyses.

463 To showcase a typical user scenario, we identified genes specifically expressed in male organs (defined
464 as, e.g., >35% reads of a gene expressed in male organs for Arabidopsis, Supplemental Figure 1). This
465 can be accomplished for one (<https://evorepro.sbs.ntu.edu.sg/search/specific/profiles>) or two
466 (https://evorepro.sbs.ntu.edu.sg/specificity_comparison/) species, where the latter option can reveal
467 specific expression profiles that are conserved across species (Fig. 7a). For this example, we selected
468 Arabidopsis and Amborella as species A and B from the drop-down menus, respectively, and used gene
469 families comprising only land plants, which uses all species found under node 8 in the species tree (Fig.
470 3a). Alternatively, the user can also select gene families constructed with seed plants (11 species found
471 under node 13, Fig. 3a) or archaeplastida (23 species found under node 1, Fig. 3a) sequences. Next, to
472 select male organs for comparisons, we specified ‘Tissue specificity’ and ‘Male’ as a method to group the
473 RNA-seq samples according to the definitions in Table 1. The slider near ‘SPM cutoff’ allows the user to
474 adjust the SPM value (the slider ranges from SPM 0.5 to 1), which interactively reveals many genes are
475 deemed organ-specific at a given SPM value cutoff. We left the slider at the default value (0.85) and
476 clicked on the ‘Compare specificities’ button. The analysis revealed that 319 gene families are expressed
477 specifically in the male organs of both Amborella and Arabidopsis (Fig. 7b), while the table below

478 showed the identity of the genes and gene families (Fig. 7c, Table S15). Interestingly, among the
479 conserved genes, we observed *GCSI/HAP2*, which is required for pollen tube guidance and fertilization
480 ⁶³. The table also contains links that redirect the user to pages dedicated to the genes and gene families.
481 For example, clicking on the Arabidopsis *GCSI/HAP2* gene identifier redirects the user to a gene page
482 containing the DNA/protein sequences (<https://evorepro.sbs.ntu.edu.sg/sequence/view/17946>), expression
483 profile (Fig. 7d), gene family, co-expression network, and Gene Ontology functional enrichment analysis
484 of the gene ⁶⁴. As expected, the interactive, exportable expression profiles confirmed that the Arabidopsis
485 *GCSI/HAP2* and the Amborella ortholog (<https://evorepro.sbs.ntu.edu.sg/sequence/view/45084>, Fig. 7e)
486 are male-specific, with the highest expression in sperm and pollen. Clicking on the gene family identifier
487 (OG_05_0008081) redirects to the gene family page
488 (<https://evorepro.sbs.ntu.edu.sg/family/view/139708>), which among others, contains an interactive
489 phylogenetic tree (Fig. 7f, <https://evorepro.sbs.ntu.edu.sg/tree/view/88288>) and heatmap (Fig 7g,
490 <https://evorepro.sbs.ntu.edu.sg/heatmap/comparative/tree/88288/row>) showcasing the male- enriched
491 expression profiles for most of the genes in this family. Therefore, this approach can be used to identify
492 conserved, organ-specific genes across two species and study family-wide expression patterns.



493

494

495

496

497

498

499

500

501

502

503

504

505

506

Fig. 7: Features of the EVOREPRO database. **a**, Compare specificities tool. The dropdown menus allow selection of the species, gene families, organs, tissues, cell types, and SPM value cutoffs. The analysis is started by clicking on the ‘Compare specificity’ button. **b**. The Venn diagram shows the number of unique and common gene families of male-specific genes in Arabidopsis and Amborella. The default SPM value cutoff of 0.85 was used for both species. **c**. The table shows the identity of genes and gene families (first column) that are specifically expressed in male organs of Arabidopsis (second column) and Amborella (third column). Each row contains a gene family, and each cell can contain multiple comma-separated genes. Red triangles containing exclamation marks indicate genes with low expression (<10TPM). **d**. Expression profile of *GCS1* from Arabidopsis. The colored columns indicate the average expression values in the different samples, while gray points indicate the minimum and maximum expression values. The y-axis indicates the TPM value. **e**. Expression profile of *GCS1*-like gene from Amborella (*AMTR_s00069p00106210*). For clarity, the gray point indicating the maximum value in the sperm sample is omitted. **f**. Phylogenetic tree of the gene family OG_05_0008081 representing *GCS1*. The branches represent genes that are color-coded by species. The heatmap to the right of the gene identifiers indicates the scaled expression

507 values in the major organ and cell types and ranges from low (yellow) to high (dark blue). Genes with TPM < 10 are
508 indicated by filled red points, while the maximum gene expression is indicated by a blue bar to the right. **g.** Heatmap
509 indicating the low (green) and high (red) expression of the *GCSI* gene family. **h.** Comparative analysis of co-
510 expression clusters significantly ($P < 0.05$) enriched for ‘pollen tube’ gene ontology term in Arabidopsis (cluster 13,
511 left) and Amborella (cluster 39, right). Nodes indicate genes, while solid gray and dashed blue edges connect co-
512 expressed and orthologous genes, respectively. We used ‘label co-occurrences’ as node options. For clarity, only
513 part of each cluster is shown.

514 Alternatively, the database can be used to identify conserved co-expression clusters of functionally
515 enriched genes. To demonstrate this tool, we navigated to
516 <https://evorepro.sbs.ntu.edu.sg/search/enriched/clusters> and entered ‘pollen’ into GO text box, selected
517 ‘pollen tube’ as query and clicked on ‘Show clusters’. The analysis revealed 5 co-expressed clusters
518 significantly ($P < 0.05$) enriched for ‘pollen tube’ gene ontology term in Arabidopsis. We clicked on one of
519 the clusters (cluster 13, <https://evorepro.sbs.ntu.edu.sg/cluster/view/113>), redirecting us to a page
520 dedicated to the cluster. As expected, the cluster is significantly ($P < 0.05$) enriched for genes involved in
521 pollen tube growth, cell wall organization and kinase activity, which are processes required to expand and
522 direct the pollen tube to the ovule. The page contains the identity of the 152 genes found in this cluster,
523 their average expression profiles, co-expression network
524 (<https://evorepro.sbs.ntu.edu.sg/cluster/graph/113>), and gene families and protein domains found in the
525 cluster.

526 Furthermore, a table labeled ‘Similar Clusters’ reveals the identity of similar (defined by Jaccard index,
527 see methods) co-expression clusters in other species, which can be used to identify functionally
528 equivalent clusters across species rapidly. To exemplify this, we first clicked on ‘Jaccard index’ table
529 header to sort the similar clusters and clicked on the ‘Compare’ link next to Cluster 39 from Amborella
530 (https://evorepro.sbs.ntu.edu.sg/graph_comparison/cluster/113/769/1). This redirected us to a co-
531 expression network page showing the genes (nodes), co-expression relationships (gray edges), and
532 orthologous genes (colored shapes of nodes connected by dashed edges) conserved in the two clusters.
533 The analysis revealed many conserved genes essential for pollen function, such as *ANX2*⁶⁵, *BUPS2*

534 (*At2g21480*)⁶⁶, *PI4K Gamma-1*⁶⁷, *PTEN1*⁶⁸, *RIC1*⁶⁹, and *ATM1*⁷⁰. To conclude, this approach can be
535 used to uncover functionally equivalent, conserved transcriptional programs.

536 **Discussion**

537 To study the evolution of plant organs and gametes, we have generated and analyzed gene expression for
538 ten land plants, comprising representatives of bryophytes, lycophytes, gymnosperms, basal angiosperms,
539 monocots and eudicots. Our analyses' main advantage is that the conclusions are drawn from comparative
540 analyses of ten species, which cover the largest collection of representatives of land plants. The
541 comparative analysis revealed that each organ type typically expressed >50% of genes, with the exception
542 of the male gametes, which showed expression of ~38% of genes, on average (Figure 1D). Conversely,
543 male gametes and roots showed the highest number (5.3% and 5.0%, respectively) of specifically
544 expressed genes (Figure 1F), suggesting that these non-photosynthesizing cell types and tissues are highly
545 unique and specialized.

546 With the surprising exception of female gametes, the corresponding transcriptomes tend to be more
547 similar across the analyzed samples (Figure 2D, Figure S3, Figure S4). Another exception is seen in the
548 leaf-like organs of bryophytes (leaflets and thallus for *Physcomitrella* and *Marchantia*, respectively),
549 indicating that these organs have evolved independently from the leaves of flowering plants or that they
550 have significantly diverged since the last common ancestor of flowering plants and bryophytes.

551 Next, we examined expression patterns of expressed gene families as a function of their age. We report a
552 clear trend of older gene families having more ubiquitous (i.e., less organ-specific) expression, while
553 younger gene families show an increasingly higher proportion of organ-specific expression (Figure 3b-c).
554 This indicates that newly-acquired genes are typically recruited to perform some specialized function in a
555 plant organ, tissue, or cell type, rather than being integrated into fundamental biological pathways. As
556 expected, male gametes show the highest expression of the youngest genes (Figure 3d-e, Figure S7),

557 which is in line with previous studies^{42,71}. Interestingly, *Physcomitrella* gametes did not show this
558 pattern, which is a finding that warrants further studies.

559 To study how new functions were gained or lost as the organs and gametes evolved, we studied which
560 phylostrata are enriched or depleted in the different organs (Figure 4a). Interestingly, we observe a
561 significant enrichment for gene families that appeared long before the corresponding organ (Figure 4a),
562 showing that the establishment of organs relies heavily on the co-option of existing genetic material, as
563 suggested previously^{45,50}. Flowers (appearance in angiosperms), stems (appearance in vascular plants)
564 and roots (appearance in vascular/seed plants) show similar patterns of enrichment and depletion of genes
565 (Fig. 4a). This is surprising, as these organs appeared at different stages of plant evolution, which
566 suggests that the co-option underlying the establishment of novel organs follows a similar pattern of gene
567 gains and losses. Based on the diverse patterns of gains and losses of organ-specific gene families (Figure
568 4b) we conclude that monocot-specific families show substantial net gains in genes that are specifically
569 expressed in male gametes, seeds, stems, roots or in apical and root meristems (Figure 4b), suggesting
570 that during monocots evolution organ-specific transcriptomes was enriched with novel functions.
571 Surprisingly, eudicots show an opposite pattern, exhibiting more net losses of organ-specific families in
572 flowers, female and male gametes, leaves, stems, roots, and apical meristems (Figure 4b). This surprising
573 pattern of loss of functions in eudicots merits investigation by further analysis, which is made possible by
574 identifying the corresponding gene families (Table S8) and genes (Table S6).

575 Our comparative analysis of male gamete development reveals that transcriptional programs of mature
576 pollen form well-defined clusters and are thus conserved across species (Figure 5f-g). The mature pollen
577 clusters are enriched for processes related to signaling (protein modification comprising protein kinases)
578 and cell wall remodeling (Figure 5e), which are likely representing processes mediating pollen
579 germination, pollen tube growth, and sperm cell delivery. Conversely, the earlier stages of male gamete
580 development showed less defined clusters and enrichment for genes with unknown function (bin 'not
581 assigned', Figure 5e), suggesting that the processes taking place in the early stages of pollen development

582 are yet to be uncovered. Furthermore, the female gametes show poor clustering, indicating overall low
583 conservation of the transcriptional programs and enrichment of genes with unknown function for most
584 clusters (Figure S9c). These results indicate that genes expressed during early male gamete and female
585 gamete formation warrant closer functional analysis, which is now made possible by our identification of
586 these genes (Table S10-11). Of particular interest are the male-specific transcription factors and kinases
587 that we identified (Figure 6c), assumingly involved in various stages of pollen development and function
588 (Table S13). As a large fraction of these genes are not yet characterized, their involvement in male
589 gametogenesis and function should be further investigated.

590 To provide easy access to the 13 expression atlases, organ-specific genes, functional enrichment analyses,
591 co-expression networks, and various comparative tools, we provide the EVOREPRO database
592 (www.evorepro.plant.tools) to the community (Figure 7). This database represents a valuable resource for
593 further study and validation of key genes involved in organogenesis and land plants reproduction.

594

595 **Methods**

596 **Physcomitrella growth conditions, RNA isolation and sequencing**

597 *Plant growth*

598 The Grandsden wild-type strain from *P. patens* Bruch & Schimp⁷² was used for this study. To initiate
599 plant growth and culture, 3 mature sporophytes were sterilized using a 5% commercial bleach solution for
600 5 minutes and rinsed twice in molecular grade water. Sterilized sporophytes were then broken using a
601 pipette tip and diluted into 5mL molecular grade water. Spore containing solution was then distributed
602 into 4 sterile peat pellets (Jiffy-7, Jiffy Products International) and two 9 cm Petri dishes containing
603 KNOPS medium (Reski and Abel, 1985) supplemented with 0.5 g/l ammonium tartrate dibasic (Sigma-

604 Aldrich Co). Petri dishes were kept at 25°C, 50% humidity, and 16 h light (light intensity 80 $\mu\text{mol}/\text{m}^2/\text{s}$).
605 Protonema samples were collected 10 days after spore germination.

606 Plants in Phytatray™ II (Sigma-Aldrich Co) containing 4 sterile peat pellets (Jiffy-7, Jiffy Products
607 International) were grown for 6-8 weeks at 25°C, 50% humidity, and 16 h light (light intensity 80
608 $\mu\text{mol}/\text{m}^2/\text{s}$). Water was supplied to the bottom of each box. Leave samples were collected after 6 weeks,
609 prior to induction of gametangia development. For gametangia and sporophyte development, water was
610 again supplied to the bottom of each box containing four pellets and were transferred to 17°C, 8 h light,
611 and 50% humidity (light intensity 50 $\mu\text{mol}/\text{m}^2/\text{s}$) to induce the development of reproductive structures⁷³.
612 Gametangia samples (archegonia, paraphysis and sperm cell packages) were collected 15 days after
613 reproductive induction. Antheridia samples were collected at several time points during their
614 development. Further development of the sporophyte was conducted under these conditions and
615 sporophyte samples were collected at different time points during sporophyte development. S1
616 sporophytes were collected 7 days after sperm cell (SC) release, S2 sporophytes 15 days after SC release,
617 S3 sporophytes 20 days after SC release (green spore capsules) and SM samples 28 days after SC release
618 (brown spore capsules).

619

620 *Sample preparation and sequencing*

621 Leaves, protonema and sporophytes were collected under a stereoscope using tweezers, placed in 2.5 μL
622 of RLT+ buffer (Qiagen), and shock frozen in liquid nitrogen. Before RNA-seq library preparation, these
623 samples were mechanically disrupted using sterile pellet pestles (Z359947, Sigma-Aldrich Co).
624 Antheridia, archegonia, paraphysis and sperm cell packages were collected using a Yokogawa CSU-W
625 Spinning Disk confocal with 10x 0.25NA objective, using the brightfield channel and an Andor Zyla 4.2
626 sCMOS camera. For each of these samples the plants were prepared under a stereoscope, isolating the
627 whole gametangia for ca. 10 shoots. They were placed in 20 μL of molecular grade water on a glass slide.

628 Using a cover slip the gametangia were disrupted into individual antheridia by applying slight pressure.
629 Slides were placed under a microscope and specific organs were identified and collected, using an
630 Eppendorf CellTram® Air/Oil/vario micromanipulator with glass capillaries (borosilicate glass with fire
631 polished ends, without filament GB100-9P) pulled with a Narishige PC-10 puller. Then they were
632 transferred to another clean slide, and subsequently excessive liquid containing possible contaminations,
633 such as cell debris, was removed. For paraphysis samples 8-15 individual paraphysis were collected
634 directly into 2 uL of RLT+ buffer and flash frozen in liquid nitrogen. For antheridia samples 5 to 20
635 individual antheridia of each specific stage (9 to 15 days after induction, distinguished by size) were
636 collected and then burst under a microscope by applying pressure on a cover slip applied to the samples
637 on the slide. The slide was washed with 4 uL of RLT+ buffer and the buffer transferred into a PCR tube,
638 subsequently flash frozen. Archegonia samples were prepared from 3-5 archegonia following the same
639 procedure. Released sperm cell packages (2-5 per sample) were collected from gametangia preparations
640 (as described above; antheridia 15 days after induction) without clean up, transferred into a tube with 2 uL
641 of RLT+ buffer, flash frozen in liquid nitrogen and subsequently used for RNA-seq library preparation.
642 RNA-seq library preparation for all samples was performed as described in ⁷⁴, with the addition of mixing
643 the PCR tubes on a Thermomixer C (Eppendorf) every 15 minutes at 200 rpm for 1 min during the RT
644 step. Libraries were sequenced on a NextSeq500 instrument with single-end 75 bp read length (SE75).

645

646 **Marchantia growth conditions, RNA isolation and sequencing**

647 Male accession of *Marchantia polymorpha* L., Takaragaike (Tak)-1 was grown on vermiculite under a
648 long-day condition (16/8 h day/night) at 22 °C. To induce sexual reproduction, thalli developed from
649 gemmae were transferred to a far-red light (700 – 780 nm, 44.3 $\mu\text{mol photons m}^{-2} \text{s}^{-1}$) supplemented
650 light condition using LabLEDs (RHENAC GreenTec Ag). Sperms were released from antheridiophores
651 by applying ddH₂O supplemented with RNasin® Ribonuclease Inhibitor (1 u/μl, Promega), collected in a

652 1.5 mL tube, and pelleted by centrifugation at 3,000 g for 5 min at 4 °C. RNA-seq libraries were
653 generated from total RNA of isolated *M. polymorpha* sperm using Smart-seq2⁷⁵ using independent
654 biological replicates. The libraries were sequenced on an Illumina HiSeq 2500 using 125 bp paired-end.

655

656 **Amborella growth conditions, RNA isolation and sequencing**

657 *Plant material and isolation procedures*

658 *Amborella trichopoda* male flowers were harvested from a male plant growing in the Botanical Garden in
659 Bonn (Germany), in a shaded place inside a greenhouse under controlled conditions of 16-18°C, constant
660 humidity of 66% and 12-hour photoperiods. Buds and fully opened male flowers were gathered in 50 ml
661 FalconTM conical tubes (Thermo Fisher), placed without lid in a hermetically sealed plastic box containing
662 a bed of silica gel.

663 Uninucleated microspores (UNM) were isolated at room temperature from flower buds of 4.5 mm length,
664 as these were found to contain 98% uninucleated microspores. In brief, three samples with each 5 g buds
665 were homogenized in 0.1 M mannitol and filtered with a 70-micron pore size PET strainer (PluriSelect).
666 The filtered solution was processed by subsequent steps of percoll gradient separation, washing and
667 centrifugation, as described previously⁷⁶.

668 *Amborella* generative cells (GC) were obtained from mature pollen grains that were purified like
669 described previously⁷⁷. Per replicate, 50 mg pollen was resuspended in 1 ml pollen germination medium
670 and transferred into a 1.5 ml vial containing glass beads (0.4 – 0.6 mm). The vial was vortexed
671 continuously at 2,200 rpm for 4 minutes to crack the pollen grains and release its contents. The solution
672 was filtered using a 15-micron PET strainer (PluriSelect). To stain the nuclei, a final concentration of 10X
673 SYBR Green I was added and GCs were identified using an inverted microscope (Nikon) equipped with
674 high-resolution 20X and 40X objectives suitable for fluorescent applications and suitable filters for SYBR

675 Green I (497 nm excitation; 520 nm emission). For RNA-seq, three replicates of each 140 GC were
676 harvested manually using an Eppendorf CellTram.

677 *Amborella* sperm cells (SC) were isolated at room temperature by adapting a method described for tomato
678 sperm cell isolation⁷⁸. In brief, three replicates with each 50 mg purified pollen were germinated as
679 described⁷⁷. 16 hours after germination, the medium was removed by filtration using a 15-micron PET
680 strainer (PluriSelect) and the pollen tubes were incubated for 10 min in a 15% mannitol solution with
681 0.4% cellulase “Onozuka” R-10 and 0.2% macerozyme R-10 to release the sperm cells. The mixture was
682 re-filtered using a 15-micron PET strainer and loaded on 5 ml 23% Percoll in 0.55 M mannitol and
683 centrifuged at 1,000 x g for 30 min. Approximately 1 ml with SC, floating on the surface of the Percoll
684 gradient, were harvested, washed with 1 ml RNeasy Protect[®] Cell Reagent (Qiagen) and centrifuged for 10
685 min at 2,500 x g. 50 µl of SC-enriched pellet (approximately 250 sperm cells each replicate) was used for
686 RNA-seq library preparation.

687 Isolation and sampling of *Amborella* ovaries, egg apparatus cells, pollen tubes, pollen grains as well as
688 male and female flowers, tepals, roots and leaves was done as described in previous studies^{77,79}.

689 *RNA isolation and sequencing*

690 RNA isolation from uninucleated microspores was performed by using the Spectrum[™] Plant Total RNA
691 Kit (Sigma-Aldrich) according to manufacturer’s instructions. Total RNA from *Amborella* generative
692 cells and sperm cells was extracted according to the “Purification of total RNA from animal and human
693 cells” protocol of the RNeasy Plus Micro Kit (QIAGEN, Hilden, Germany). In brief, cells were stored
694 and shipped on dry ice. After adding RLT Plus containing β-mercaptoethanol the samples were
695 homogenized by vortexing for 30 sec. Genomic DNA contamination was removed using gDNA
696 Eliminator spin columns. Next ethanol was added and the samples were applied to RNeasy MinElute spin
697 columns followed by several wash steps. Finally total RNA was eluted in 12 µl of nuclease free water.

698 Purity and integrity of the RNA was assessed on the Agilent 2100 Bioanalyzer with the RNA 6000 Pico
699 LabChip reagent set (Agilent, Palo Alto, CA, USA).

700 The SMARTer Ultra Low Input RNA Kit for Sequencing v4 (Takara) was used to generate first strand
701 cDNA from 2.5 ng UNM, 0.8 ng GC and 0.5 ng SC total RNA. Double stranded cDNA was amplified by
702 LD PCR (10 for UNM, 13 cycles for GC and 15 cycles for SC) and purified via magnetic bead clean-up.
703 Library preparation was carried out as described in the Illumina Nextera XT Sample Preparation Guide
704 (Illumina, Inc., San Diego, CA, USA). 150 pg of input cDNA were tagged by the Nextera XT
705 transposome. The products were purified and amplified via a limited-cycle PCR program to generate
706 multiplexed sequencing libraries. For the PCR step 1:5 dilutions of index 1 (i7) and index 2 (i5) primers
707 were used. The libraries were quantified using the KAPA SYBR FAST ABI Prism Library Quantification
708 Kit. Equimolar amounts of each library were used for cluster generation on the cBot (TruSeq SR Cluster
709 Kit v3). The sequencing run was performed on a HiSeq 1000 instrument using the indexed, 2x100 cycles
710 paired end (PE) protocol and the TruSeq SBS v3 Kit. Image analysis and base calling resulted in .bcl
711 files, which were converted into .fastq files by the CASAVA1.8.2 software. Library preparation and
712 RNA-seq were performed at the service facility “Center of Excellence for Fluorescent Bioanalytics
713 (KFB)” (Regensburg, Germany; www.kfb-regensburg.de).

714

715 **Arabidopsis growth conditions, RNA isolation and sequencing**

716 *Arabidopsis thaliana* accession Columbia-0 (Col-0) plants were grown in controlled-environment
717 cabinets at 22°C under illumination of 150 µmol/m²/sec with a 16-h photoperiod. Mature pollen grains
718 (MPG) were harvested from open flowers of 5 to 6-week old plants by shaking into liquid medium (0.1 M
719 D-mannitol) as described previously⁷⁹. Microspores and developing pollen grains were released from
720 anthers of closed flower buds and purified by Percoll density gradient centrifugation as described^{76,80}.

721 Populations of spores at five stages of development were isolated: uninucleate microspores (UNM),
722 bicellular pollen (BCP), late bicellular pollen (LBC), tricellular pollen (TCP) and mature pollen (MPG).

723 For semi in vivo pollen tube growth, a transgenic marker line harboring MGH3p::MGH3-eGFP and
724 ACT11p::H2B-mRFP²¹ was used to pollinate WT emasculated pistils. After 2 hours, the pollinated pistil
725 was excised and placed on double sided tape. The excised pistil was then cut at the junction of style and
726 ovary and placed gently on solidified agarose pollen germination medium⁸¹. The pistil was incubated for
727 an additional 4 hours for the pollen tubes to emerge from the cut end of the style. The pollen tubes were
728 harvested using a 25G needle and immediately frozen in liquid nitrogen and subsequently used for the
729 RNA-seq library preparation as described in⁷⁴.

730 Total RNA was isolated from each sample using the RNeasy Plant Kit (Qiagen) according to the
731 manufacturer's instructions. RNA was DNase-treated (DNA-freeTM Kit Ambion, Life Technologies)
732 according to the manufacturer's protocol. RNA yield and purity were determined spectrophotometrically
733 and using an Agilent 2100 Bioanalyzer. cDNA was prepared using a slightly modified SmartSeq2
734 protocol in which cDNA is synthesized from poly(A)+ RNA with an oligo(dT)-tailed primer^{75,82}. The
735 final libraries were prepared using a low-input Nextera protocol⁸³. Libraries were sequenced on a
736 NextSeq500 instrument with single-end 75 bp read length (SE75).

737 A transgenic line expressing EC1.1p:NLS-3xGFP was cultured and used for Arabidopsis egg cell
738 isolation as previously described⁸⁴. Three replicates of 25 to 30 pooled egg cells were used for RNA
739 extraction, RNA-seq library preparation and Illumina Next Generation Sequencing⁸⁵.

740

741 **Tomato growth conditions, RNA isolation and sequencing**

742 *Solanum lycopersicum* (tomato accession Nagcarlang, LA2661) seeds were obtained from the Tomato
743 Genetics Resource Center (TGRC, <https://tgrc.ucdavis.edu/>) and grown in the Brown University

744 Greenhouse (Providence, RI, USA). Dry pollen grains were collected from stage 15 flowers⁸⁶ into 500µl
745 eppendorf tubes. Pollen tubes were grown in 300µl of pollen growth medium in a 750µl eppendorf tube
746 that was incubated in a 28°C water bath. Pollen tubes were grown at a density of ~1000 pollen grains/µl.
747 The pollen germination medium⁸⁷ comprised 24% (w/v) polyethylene glycol (PEG) 4000, 0.01% (w/v)
748 boric acid, 2% (w/v) Suc, 20 mM MES buffer, pH 6.0, 3 mM Ca(NO₃)₂·4H₂O, 0.02% (w/v)
749 MgSO₄·7H₂O, and 1 mM KNO₃. Pollen tubes were grown for 1.5 hours, 3 hours, or 9 hours before they
750 were collected by centrifugation (1000 x g) for 1 minute. Pollen germination medium was carefully
751 removed by pipetting to avoid disrupting the loose pollen tube pellet. Independent pollen collections were
752 made for each of three biological replicates at each time point. Eppendorf tubes containing pollen tubes
753 were immediately flash frozen in liquid N₂, then stored at -80°C, or put directly on a dry-ice cooled metal
754 block for cell disruption by grinding with a frozen plastic pestle (Kontes). Total RNA was extracted using
755 the RNeasy Plant Kit (Qiagen). RNA samples were evaluated by Agilent 2100 Bioanalyzer (Brown
756 University Genomics Core Facility) before RNA-seq library preparation (polyA selection) and Illumina
757 HiSeq, (150bp, paired end) sequencing were performed by Genewiz (South Plainfield, New Jersey, USA).

758

759 **Maize growth conditions, RNA isolation and sequencing**

760 Maize plants (inbred line B73) were grown in an air-conditioned greenhouse at 26°C under illumination
761 of about 400 µmol/m²/sec with a 16-h photoperiod (21°C night temperature) and air humidity between
762 60-65%. Fresh mature pollen grains were harvested as described⁸⁸. Pollen tubes were germinated and
763 grown for 2 hours *in vitro* using liquid pollen germination medium⁸⁹. Total RNA was extracted from
764 each three biological replicates of 100 mg pollen grains/pollen tubes by using a SpectrumTM Plant Total
765 RNA Kit (Sigma-Aldrich) according to manufacturer's instructions. 250 ng of total RNA was each used
766 for library construction. RNA-seq was carried out as described in the Illumina TruSeq Stranded mRNA
767 Sample Preparation Guide for the Illumina HiSeq 1000 System (Illumina) and the KAPA Library

768 Quantification Kit (Kapa Biosystems). Data from sperm cells, egg cells and various zygote stages were
769 taken from published data ⁸⁸.

770

771 **Compiling gene expression atlases**

772 RNA data of different samples from nine species (*Physcomitrium patens*, *Marchantia polymorpha*,
773 *Ginkgo biloba*, *Picea abies*, *Amborella trichopoda*, *Oryza sativa*, *Zea mays*, *Arabidopsis thaliana*,
774 *Solanum lycopersicum*) were grouped in ten different classes (flower, female, male, seeds, spore, leaf,
775 stem, apical meristem, root meristem, root) (Table 1, Supplementary Table 1). For male and female
776 reproductive organs samples we also included different sub-samples (female: egg cell, ovary, ovule;
777 Male: microspore, bicellular pollen, tricellular pollen, mature pollen, pollen tube, generative cell, sperm)
778 for each species (Table 1, Supplementary Table 1). A total of 4,806 different RNA sequencing samples
779 were used, from which 4,672 were downloaded from the SRA database and 134 obtained from our
780 experiments (see above). Publicly available RNA-seq experiments data were downloaded from ENA ⁹⁰, as
781 described in CoNekt-Plants ⁶⁴. Proteomes and CDSs of each species were downloaded from different
782 sources (Supplementary Table 14). The raw reads of each sample were mapped to the coding sequences
783 (CDS) with Kallisto v.0.46.1 ²⁶ to obtain transcripts per million (TPM) gene expression values. If the
784 reads came from single cell samples (egg cell, ovule, sperm, generative cell), we removed the samples
785 that have <1M reads mapped, and for the other samples we removed those with <5M reads mapped
786 (Supplementary Table 1). All those samples were used to calculate Highest Reciprocal Rank (HRR)
787 networks, where two genes with HRR<100 were connected ⁹¹. For comparative expression analysis, an
788 additional filter was applied by keeping only samples with a Pearson correlation coefficient (PCC) ≥ 0.8
789 to at least one other sample of the same type (e.g. flower to flower) (Supplementary Table 1).
790 Additionally, we included the expression matrix of *Selaginella moellendorffii* which has 18 samples
791 (Supplementary Table 1), and exclusively for the database (see section Constructing the co-expression

792 network and establishing the EVOREPRO database) the expression matrices of two unicellular algae
793 (*Chlamydomonas reinhardtii* and *Cyanophora paradoxa*) and *Vitis vinifera*⁹². Finally, genes with median
794 expression levels >2 TPM were considered as expressed⁹³. All expression matrices are available for
795 download from <http://www.gene2function.de/download.html>.

796

797 **Identifying sample-specific genes**

798 Sample-specific genes based on expression data were detected by calculating the specificity measure
799 (SPM), using a similar method as described in⁹⁴. For each gene, we calculated the average TPM value in
800 each sample (e.g., root, leaf, seeds). Then, the SPM value of a gene in a sample was computed by dividing
801 the average TPM in the sample by the sum of the average TPM values of all samples. The SPM value
802 ranges from 0 (a gene is not expressed in a sample) to 1 (a gene is fully sample-specific). To identify
803 sample-specific genes, for each of the ten species, we first identified a SPM value threshold above which
804 the top 5% SMP values were found (Supplementary Fig. S1, red line). Then, if a gene's SPM value in a
805 sample was equal to or larger than the threshold, the gene was deemed to be specifically expressed in this
806 sample.

807

808 **Similarity of sample-specific transcriptomes between samples and species**

809 To estimate whether sample-specific transcriptomes (see above) are similar, we calculated Jaccard
810 distance d_j between orthogroup sets. These orthogroup sets were found by identifying the orthogroups of
811 sample-specific genes per each species. Then pairwise d_j was calculated for all the samples and used as
812 input for the clustermap. The d_j ranges between 0 (the two sets of orthogroups are identical) to 1 (the two
813 sets have no orthogroups in common).

814 To estimate whether a species' sample-specific transcriptome was significantly similar to a corresponding
815 sample in the other species (e.g. Arabidopsis root vs. rice root, tomato root), we tested whether the
816 d_j values comparing the same sample were smaller (i.e. more similar) than d_j values comparing the
817 sample to the other samples (e.g., Arabidopsis root vs. rice flower, rice leaf, tomato flower, tomato leaf).
818 We used Wilcoxon rank-sum to obtain the p-values, which were adjusted using a false discovery rate
819 (FDR) correction ⁹⁵.

820

821 **Phylogenomic and phylostratigraphic analysis**

822 We used proteomes of 23 species representing key phylogenetic positions in the plant kingdom (see
823 Supplementary Table 14), to construct orthologous gene groups (orthogroups) with Orthofinder v2.4.0 ⁹⁶,
824 where Diamond v0.9.24.125 ⁹⁷ was used as sequence aligner. A species tree based on a recent phylogeny
825 including more than 1000 species ⁹⁸ was used for the phylostratigraphic analysis. The phylostratum (node)
826 of an orthogroup was assessed by identifying the oldest clade found in the orthogroup ⁹⁹ using ETE v3.0
827 ¹⁰⁰. To test whether a specific phylostratum is enriched in a sample, we randomly selected (without
828 replacement) the number of observed sample-specific genes 1000 times. The empirical p-values were
829 obtained by calculating whether the observed number of gene families for each phylostratum was larger
830 (when testing for enrichment) or smaller than (testing for depletion) than the number obtained from the
831 1000 sampling procedure. The p values were FDR corrected ⁹⁵.

832

833 **Transcriptomic age index calculation**

834 Transcriptome age index (TAI) is the weighted mean of phylogenetic ranks (phylostrata) and we
835 calculated it for every sample ⁷¹. We used the species tree from ⁹⁸. The nodes in the tree were assigned
836 numbers ranging from 1 (oldest node) to 22 (youngest node, Fig. 3a) by traversing the tree using ETE

837 v3.0 (Huerta-Cepas et al. 2016) with default parameters. The age (phylostratum) of an orthogroup and all
838 genes belonging to the orthogroup, were derived by identifying the last common ancestor found in the
839 orthogroup using ETE v3.0¹⁰⁰. In the case of species-specific orthogroups the age of the orthogroup was
840 assigned as 23. Finally, all genes with TPM values <2 were excluded and the TAI was calculated for the
841 remaining genes by dividing the product of the gene's TPM value and the node number by the sum of
842 TPM values.

843

844 **Functional annotation of genes and identification of transcription factor and kinase families**

845 The proteomes of the ten species included in the transcriptome dataset were annotated using the online
846 tool Mercator4 v2.0 ([https://www.plabipd.de/portal/web/guest/mercator4/-](https://www.plabipd.de/portal/web/guest/mercator4/-/wiki/Mercator4/recent_changes)
847 [/wiki/Mercator4/recent_changes](#)). This tool assigns Mapman4 bins to genes¹⁰¹. Transcription factors and
848 kinases were predicted using iTAK v1.7a¹⁰². Additional transcription factors were identified using the
849 online tool PlantTFDB v5.0 (<http://planttfdb.cbi.pku.edu.cn/prediction.php>)¹⁰³.

850

851 **Functional enrichment analysis**

852 Functional enrichment of the list of sample-specific and cluster-specific genes of each species, and genes
853 gained in each node, was calculated using the bins predicted with Mercator 4 v2.0. Briefly, for a group of
854 m genes (e.g., genes specifically expressed in Arabidopsis root), we first counted the number of Mapman
855 bins present in the group, and then evaluated if these bins were significantly enriched or depleted by
856 calculating an empirical p -value. The empirical p -value that tests whether a Mapman bin (term) is
857 enriched in a collection of m genes is defined as:

$$858 \quad P - \text{value}_{\text{term}} = \frac{\sum_{n=1}^N I(\text{pred}_{\text{observed}} \leq \text{pred}_{\text{sampled}})}{N}$$

859 Where $pred_{observed}$ is the number of times a term is observed, $pred_{sampled}$ is the number of times the
860 term is observed when the terms of m genes are randomly sampled (without replacement) from the all
861 genes in the genome. N is the number of permutations, which was set to 1000. I is an indicator function,
862 which takes a value of 1 when the event (in this case $pred_{observed} \leq pred_{sampled}$) is true, and 0 when it
863 is not. For functional depletion analysis a similar approach was followed, with I taking a value of 1 when
864 $pred_{observed} \geq pred_{sampled}$. To account for multiple hypothesis testing, we applied a false discovery
865 rate (FDR) correction to the p-values⁹⁵. Transcription factor and kinase enrichment was calculated
866 following the same procedure.

867

868 **Identification of orthogroup expression profiles**

869 In order to analyse the expression profiles at phylostrata level, orthogroups were classified as ‘sample-
870 specific’, ‘ubiquitous’, and ‘not conserved’. ‘Sample-specific’ orthogroups are orthogroups containing
871 sample-specific genes and can be sub-classified according to the organ (flower-, female-, male-, seeds-,
872 spore-, leaf-, apical meristem-, stems-, root meristem-, root-specific). ‘Ubiquitous’ are orthogroups that
873 are expressed in different samples for each species, i.e., they do not show a ‘sample-specific’ expression
874 profile. ‘Not conserved’ are orthogroups that have different sample-specific expression profiles in
875 different species (e.g., orthogroups containing root-specific genes for *Arabidopsis* and male-specific
876 genes for *Solanum*). Only orthogroups with species with sufficient expression data were used. More
877 specifically, we only analyzed orthogroups that were: i) species-specific with transcriptome data or, ii)
878 contained at least two species with transcriptome data. To identify sample-specific orthogroups, we
879 required, iii) >50% of genes of the orthogroup should support the expression profile, iv) >=50% of the
880 species with transcriptome data present in the node should support the expression profile.

881

882 **Gene enrichment analysis per phylostrata**

883 In order to analyse gene enrichment of specific samples across the different phylostrata in the species tree
884 (Fig. 3a), we used all the sample-specific genes of the ten species included. For each species and for each
885 defined sample (ubiquitous, flower, female, male, seeds, spore, leaf, stem, apical meristem, root
886 meristem, root) we counted the number of genes present in each node of the species tree, and then
887 evaluated if the number of sample-specific genes were significantly enriched or depleted by calculating an
888 empirical p-value as described for functional enrichment analysis. Then, we evaluated each sample and
889 counted the number of species that show significant enrichment/depletion ($p < 0.05$) in each node of the
890 species tree. We obtained a normalized value per each node by calculating the difference of species
891 showing enrichment and species showing depletion and dividing it by the total number of species that
892 show enrichment/depletion. These results were used to plot a heatmap using the seaborn python package
893 ¹⁰⁴.

894

895 **Gene family comparisons**

896 For each sample-specific (flower, female, male, seeds, spore, leaf, stem, apical meristem, root meristem,
897 root) and ubiquitous expression profiles we mapped loss and gain of organ-specific gene families onto the
898 species tree (Fig. 3a). All the orthogroups classified as sample-specific (see above) were analysed
899 independently and gain and loss was computed using the approach described in ¹⁰⁵ with ETE v3.0 ¹⁰⁰.
900 Briefly, a gene family gain was inferred at the last common ancestor of all the species included in the
901 family and a loss when a species did not have orthologs in the particular gene family. Groups of
902 monophyletic species that have lost the gene were counted as one loss. Then, we collapsed the values of
903 the nodes of the species tree to fit the different clades included (Fig. 4b), and we calculated the difference
904 between the total gains and the total losses to obtain an absolute value for each node. The values of each
905 expression profile were normalized dividing the values by the maximum absolute value in a way that we
906 got a range from -1 to 1 (negative values for losses and positive values for gains). Finally, per each

907 expression profile (ubiquitous, flower, female, male, seeds, spore, leaf, stem, apical meristem, root
908 meristem, root) a graphical representation of the different clades showing the nodes with a intensity of
909 color proportional to the normalized values of gains and losses was plotted using ETE v3.0¹⁰⁰.

910

911 **Identification of gamete-specific transcriptional profiles by clustering analysis**

912 We analyzed the male and female sample-specific genes and their different sub-samples (Supplementary
913 Table 1), to identify transcriptional profiles by clustering analysis. For the clustering analysis we only
914 included species with at least 2 subsamples (*Amborella trichopoda*, *Oryza sativa*, *Zea mays*, *Arabidopsis*
915 *thaliana*, *Solanum lycopersicum*). The male samples were divided into: microspore, bicellular pollen,
916 tricellular pollen, mature pollen, pollen tube, generative cell, and sperm cell for Angiosperms; and sperm
917 for bryophytes. The female samples were divided into egg cell, ovary, and ovule. For each gene, the
918 average TPM in each sub-sample was calculated, and the average TPM values were scaled by dividing
919 with the highest average TPM value for the gene. The k-means clustering method from the sklearn.cluster
920 package¹⁰⁶ was used to fit the scaled average TPM values to the number of clusters (k) ranging from 1 to
921 20. The optimal number of k for each species was estimated by using the elbow method, where k that
922 produced a sum of squared distances < 80% of $k=1$ was selected (Supplementary Fig. 10). Seaborn¹⁰⁴
923 python package was used for plotting the figures.

924

925 **Constructing the co-expression network and establishing the EVOREPRO database**

926 Coexpression networks were calculated by using Highest Reciprocal Rank (HRR) value⁹¹, which is a
927 distance-based metric that ranges from 0 (two genes are strongly coexpressed) to 100 (two genes are
928 weakly coexpressed). The networks were constructed by a CoNekT framework⁶⁴, which was also used to
929 establish the EVOREPRO database available at www.evorepro.plant.tools. For each species, all the genes

930 that were co-expressed in each male cluster were analysed to test whether the number of connections
931 observed is similar to the expected number. For this, we divided the number of observed connections
932 between the genes of two clusters (eg. cluster 1 and cluster 2) by the expected value (product of the
933 number of genes in cluster 1 x number of genes in cluster 2). These values were used to perform a
934 pearson correlation analysis and the results were presented in heatmaps. The networks present in the male
935 clusters were visualized using Cytoscape v3.8.0¹⁰⁷. The network files are available from
936 www.evorepro.plant.tools/species/.

937 **Data availability**

938 The fastq files are available for Arabidopsis (E-MTAB-9456), Amborella (E-MTAB-9190), Marchantia
939 (E-MTAB-9457), Physcomitrella (E-MTAB-9466), maize (E-MTAB-9692) and tomato (E-MTAB-9725).

940

941 **References**

- 942 1. Brown, R. C. & Lemmon, B. E. Spores before sporophytes: hypothesizing the origin of
943 sporogenesis at the algal-plant transition. *New Phytol.* **190**, 875–881 (2011).
- 944 2. Wellman, C. H., Osterloff, P. L. & Mohiuddin, U. Fragments of the earliest land plants. *Nature*
945 **425**, 282–285 (2003).
- 946 3. Edwards, D., Morris, J. L., Richardson, J. B. & Kenrick, P. Cryptospores and cryptophytes reveal
947 hidden diversity in early land floras. *New Phytol.* **202**, 50–78 (2014).
- 948 4. Jill Harrison, C. Development and genetics in the evolution of land plant body plans. *Philos. Trans.*
949 *R. Soc. Lond. B. Biol. Sci* **372**, (2017).
- 950 5. Kenrick, P. & Crane, P. R. The origin and early evolution of plants on land. *Nature* **389**, 33–39
951 (1997).
- 952 6. Friend, P. F. & House, M. R. The Devonian period. *Geological Society, London, Special*
953 *Publications* **1**, 233–236 (1964).
- 954 7. Berner, R. A. GEOCARBSULF: A combined model for Phanerozoic atmospheric O₂ and CO₂.

- 955 *Geochim. Cosmochim. Acta* **70**, 5653–5664 (2006).
- 956 8. Beerling, D. J., Osborne, C. P. & Chaloner, W. G. Evolution of leaf-form in land plants linked to
957 atmospheric CO₂ decline in the Late Palaeozoic era. *Nature* **410**, 352–354 (2001).
- 958 9. Menand, B. *et al.* An ancient mechanism controls the development of cells with a rooting function
959 in land plants. *Science* **316**, 1477–1480 (2007).
- 960 10. Hater, F., Nakel, T. & Groß-Hardt, R. Reproductive multitasking: the female gametophyte. *Annu.*
961 *Rev. Plant Biol.* **71**, 517–546 (2020).
- 962 11. Hackenberg, D. & Twell, D. The evolution and patterning of male gametophyte development.
963 *Curr. Top. Dev. Biol.* **131**, 257–298 (2019).
- 964 12. Johnson, M. A., Harper, J. F. & Palanivelu, R. A Fruitful Journey: Pollen Tube Navigation from
965 Germination to Fertilization. *Annu. Rev. Plant Biol.* **70**, 809–837 (2019).
- 966 13. Zhou, L.-Z. & Dresselhaus, T. Friend or foe: Signaling mechanisms during double fertilization in
967 flowering seed plants. *Curr. Top. Dev. Biol.* **131**, 453–496 (2019).
- 968 14. Dresselhaus, T., Sprunck, S. & Wessel, G. M. Fertilization mechanisms in flowering plants. *Curr.*
969 *Biol.* **26**, R125-39 (2016).
- 970 15. Sprunck, S. Twice the fun, double the trouble: gamete interactions in flowering plants. *Curr. Opin.*
971 *Plant Biol.* **53**, 106–116 (2020).
- 972 16. Borg, M. *et al.* The R2R3 MYB transcription factor DUO1 activates a male germline-specific
973 regulon essential for sperm cell differentiation in Arabidopsis. *Plant Cell* **23**, 534–549 (2011).
- 974 17. Favery, B. *et al.* KOJAK encodes a cellulose synthase-like protein required for root hair cell
975 morphogenesis in Arabidopsis. *Genes Dev.* **15**, 79–89 (2001).
- 976 18. Denninger, P. *et al.* Male-female communication triggers calcium signatures during fertilization in
977 Arabidopsis. *Nat. Commun.* **5**, 4645 (2014).
- 978 19. Leydon, A. R. *et al.* Pollen Tube Discharge Completes the Process of Synergid Degeneration That
979 Is Initiated by Pollen Tube-Synergid Interaction in Arabidopsis. *Plant Physiol.* **169**, 485–496
980 (2015).

- 981 20. Erbasol Serbes, I., Palovaara, J. & Groß-Hardt, R. Development and function of the flowering plant
982 female gametophyte. *Curr. Top. Dev. Biol.* **131**, 401–434 (2019).
- 983 21. Borges, F. *et al.* FACS-based purification of Arabidopsis microspores, sperm cells and vegetative
984 nuclei. *Plant Methods* **8**, 44 (2012).
- 985 22. Borg, M. *et al.* An EAR-Dependent Regulatory Module Promotes Male Germ Cell Division and
986 Sperm Fertility in Arabidopsis. *Plant Cell* **26**, 2098–2113 (2014).
- 987 23. Sprunck, S. *et al.* Egg cell-secreted EC1 triggers sperm cell activation during double fertilization.
988 *Science* **338**, 1093–1097 (2012).
- 989 24. Cyprys, P., Lindemeier, M. & Sprunck, S. Gamete fusion is facilitated by two sperm cell-expressed
990 DUF679 membrane proteins. *Nat. Plants* **5**, 253–257 (2019).
- 991 25. Rhee, S. Y. & Mutwil, M. Towards revealing the functions of all genes in plants. *Trends Plant Sci.*
992 **19**, 212–221 (2014).
- 993 26. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq
994 quantification. *Nat. Biotechnol.* **34**, 525–527 (2016).
- 995 27. Honys, D. & Twell, D. Comparative analysis of the Arabidopsis pollen transcriptome. *Plant*
996 *Physiol.* **132**, 640–652 (2003).
- 997 28. Pina, C., Pinto, F., Feijó, J. A. & Becker, J. D. Gene family analysis of the Arabidopsis pollen
998 transcriptome reveals biological implications for cell growth, division control, and gene expression
999 regulation. *Plant Physiol.* **138**, 744–756 (2005).
- 1000 29. Steffen, J. G., Kang, I.-H., Macfarlane, J. & Drews, G. N. Identification of genes expressed in the
1001 Arabidopsis female gametophyte. *Plant J.* **51**, 281–292 (2007).
- 1002 30. Wuest, S. E. *et al.* Arabidopsis female gametophyte gene expression map reveals similarities
1003 between plant and animal gametes. *Curr. Biol.* **20**, 506–512 (2010).
- 1004 31. Bowman, J. L. The YABBY gene family and abaxial cell fate. *Curr. Opin. Plant Biol.* **3**, 17–22
1005 (2000).
- 1006 32. Kim, J. H. & Lee, B. H. GROWTH-REGULATING FACTOR4 of Arabidopsis thaliana is required

- 1007 for development of leaves, cotyledons, and shoot apical meristem. *J. Plant Biol.* **49**, 463–468
1008 (2006).
- 1009 33. Lee, T. G. *et al.* A Myb transcription factor (TaMyb1) from wheat roots is expressed during
1010 hypoxia: roles in response to the oxygen concentration in root environment and abiotic stresses.
1011 *Physiol. Plant.* **129**, 375–385 (2006).
- 1012 34. Chen, D., Chai, S., McIntyre, C. L. & Xue, G.-P. Overexpression of a predominantly root-
1013 expressed NAC transcription factor in wheat roots enhances root length, biomass and drought
1014 tolerance. *Plant Cell Rep.* **37**, 225–237 (2018).
- 1015 35. Ding, Z. J. *et al.* Transcription factor WRKY46 modulates the development of Arabidopsis lateral
1016 roots in osmotic/salt stress conditions via regulation of ABA signaling and auxin homeostasis.
1017 *Plant J.* **84**, 56–69 (2015).
- 1018 36. Long, T. A. *et al.* The bHLH transcription factor POPEYE regulates response to iron deficiency in
1019 Arabidopsis roots. *Plant Cell* **22**, 2219–2236 (2010).
- 1020 37. Ding, W. *et al.* A transcription factor with a bHLH domain regulates root hair development in rice.
1021 *Cell Res.* **19**, 1309–1311 (2009).
- 1022 38. Betrán, E., Thornton, K. & Long, M. Retroposed new genes out of the X in Drosophila. *Genome*
1023 *Res.* **12**, 1854–1859 (2002).
- 1024 39. Begun, D. J., Lindfors, H. A., Kern, A. D. & Jones, C. D. Evidence for de novo evolution of testis-
1025 expressed genes in the Drosophila yakuba/Drosophila erecta clade. *Genetics* **176**, 1131–1137
1026 (2007).
- 1027 40. Dubruille, R., Marais, G. A. B. & Loppin, B. Repeated evolution of testis-specific new genes: the
1028 case of telomere-capping genes in Drosophila. *Int. J. Evol. Biol.* **2012**, 708980 (2012).
- 1029 41. Gossmann, T. I., Saleh, D., Schmid, M. W., Spence, M. A. & Schmid, K. J. Transcriptomes of
1030 Plant Gametophytes Have a Higher Proportion of Rapidly Evolving and Young Genes than
1031 Sporophytes. *Mol. Biol. Evol.* **33**, 1669–1678 (2016).
- 1032 42. Cui, X. *et al.* Young Genes out of the Male: An Insight from Evolutionary Age Analysis of the

- 1033 Pollen Transcriptome. *Mol. Plant* **8**, 935–945 (2015).
- 1034 43. Doyle, J. A. in *Annual Plant Reviews* (eds. Roberts, J. A., Evan, D., McManus, M. T. & Rose, J. K.
1035 C.) 1–50 (John Wiley & Sons, Ltd, 2018). doi:10.1002/9781119312994.apr0486
- 1036 44. Beerling, D. J. Leaf evolution: gases, genes and geochemistry. *Ann. Bot.* **96**, 345–352 (2005).
- 1037 45. Pires, N. D. & Dolan, L. Morphological evolution in land plants: new designs with old genes.
1038 *Philos. Trans. R. Soc. Lond. B. Biol. Sci* **367**, 508–518 (2012).
- 1039 46. Cardona, T. Thinking twice about the evolution of photosynthesis. *Open Biol.* **9**, 180246 (2019).
- 1040 47. Harrison, C. J. & Morris, J. L. The origin and early evolution of vascular plant shoots and leaves.
1041 *Philos. Trans. R. Soc. Lond. B. Biol. Sci* **373**, (2018).
- 1042 48. Hetherington, A. J. & Dolan, L. Stepwise and independent origins of roots among land plants.
1043 *Nature* **561**, 235–238 (2018).
- 1044 49. Specht, C. D. & Bartlett, M. E. Flower Evolution: The Origin and Subsequent Diversification of
1045 the Angiosperm Flower. *Annu. Rev. Ecol. Evol. Syst.* **40**, 217–243 (2009).
- 1046 50. Pires, N. D. *et al.* Recruitment and remodeling of an ancient gene regulatory network during land
1047 plant evolution. *Proc Natl Acad Sci USA* **110**, 9571–9576 (2013).
- 1048 51. Huang, L. & Schiefelbein, J. Conserved Gene Expression Programs in Developing Roots from
1049 Diverse Plants. *Plant Cell* **27**, 2119–2132 (2015).
- 1050 52. He, C., Si, C., Teixeira da Silva, J. A., Li, M. & Duan, J. Genome-wide identification and
1051 classification of MIKC-type MADS-box genes in Streptophyte lineages and expression analyses to
1052 reveal their role in seed germination of orchid. *BMC Plant Biol.* **19**, 223 (2019).
- 1053 53. Tanabe, Y. *et al.* Characterization of MADS-box genes in charophycean green algae and its
1054 implication for the evolution of MADS-box genes. *Proc Natl Acad Sci USA* **102**, 2436–2441
1055 (2005).
- 1056 54. Brodribb, T. J., Carriquí, M., Delzon, S., McAdam, S. A. M. & Holbrook, N. M. Advanced
1057 vascular function discovered in a widespread moss. *Nat. Plants* **6**, 273–279 (2020).
- 1058 55. Ruprecht, C. *et al.* Phylogenomic analysis of gene co-expression networks reveals the evolution of

- 1059 functional modules. *Plant J.* **90**, 447–465 (2017).
- 1060 56. Rao, X. & Dixon, R. A. Co-expression networks for plant biology: why and how. *Acta Biochim*
1061 *Biophys Sin (Shanghai)* **51**, 981–988 (2019).
- 1062 57. Mutwil, M. Computational approaches to unravel the pathways and evolution of specialized
1063 metabolism. *Curr. Opin. Plant Biol.* **55**, 38–46 (2020).
- 1064 58. Borges, F. *et al.* Comparative transcriptomics of Arabidopsis sperm cells. *Plant Physiol.* **148**,
1065 1168–1181 (2008).
- 1066 59. Liu, L. *et al.* Transcriptomics analyses reveal the molecular roadmap and long non-coding RNA
1067 landscape of sperm cell lineage development. *Plant J.* **96**, 421–437 (2018).
- 1068 60. Anderson, S. N. *et al.* Transcriptomes of isolated *Oryza sativa* gametes characterized by deep
1069 sequencing: evidence for distinct sex-dependent chromatin and epigenetic states before
1070 fertilization. *Plant J.* **76**, 729–741 (2013).
- 1071 61. Borg, M. *et al.* Targeted reprogramming of H3K27me3 resets epigenetic memory in plant paternal
1072 chromatin. *Nat. Cell Biol.* **22**, 621–629 (2020).
- 1073 62. Becker, J. D., Takeda, S., Borges, F., Dolan, L. & Feijó, J. A. Transcriptional profiling of
1074 Arabidopsis root hairs and pollen defines an apical cell growth signature. *BMC Plant Biol.* **14**, 197
1075 (2014).
- 1076 63. von Besser, K., Frank, A. C., Johnson, M. A. & Preuss, D. Arabidopsis HAP2 (GCS1) is a sperm-
1077 specific gene required for pollen tube guidance and fertilization. *Development* **133**, 4761–4769
1078 (2006).
- 1079 64. Proost, S. & Mutwil, M. CoNekT: an open-source framework for comparative genomic and
1080 transcriptomic network analyses. *Nucleic Acids Res.* **46**, W133–W140 (2018).
- 1081 65. Boisson-Dernier, A. *et al.* Disruption of the pollen-expressed FERONIA homologs ANXUR1 and
1082 ANXUR2 triggers pollen tube discharge. *Development* **136**, 3279–3288 (2009).
- 1083 66. Zhu, L. *et al.* The Arabidopsis CrRLK1L protein kinases BUPS1 and BUPS2 are required for
1084 normal growth of pollen tubes in the pistil. *Plant J.* **95**, 474–486 (2018).

- 1085 67. Alves-Ferreira, M. *et al.* Global expression profiling applied to the analysis of Arabidopsis stamen
1086 development. *Plant Physiol.* **145**, 747–762 (2007).
- 1087 68. Gupta, R., Ting, J. T. L., Sokolov, L. N., Johnson, S. A. & Luan, S. A tumor suppressor homolog,
1088 AtPTEN1, is essential for pollen development in Arabidopsis. *Plant Cell* **14**, 2495–2507 (2002).
- 1089 69. Zhou, Z. *et al.* Arabidopsis RIC1 severs actin filaments at the apex to regulate pollen tube growth.
1090 *Plant Cell* **27**, 1140–1161 (2015).
- 1091 70. Liang, Y. *et al.* MYB97, MYB101 and MYB120 function as male factors that control pollen tube-
1092 synergid interaction in Arabidopsis thaliana fertilization. *PLoS Genet.* **9**, e1003933 (2013).
- 1093 71. Domazet-Lošo, T. & Tautz, D. A phylogenetically based transcriptome age index mirrors
1094 ontogenetic divergence patterns. *Nature* **468**, 815–818 (2010).
- 1095 72. Ashton, N. W. & Cove, D. J. The isolation and preliminary characterisation of auxotrophic and
1096 analogue resistant mutants of the moss, *Physcomitrella patens*. *Molec. Gen. Genet.* **154**, 87–95
1097 (1977).
- 1098 73. Hohe, A., Rensing, S. A., Mildner, M., Lang, D. & Reski, R. Day Length and Temperature
1099 Strongly Influence Sexual Reproduction and Expression of a Novel MADS-Box Gene in the Moss
1100 *Physcomitrella patens*. *Plant Biol (Stuttg)* **4**, 595–602 (2002).
- 1101 74. Misra, C. S. *et al.* Transcriptomics of Arabidopsis sperm cells at single-cell resolution. *Plant*
1102 *Reprod.* **32**, 29–38 (2019).
- 1103 75. Picelli, S. *et al.* Full-length RNA-seq from single cells using Smart-seq2. *Nat. Protoc.* **9**, 171–181
1104 (2014).
- 1105 76. Dupláková, N., Dobrev, P. I., Reňák, D. & Honys, D. Rapid separation of Arabidopsis male
1106 gametophyte developmental stages using a Percoll gradient. *Nat. Protoc.* **11**, 1817–1832 (2016).
- 1107 77. Flores-Tornero, M. *et al.* Transcriptomic and Proteomic Insights into Amborella trichopoda Male
1108 Gametophyte Functions. *Plant Physiol.* (2020). doi:10.1104/pp.20.00837
- 1109 78. Lu, Y., Wei, L. & Wang, T. Methods to isolate a large amount of generative cells, sperm cells and
1110 vegetative nuclei from tomato pollen for “omics” analysis. *Front. Plant Sci.* **6**, 391 (2015).

- 1111 79. Flores-Tornero, M. *et al.* Transcriptomics of manually isolated *Amborella trichopoda* egg
1112 apparatus cells. *Plant Reprod.* **32**, 15–27 (2019).
- 1113 80. Honys, D. & Twell, D. Transcriptome analysis of haploid male gametophyte development in
1114 *Arabidopsis*. *Genome Biol.* **5**, R85 (2004).
- 1115 81. Boavida, L. C. & McCormick, S. Temperature as a determinant factor for increased and
1116 reproducible in vitro pollen germination in *Arabidopsis thaliana*. *Plant J.* **52**, 570–582 (2007).
- 1117 82. Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat.*
1118 *Methods* **10**, 1096–1098 (2013).
- 1119 83. Baym, M. *et al.* Inexpensive multiplexed library preparation for megabase-sized genomes. *PLoS*
1120 *ONE* **10**, e0128036 (2015).
- 1121 84. Enghart, M., Šoljić, L. & Sprunck, S. Manual Isolation of Living Cells from the *Arabidopsis*
1122 *thaliana* Female Gametophyte by Micromanipulation. *Methods Mol. Biol.* **1669**, 221–234 (2017).
- 1123 85. Sprunck, S. *et al.* Elucidating small RNA pathways in *Arabidopsis thaliana* egg cells. *BioRxiv*
1124 (2019). doi:10.1101/525956
- 1125 86. Brukhin, V., Hernould, M., Gonzalez, N., Chevalier, C. & Mouras, A. Flower development
1126 schedule in tomato *Lycopersicon esculentum* cv. sweet cherry. *Sex. Plant Reprod.* **15**, 311–320
1127 (2003).
- 1128 87. Covey, P. A. *et al.* A pollen-specific RALF from tomato that regulates pollen tube elongation.
1129 *Plant Physiol.* **153**, 703–715 (2010).
- 1130 88. Chen, J. *et al.* Zygotic Genome Activation Occurs Shortly after Fertilization in Maize. *Plant Cell*
1131 **29**, 2106–2125 (2017).
- 1132 89. Schreiber, D. N., Bantin, J. & Dresselhaus, T. The MADS box transcription factor ZmMADS2 is
1133 required for anther and pollen maturation in maize and accumulates in apoptotic bodies during
1134 anther dehiscence. *Plant Physiol.* **134**, 1069–1079 (2004).
- 1135 90. Harrison, P. W. *et al.* The european nucleotide archive in 2018. *Nucleic Acids Res.* **47**, D84–D88
1136 (2019).

- 1137 91. Mutwil, M. *et al.* Assembly of an interactive correlation network for the Arabidopsis genome using
1138 a novel heuristic clustering algorithm. *Plant Physiol.* **152**, 29–43 (2010).
- 1139 92. Ferrari, C. *et al.* Expression Atlas of *Selaginella moellendorffii* Provides Insights into the Evolution
1140 of Vasculature, Secondary Metabolism, and Roots. *Plant Cell* **32**, 853–870 (2020).
- 1141 93. Wagner, G. P., Kin, K. & Lynch, V. J. A model based criterion for gene expression calls using
1142 RNA-seq data. *Theory Biosci.* **132**, 159–164 (2013).
- 1143 94. Xiao, S.-J., Zhang, C., Zou, Q. & Ji, Z.-L. TiSGeD: a database for tissue-specific genes.
1144 *Bioinformatics* **26**, 1273–1275 (2010).
- 1145 95. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful
1146 approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Methodological)*
1147 **57**, 289–300 (1995).
- 1148 96. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative
1149 genomics. *Genome Biol.* **20**, 238 (2019).
- 1150 97. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat.*
1151 *Methods* **12**, 59–60 (2015).
- 1152 98. One Thousand Plant Transcriptomes Initiative. One thousand plant transcriptomes and
1153 the phylogenomics of green plants. *Nature* **574**, 679–685 (2019).
- 1154 99. Domazet-Loso, T., Brajković, J. & Tautz, D. A phylostratigraphy approach to uncover the genomic
1155 history of major adaptations in metazoan lineages. *Trends Genet.* **23**, 533–539 (2007).
- 1156 100. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: reconstruction, analysis, and visualization of
1157 phylogenomic data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).
- 1158 101. Schwacke, R. *et al.* MapMan4: A Refined Protein Classification and Annotation Framework
1159 Applicable to Multi-Omics Data Analysis. *Mol. Plant* **12**, 879–892 (2019).
- 1160 102. Zheng, Y. *et al.* iTAK: A Program for Genome-wide Prediction and Classification of Plant
1161 Transcription Factors, Transcriptional Regulators, and Protein Kinases. *Mol. Plant* **9**, 1667–1670
1162 (2016).

- 1163 103. Tian, F., Yang, D.-C., Meng, Y.-Q., Jin, J. & Gao, G. PlantRegMap: charting functional regulatory
1164 maps in plants. *Nucleic Acids Res.* **48**, D1104–D1113 (2020).
- 1165 104. Waskom, M. *et al.* Seaborn: V0.5.0 (November 2014). *Zenodo* (2014). doi:10.5281/zenodo.12710
- 1166 105. Ballester, A.-R. *et al.* Genome, Transcriptome, and Functional Analyses of *Penicillium expansum*
1167 Provide New Insights Into Secondary Metabolism and Pathogenicity. *Mol. Plant Microbe Interact.*
1168 **28**, 232–248 (2015).
- 1169 106. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine Learning*
1170 *Research* (2011).
- 1171 107. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular
1172 interaction networks. *Genome Res.* **13**, 2498–2504 (2003).

1173 **Acknowledgments**

1174 I.J is supported by Singaporean Ministry of Education grant MOE2018-T2-2-053, while M.M is
1175 supported by NTU Start-Up Grant. ERA-CAPS EVO-REPRO I2163 to F.B.; ERA-CAPS-0001-2014 to
1176 J.D.B; ERA-CAPS EVO-REPRO DR 334/12-1 to S.S. and T.D. DH was supported by ERA-CAPS UK
1177 Biotechnology and Biological Research Council Grant BB/N005090 awarded to DT; M.B. was supported
1178 through the FWF Lise Meitner fellowship M1818. The Vienna BioCenter Core Facilities GmbH (VBCF)
1179 Plant Sciences Facility acknowledges funding from the Austrian Federal Ministry of Education, Science
1180 and Research and the City of Vienna. L.S was supported by CFS grant 17-23183S. C.M. and D.Ho. were
1181 supported by Czech Ministry of Education, Youth and Sport (LTC18034 and LTAIN19030) through the
1182 European Regional Development Fund-Project “Centre for Experimental Plant Biology”: No.
1183 CZ.02.1.01/0.0/0.0/16_019/0000738. The Genomics Unit of Instituto Gulbenkian de Ciência was partially
1184 supported by ONEIDA Project (LISBOA-01-0145-FEDER-016417) co-funded by FEEI - “Fundos
1185 Europeus Estruturais e de Investimento” from “Programa Operacional Regional Lisboa 2020” and by
1186 national funds from FCT - “Fundação para a Ciência e a Tecnologia”. C.S.M acknowledges a doctoral
1187 fellowship from FCT (PD/BD/114362/2016) under the Plants for Life PhD Program. J.D.B received
1188 salary support from FCT through an “Investigador FCT” position. MJ and JG were supported by a US
1189 National Science Foundation grant (IOS-1540019).

1190 Help with sample generation: Lenka Závěská Drábková and David Reňák. *Marchantia* growth was
1191 performed by the Plant Sciences Facility at Vienna BioCenter Core Facilities GmbH (VBCF), member of
1192 the Vienna BioCenter (VBC), Austria. Maximilian Weigend, Cornelia Löhne and Bernhard Reinken
1193 (Botanical Garden of the University of Bonn, Germany) are acknowledged for providing *Amborella*
1194 *trichopoda* plant material.

1195 We would like to thank Debbie Maizels (<http://www.scientificart.com>) for the illustrations on Fig.1 and
1196 Fig. 5.

1197

1198

1199 **Author Contributions**

1200 Conceived and designed the analysis: JDB, MM

1201 Collected the data: ACL, MFT, SGP, CSM, IJ, LS, CM, DHo, DH

1202 Contributed data or analysis tools: FB, MB, SS, TD, DT

1203 Performed the analysis: IJ, CF, SP, ACL, MM

1204 Wrote the paper: IJ, JDB, MM

1205

1206 **Competing interests**

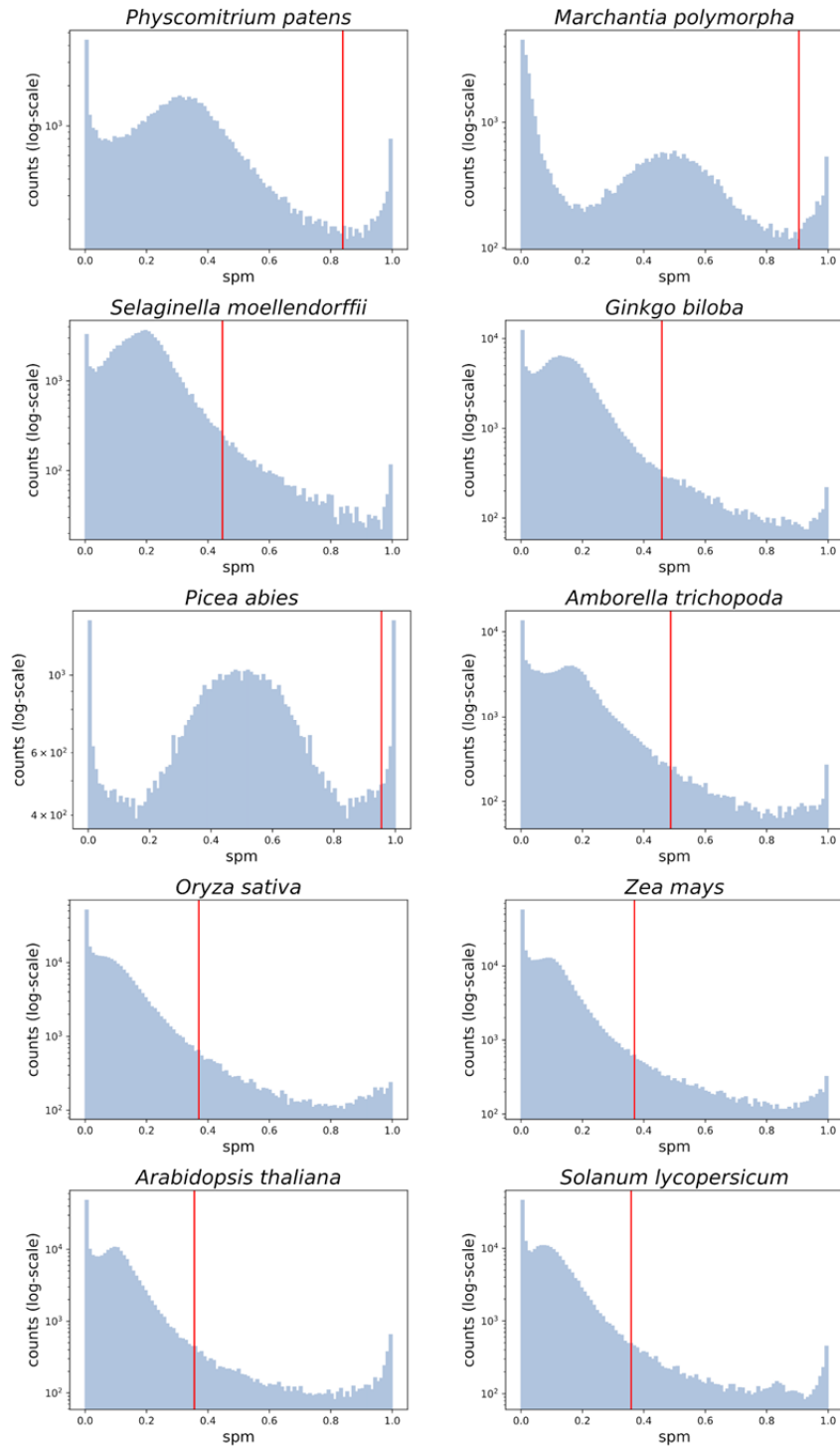
1207 The authors declare no competing interests.

1208

1209 **Supplementary information**

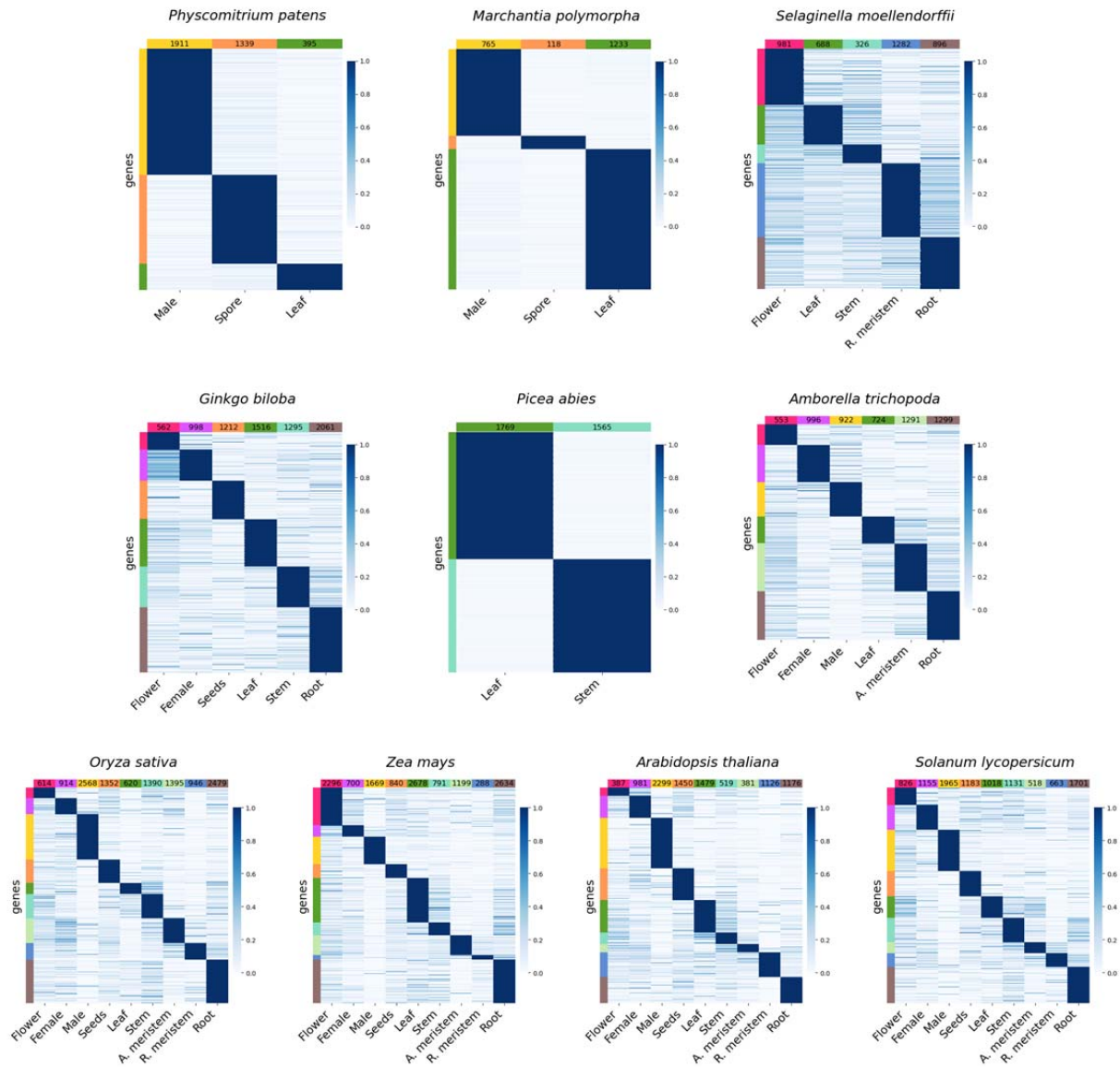
1210

1211



1212

1213 **Supplementary Fig. 1: Distribution of SPM values in the ten species.** The x-axis indicates the
1214 specificity measure (SPM), while the y-axis indicates the log₁₀-transformed frequency of the SPM values
1215 observed for all genes across the samples. The vertical red line indicates the SPM value cutoff, below
1216 which 95% of values are found.



1217

1218 **Supplementary Fig. 2: Expression profiles of the genes that were deemed to be specifically**

1219 **expressed in one of the organs/tissues/cells (sample) of the ten species used in this study. Genes are in**

1220 rows, samples in columns, and the genes are sorted according to the expression profiles (e.g., flower,

1221 female). The numbers at the top of each column indicate the total number of specific genes in each

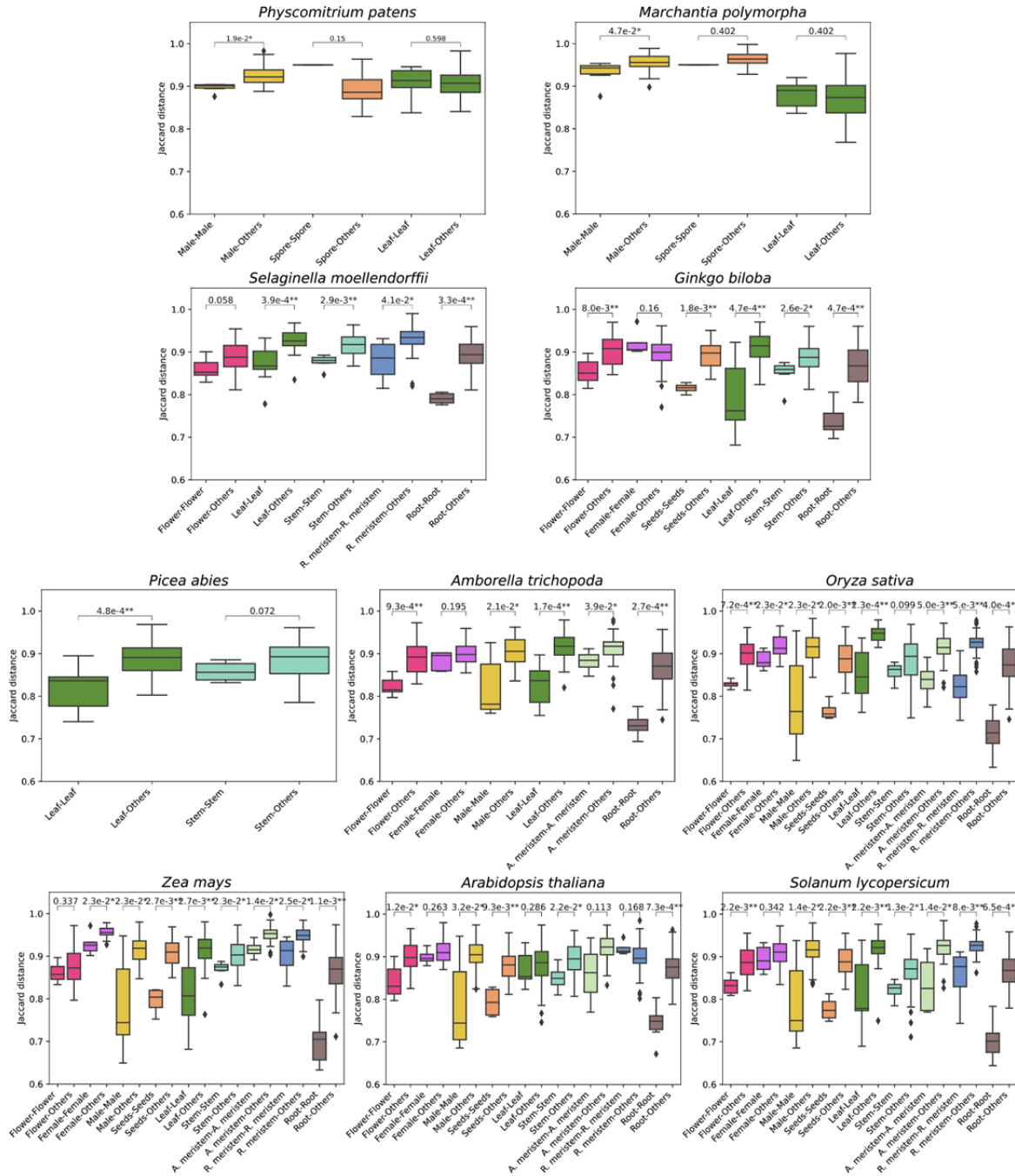
1222 sample. Gene expression is scaled to range from 0-1. Bars on the left of each heatmap show the sample-

1223 sample-specific genes and correspond to the samples on the bottom: pink - Flower, purple - Female, yellow -

1224 Male, orange - Seeds/Spore, dark-green - Leaf, medium-green - Stem, light-green - Apical meristem,

1225 blue - Root meristem, brown - Root.

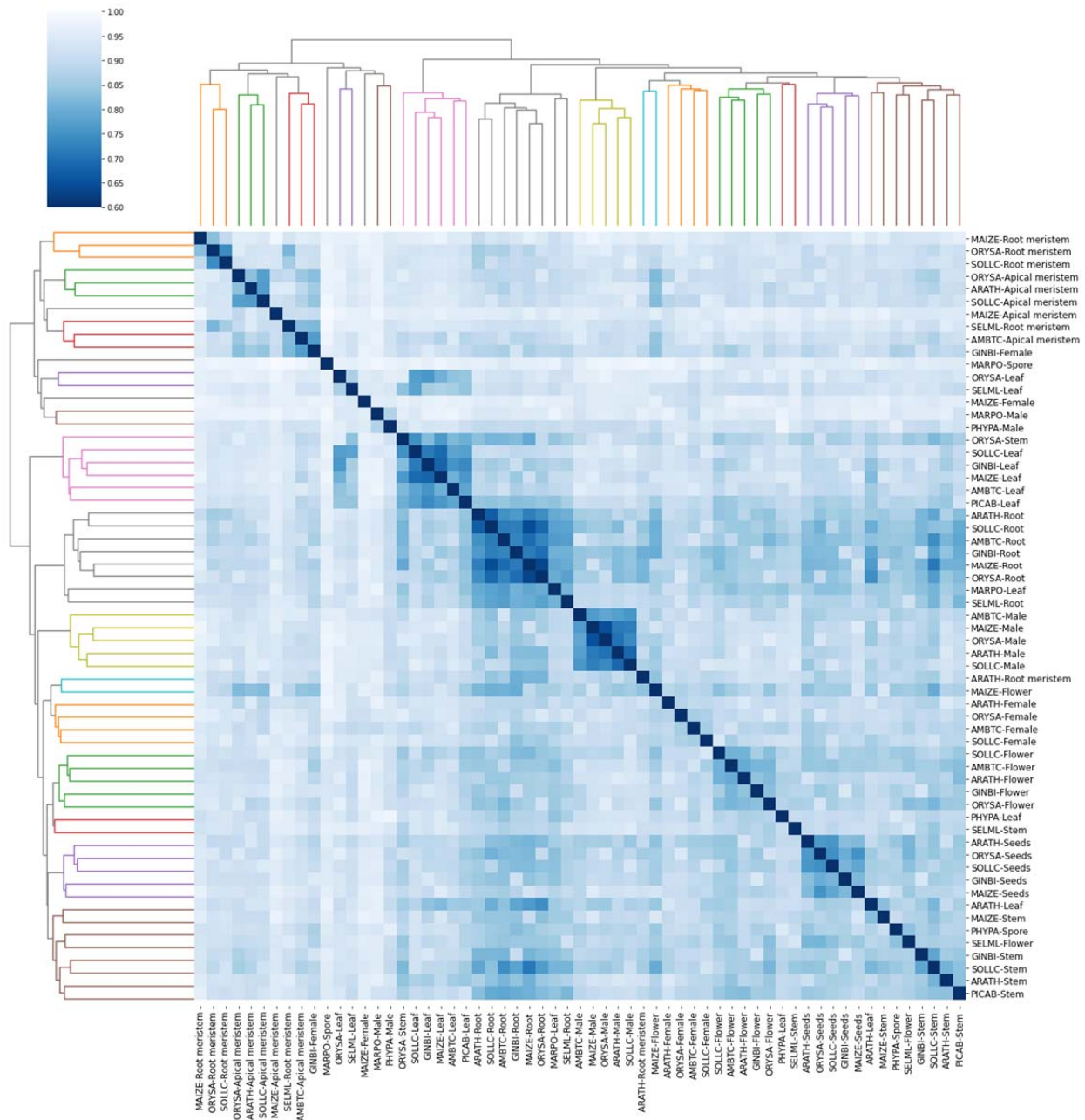
1226



1227

1228 **Supplementary Fig. 3:** Bar plot showing the Jaccard distances when comparing the same samples (i.e.,
 1229 male-male) and one sample versus the others (i.e., male-others) for the ten species included in this study.

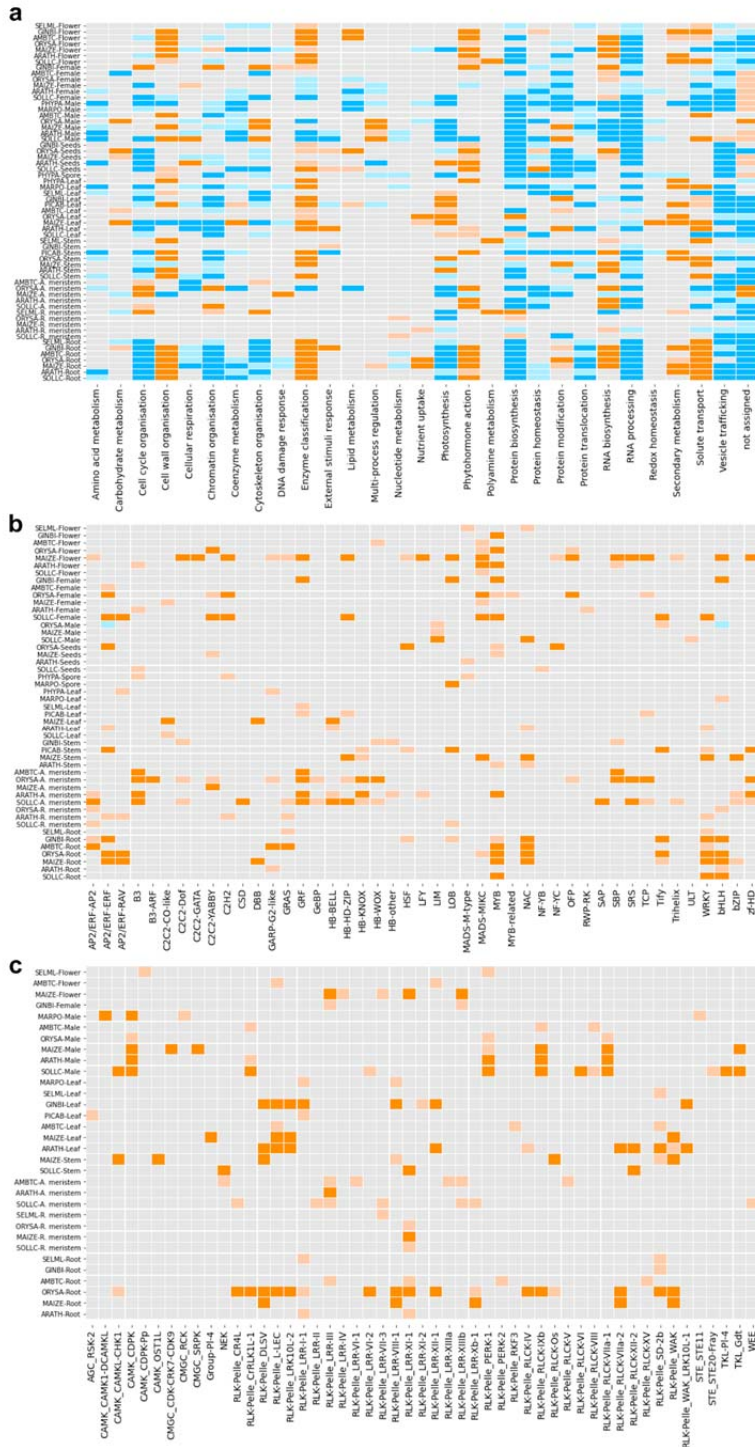
1230



1231

1232 **Supplementary Fig. 4:** Comparing transcriptome similarities of the samples of the ten species. We used
 1233 the Jaccard Index to calculate the similarity of transcriptomes of all samples in the dataset. The heatmap
 1234 shows which transcriptomes of samples across species are similar by hierarchical clustering (dark blue).
 1235 A lower value indicates a stronger similarity between two samples (white). For example, when comparing
 1236 *Arabidopsis* root to roots from other species, we observe more similar transcriptomes than *Arabidopsis*
 1237 root to non-root samples. The dendrograms on top and the left show the different clusters formed when
 1238 the distance is <1.3.

1239



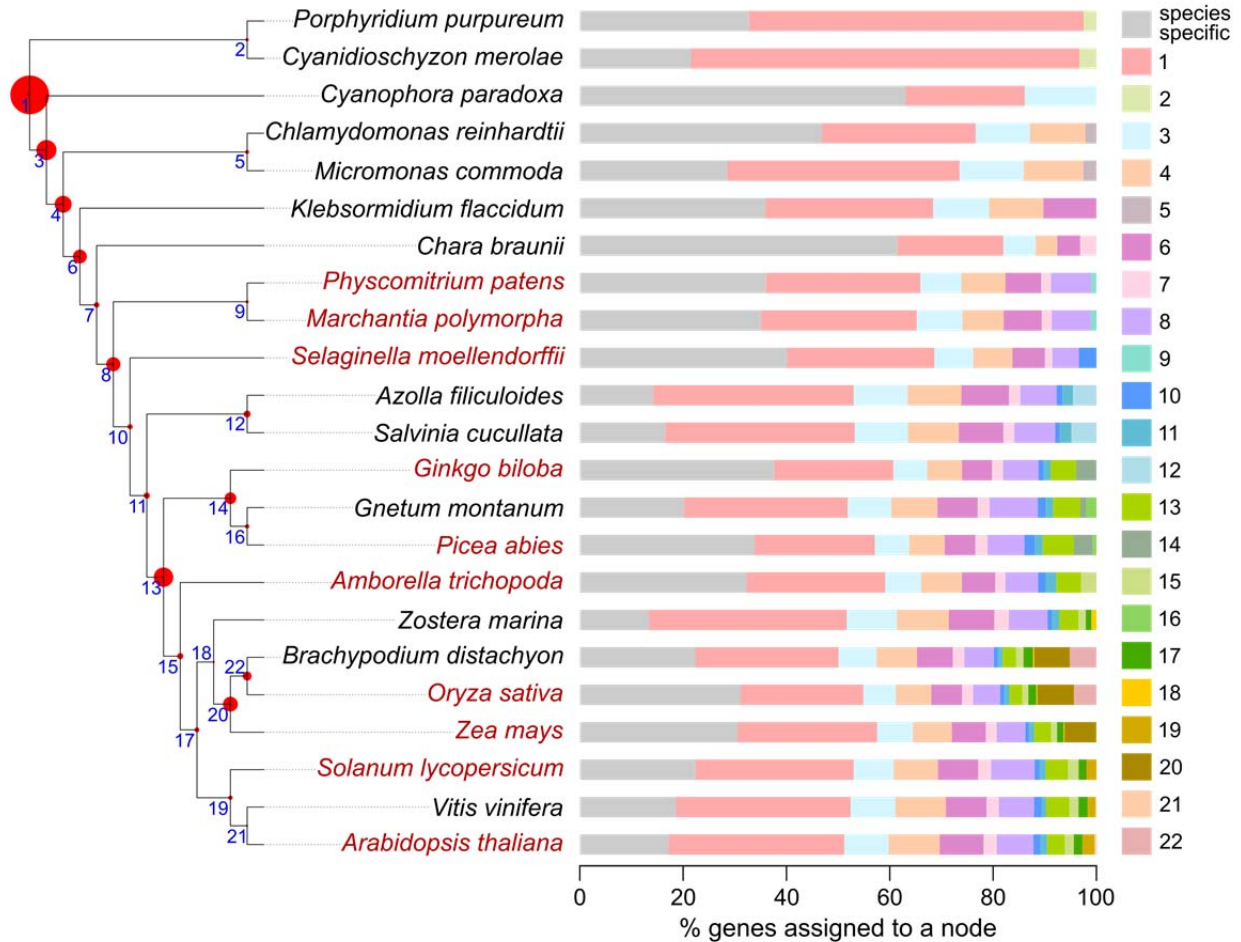
1240

1241 **Supplementary Fig. 5. Functional enrichment analysis of the sample-specific transcriptomes.**

1242 Samples are shown on the y-axis and functions in the x-axis for MapMan bins (a), transcription factors

1243 (b), and kinases (c). Orange and blue colors indicate enrichment and depletion, respectively. The intensity

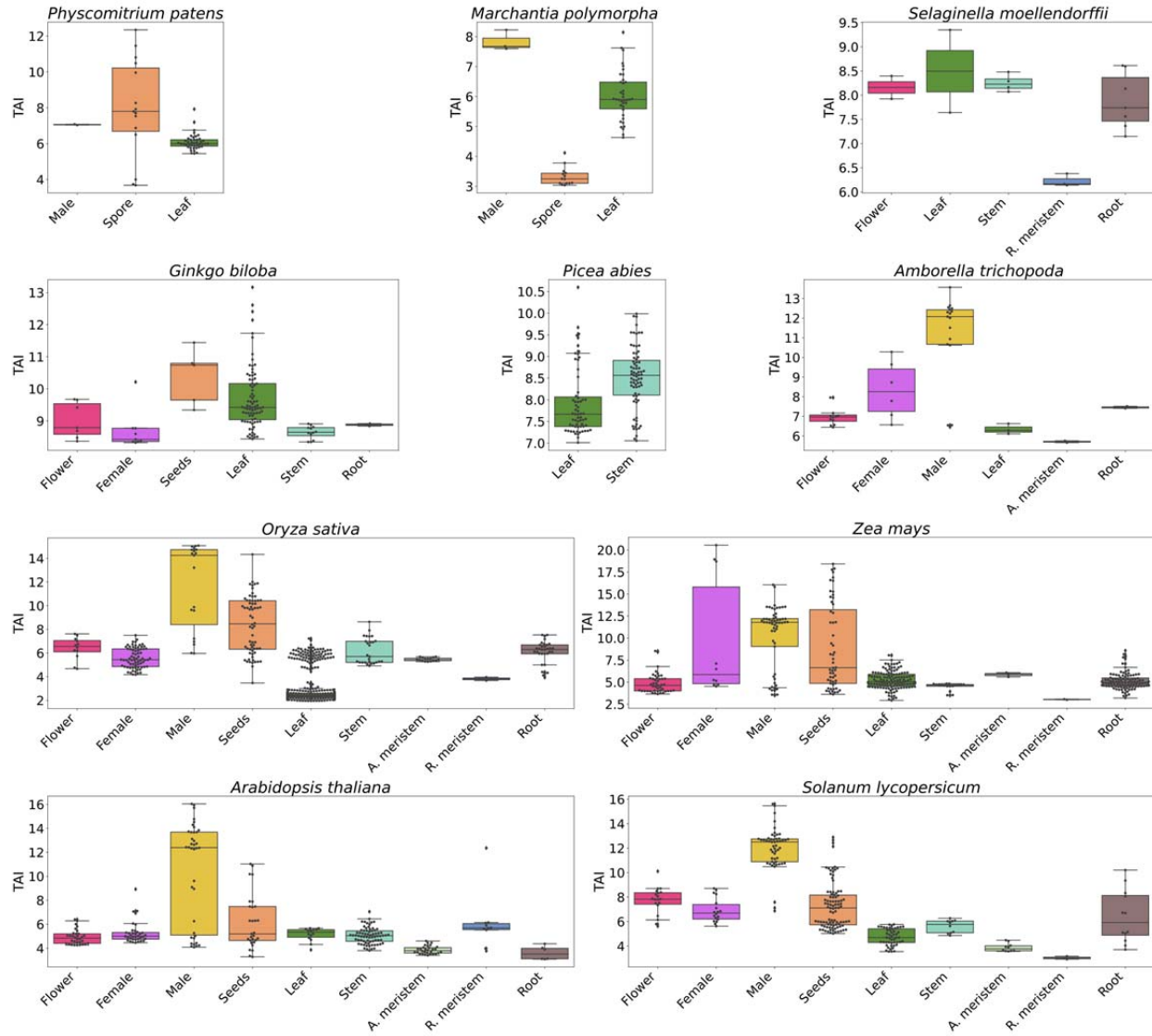
1244 of the color is in correlation with the p-value (dark orange/blue: $p < 0.01$, light orange/blue: $p < 0.05$).



1245

1246 **Supplementary Fig. 6: Cladogram of the 23 species included in the analysis.** The phylogenetic
 1247 relationship was based on One Thousand Plant Transcriptomes Initiative, 2019. Species in red are
 1248 associated with transcriptomic data in this study. Blue numbers in the nodes indicate the node number
 1249 (e.g., 1: NODE_1). The tree's red circles show the percentage of orthogroups found in each node (largest
 1250 and smallest amounts: Node_1 - 24% and NODE_21 - 0.1%). Bars on the right show the percentage of
 1251 genes per species that are present in each node. The nodes are shown in different colors, as indicated in
 1252 the right bar.

1253

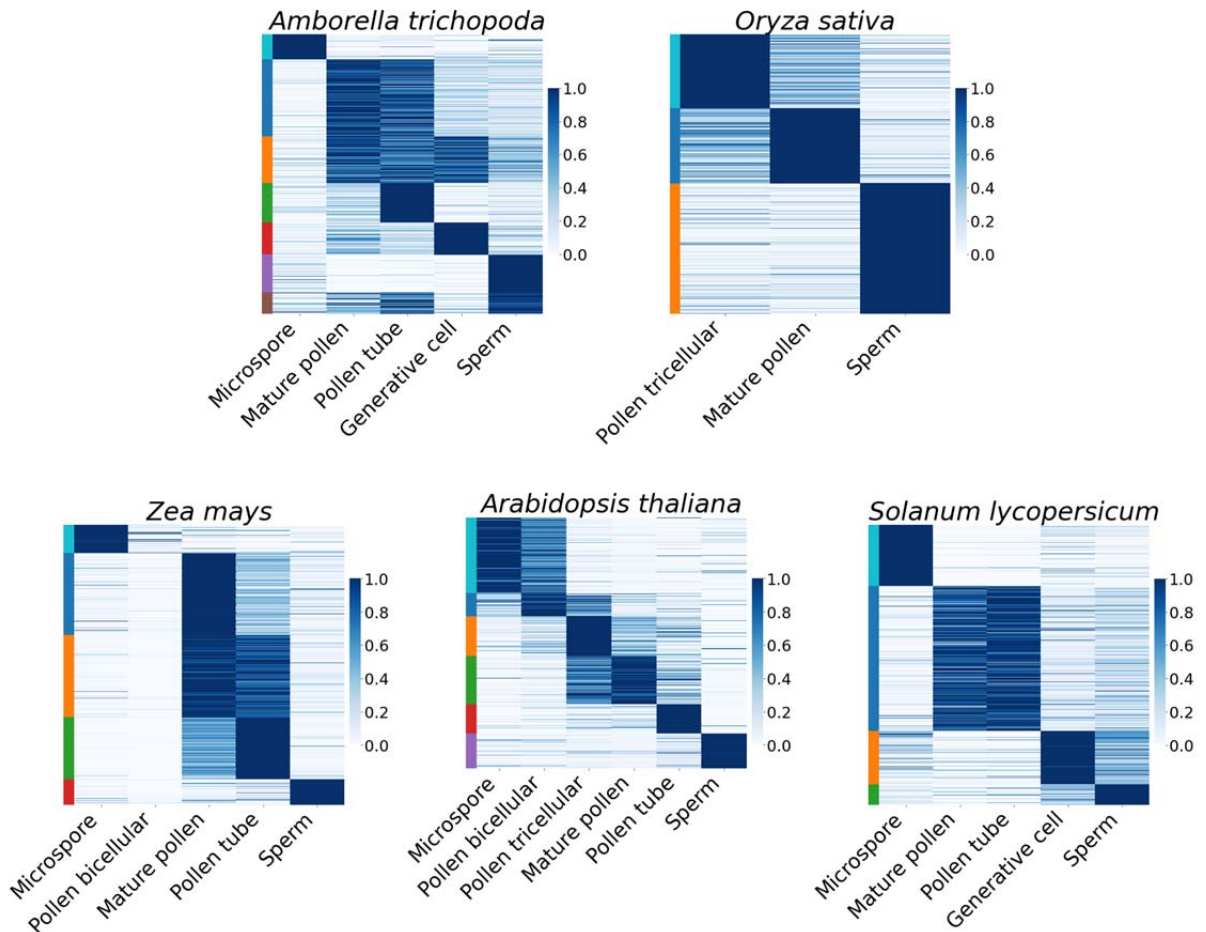


1254

1255 **Supplementary Fig. 7: Transcriptomic age index in the ten species.**

1256

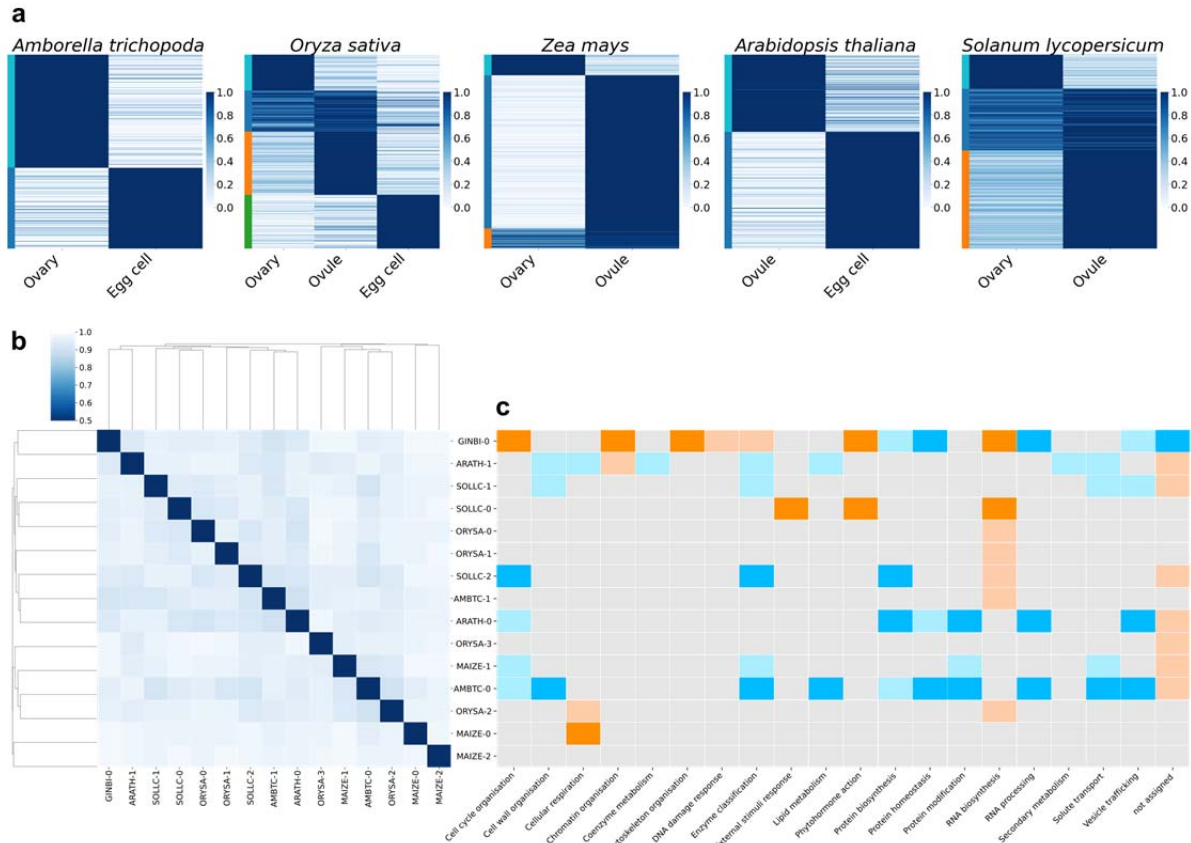
1257



1258

1259 **Supplementary Fig. 8: Expression of male developmental stages genes for five species.** Genes are in
1260 rows, developmental stages in columns. Gene expression is scaled to range from 0-1. Darker color
1261 corresponds to a stronger positive correlation. Bars in the left mark the different clusters.

1262



1263

1264 **Supplementary Fig. 9: Analysis of the expression profile in different development stages of female**

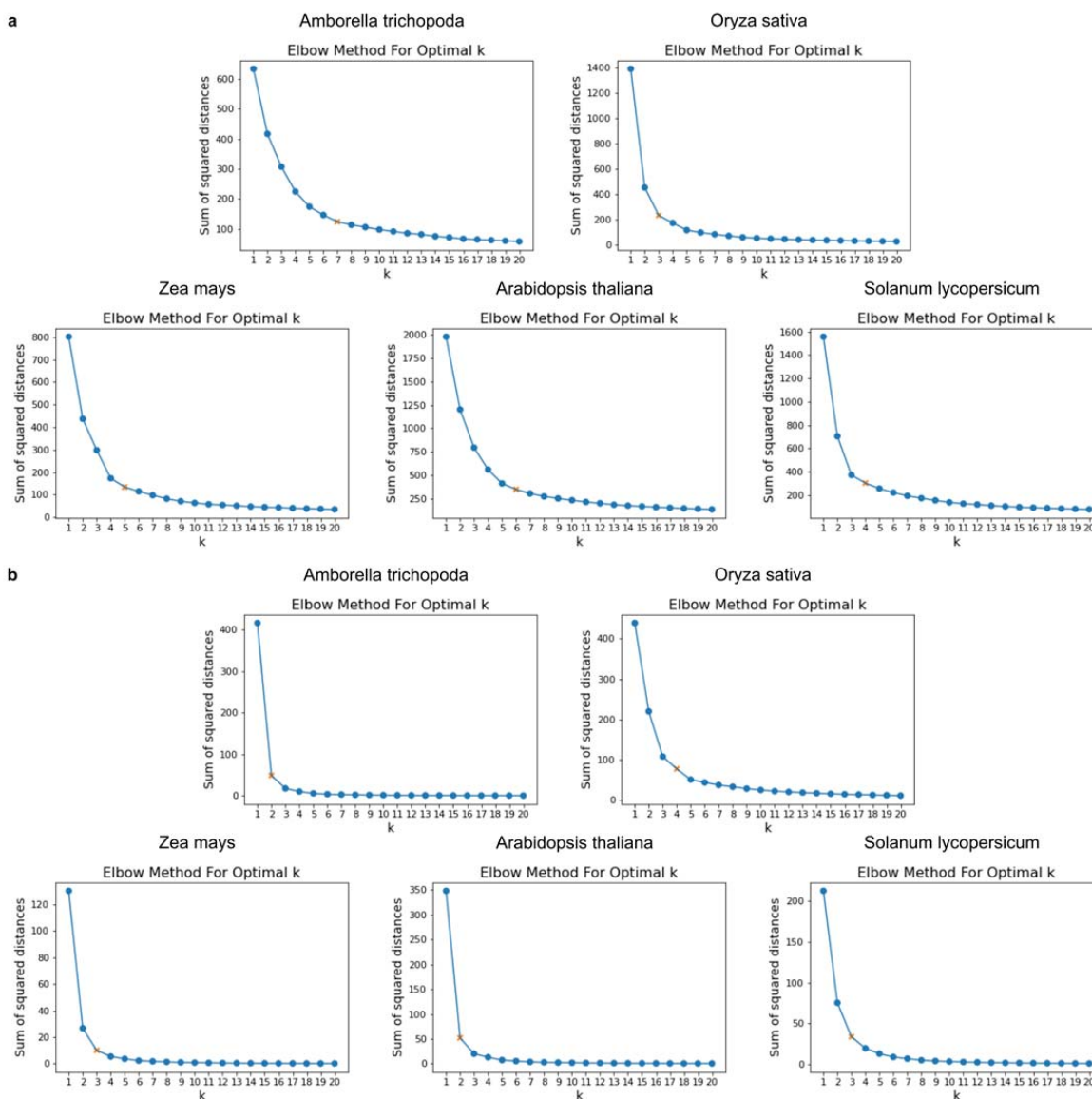
1265 **organs.** Heat map showing the normalized TMP of genes per each development stage for five species.

1266 Bars on the left indicate the clusters. **b**, Jaccard distance between the clusters. **c**, Heatmap showing

1267 enrichment and depletion of functions. Orange and blue indicate enrichment and depletion, respectively

1268 (light colors: $p < 0.05$, dark colors: $p < 0.01$).

1269



1270

1271 **Supplementary Fig. 10: Identifying k value with the elbow method.** The orange mark indicates the k
1272 value where the sum of squared distances was less than 80% of the highest value found at $k=1$. **a**, For the
1273 male samples, and **b**, for the female samples.

1274

1275

1276

1277

1278

1279 **Supplementary tables:**

1280 **Supplementary Table 1.** Samples included per each species. The columns in order show: mnemonic of
1281 the species, sample ID, original annotation, Organ name, Subsample name, source of the sample, number
1282 of fragments that could be pseudoaligned using, percentage of fragments that could be pseudoaligned, tag
1283 for the samples that pass/fail Kallisto stats, tag for the samples that pass/fail PCC filter.

1284 **Supplementary Table 2.** The number of expressed and organ-specific expressed genes in the ten species.

1285 **Supplementary Table 3.** Organ-specific genes in the ten land plants. Columns show species mnemonic,
1286 sample name, number of genes, gene names.

1287 **Supplementary Table 4.** Organ-specific transcription factors in the ten land plants. Mnemonics indicate
1288 the species. The columns indicate transcription factor families.

1289 **Supplementary Table 5.** Organ-specific kinases in the ten land plants. The species are indicated by
1290 mnemonics, while the organs are given after the species name. The different families of kinases are given
1291 in columns.

1292 **Supplementary Table 6.** List of orthogroups identified in the 23 species included. The columns show the
1293 orthogroup name, node in the species tree (Fig. 3a), expression profile, pass/fail filter of the expression
1294 profile, list of species (mnemonic). The following columns show the list of genes per species.

1295 **Supplementary Table 7.** Sample-specific gene enrichment for species and for node in the species tree
1296 (Fig. 3a). The columns show: species mnemonic, sample name, node of the species tree, p-value, tag
1297 (enrichment or depletion).

1298 **Supplementary Table 8.** Gain/loss of gene families. The columns show the sample, node, number of
1299 total gains, number of total losses, orthogroups gained, orthogroups lost.

1300 **Supplementary Table 9.** List of enriched functions in gained organ-specific and ubiquitous gene families
1301 per each node.

1302 **Supplementary Table 10.** List of male cluster-specific genes. The first column shows the mnemonic of
1303 the species. The second, the cluster number. The third to ninth column: the average TPM per each male
1304 sample included. The last column: the list of genes of the cluster.

1305 **Supplementary Table 11.** List of female cluster-specific genes. The first column shows the mnemonic of
1306 the species. The second, the cluster number. The third to fifth column: the average TPM per each male
1307 sample included. The sixth column: the list of genes of the cluster.

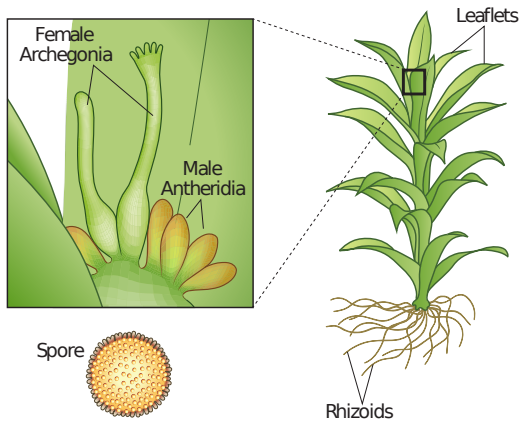
1308 **Supplementary Table 12.** Features of male cluster-specific genes. The columns show the mnemonic of
1309 the species, gene name, cluster number, if it is co-expressed (Yes/No), transcription factor, or kinase
1310 name if reported in the annotation.

1311 **Supplementary Table 13.** Annotation of the male cluster-specific genes of *A. thaliana*. The columns
1312 show: cluster name, gene, tag for transcription factor (TF) or kinase (KIN), name of the transcription
1313 factor or kinase, if it is co-expressed (Y/N), name of the sample that the known mutant affect, mutant, if
1314 the gene is involved in pollen, references.

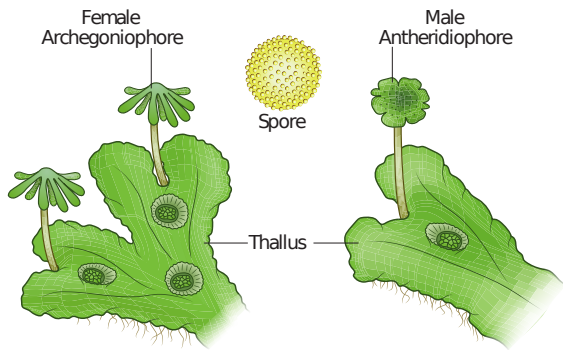
1315 **Supplementary Table 14.** List of species included in this study and the source of their proteomes and
1316 CDSs. Columns show: mnemonic, taxon identifier, species name, genome version, and source.

1317 **Supplementary Table 15.** Gene families (first column), Arabidopsis male-specific genes (second
1318 column) and Amborella male-specific genes (third column). Gene and family IDs are clickable and will
1319 redirect the user to a corresponding page. The fourth column indicates gene families found in common in
1320 Arabidopsis and Amborella (intersection), only in Arabidopsis (left) or only in Amborella (right).

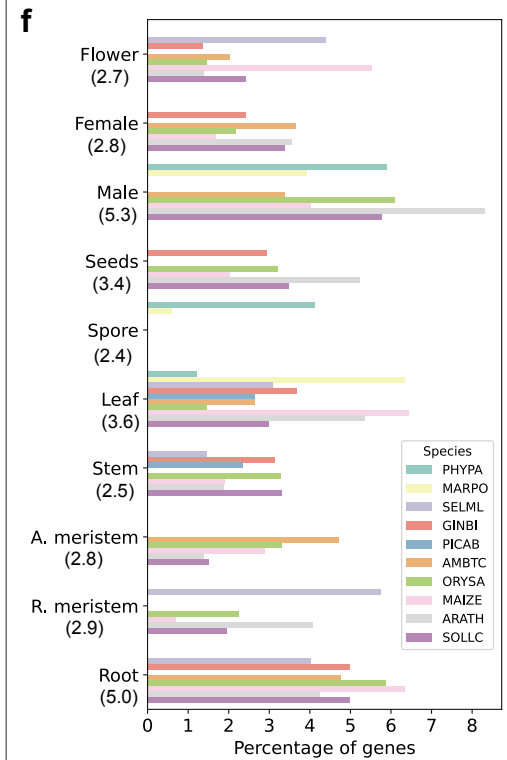
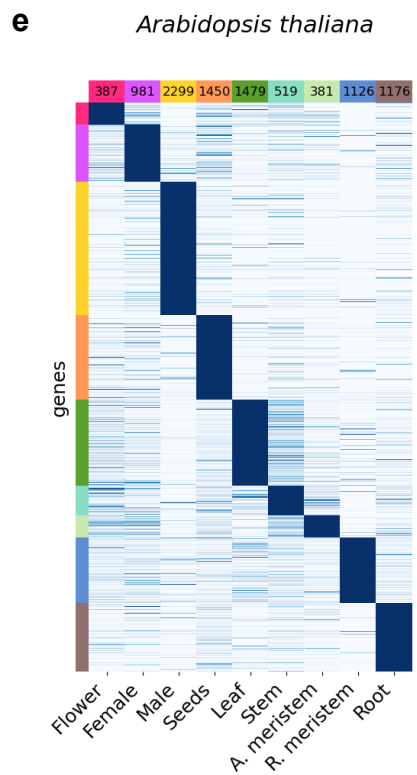
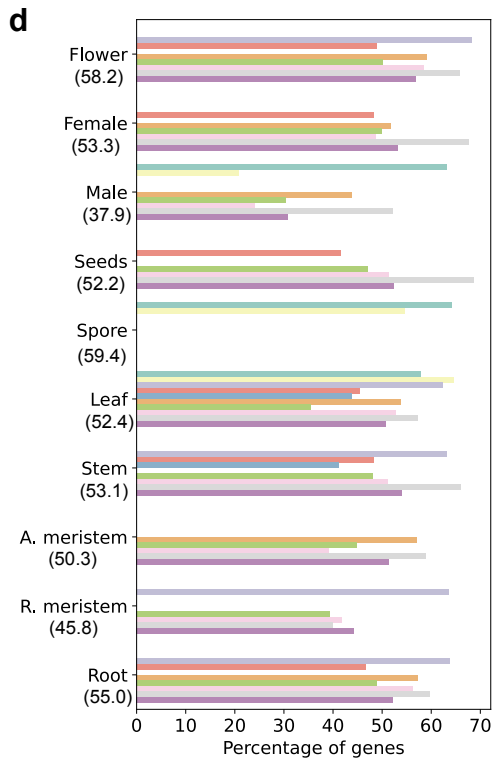
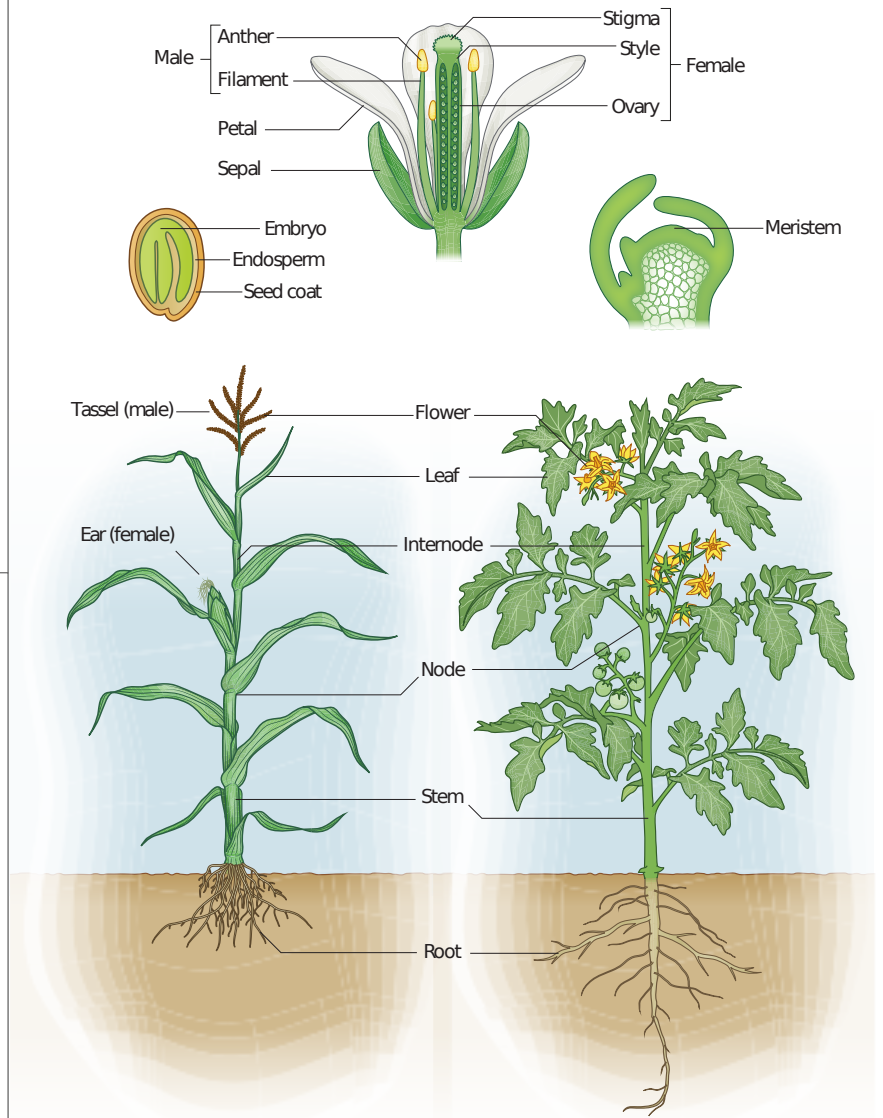
a *Physcomitrella*

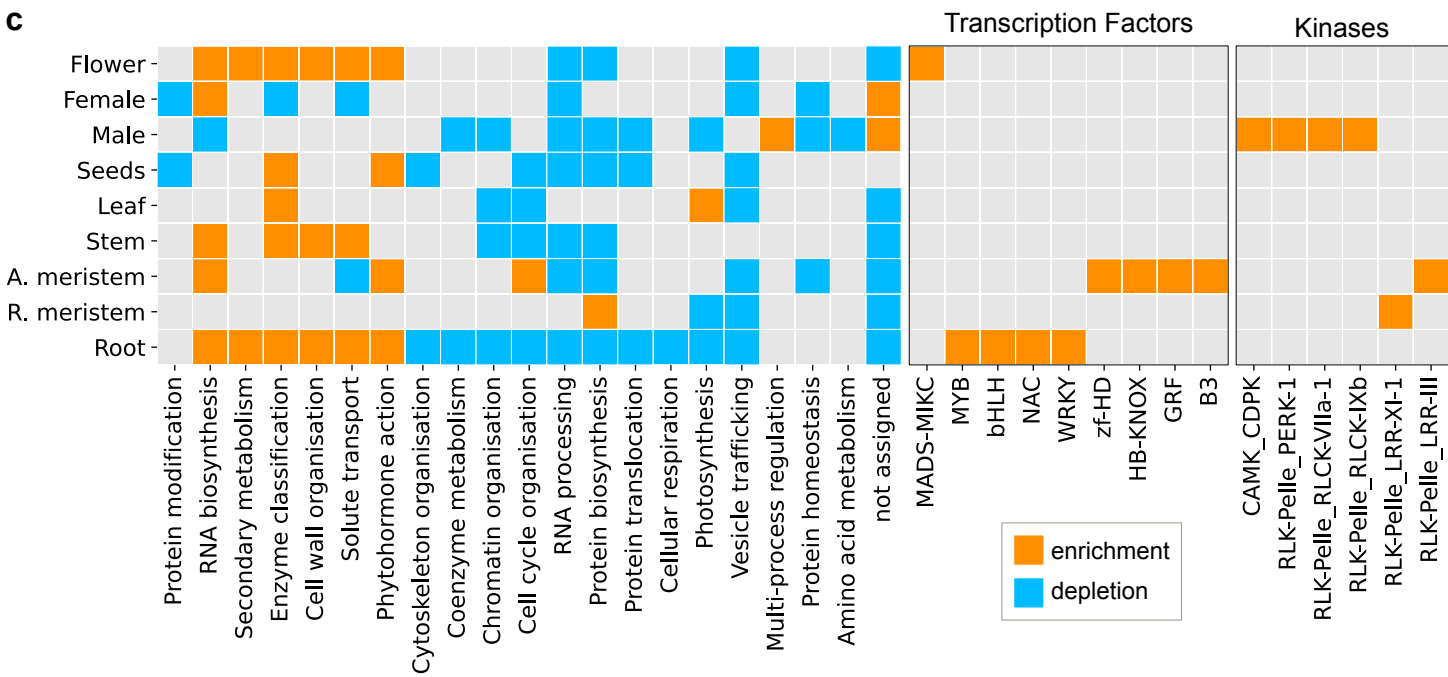
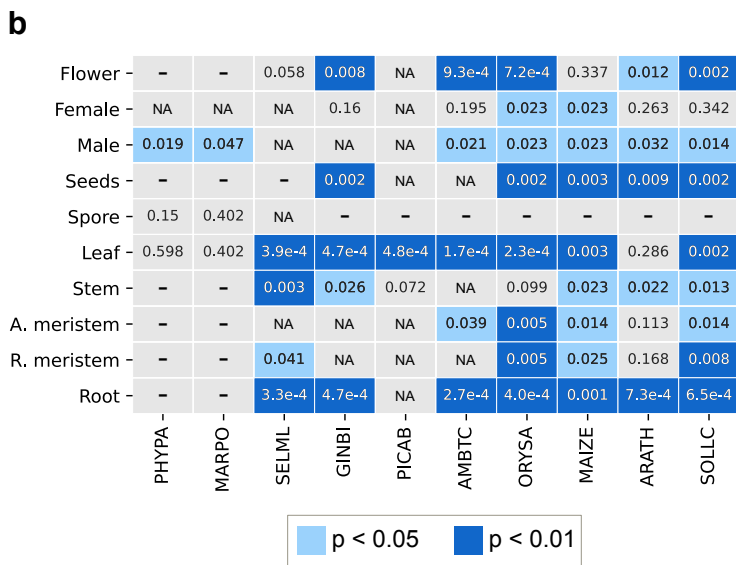
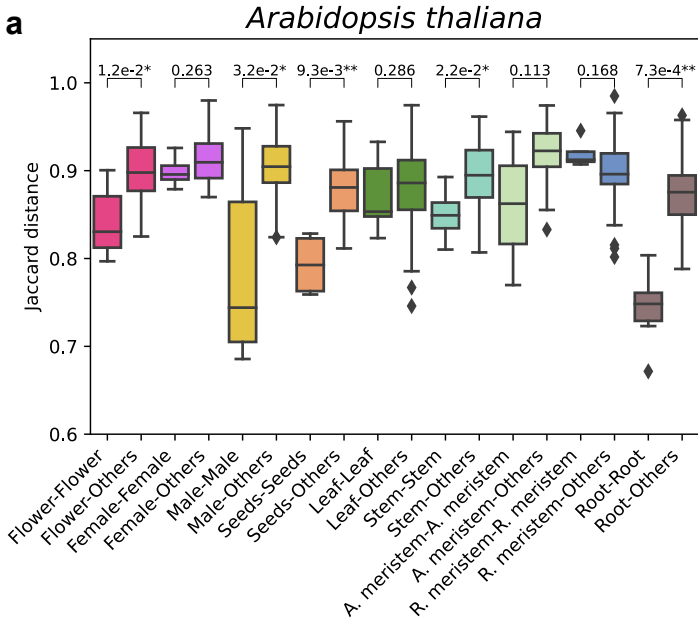


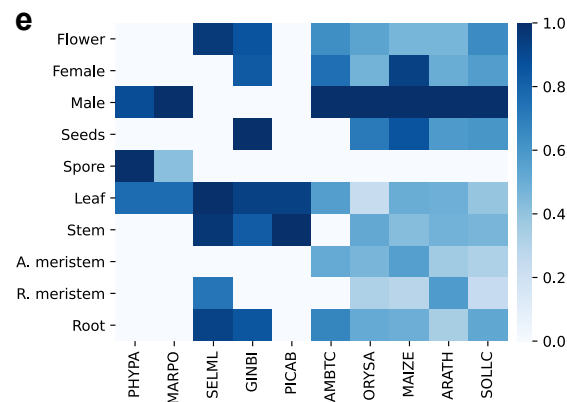
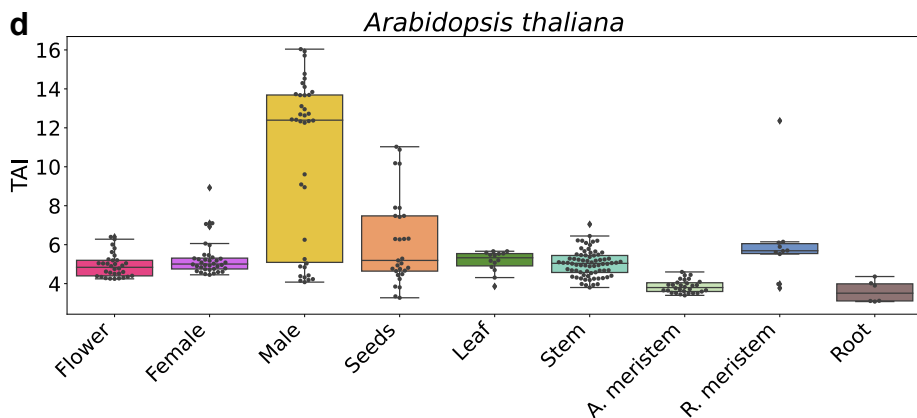
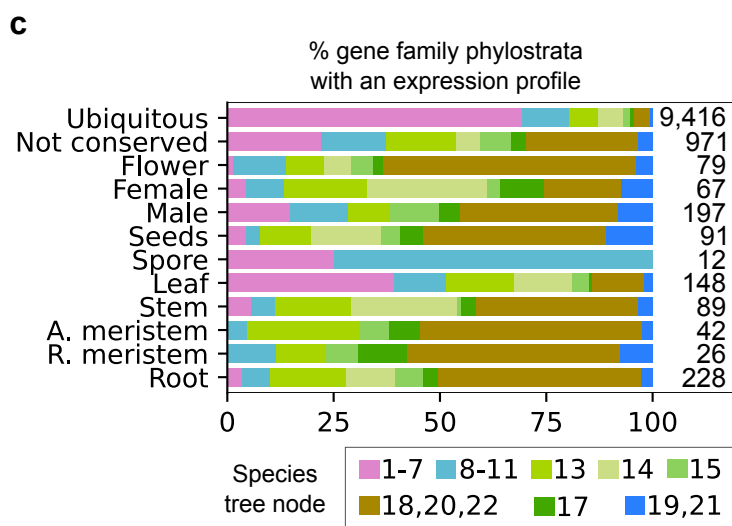
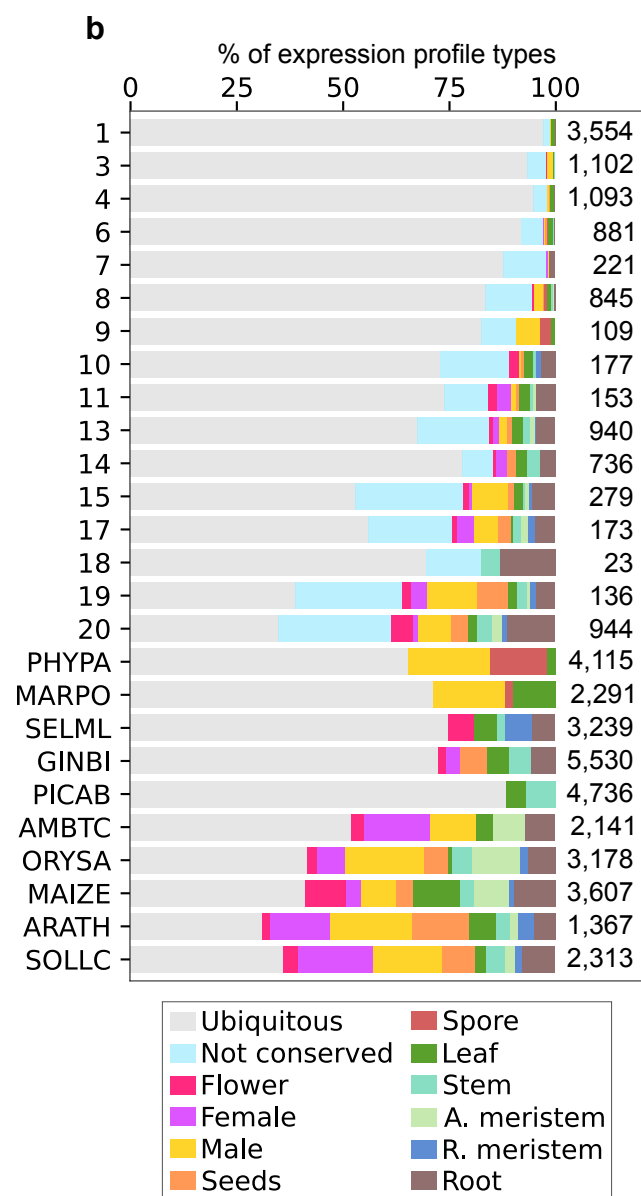
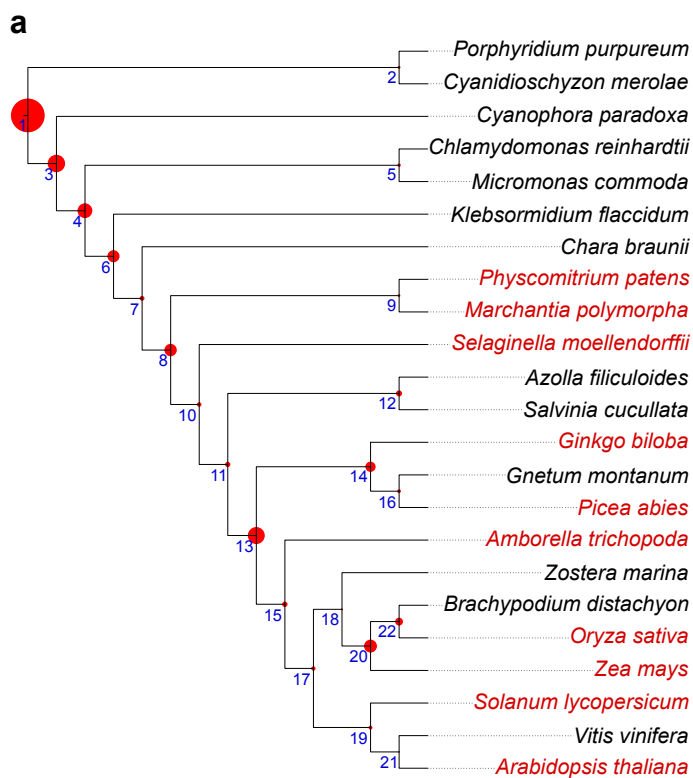
b *Marchantia*

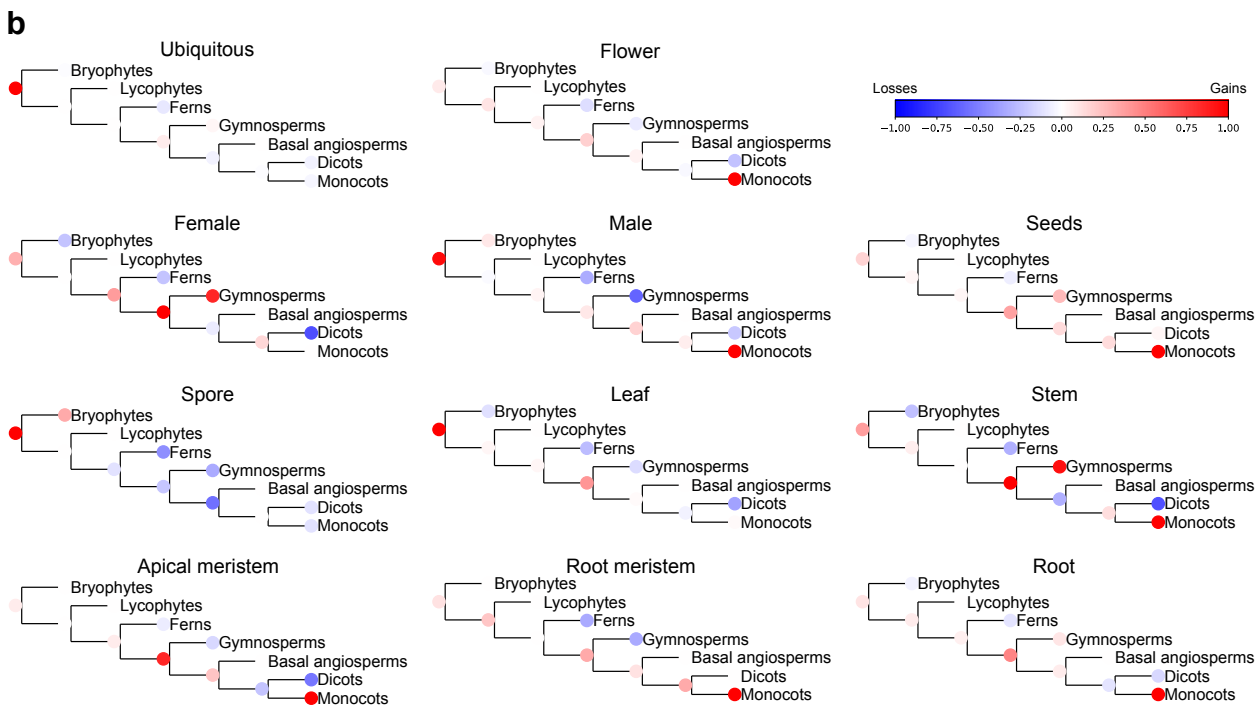
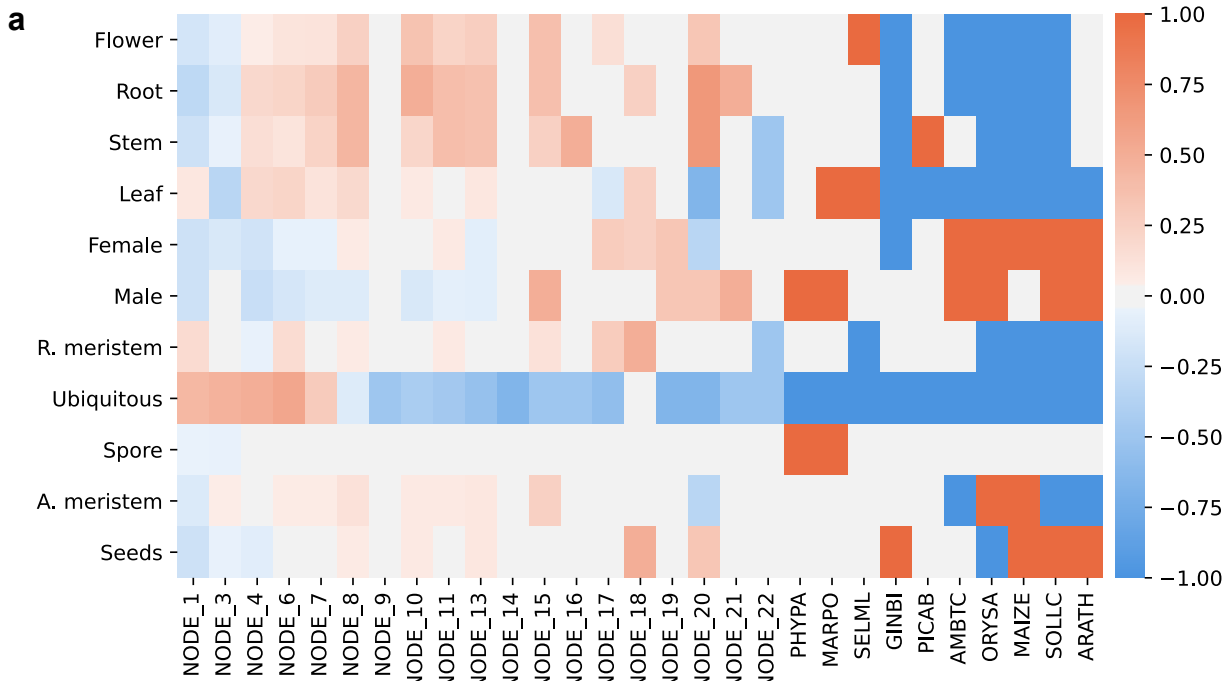


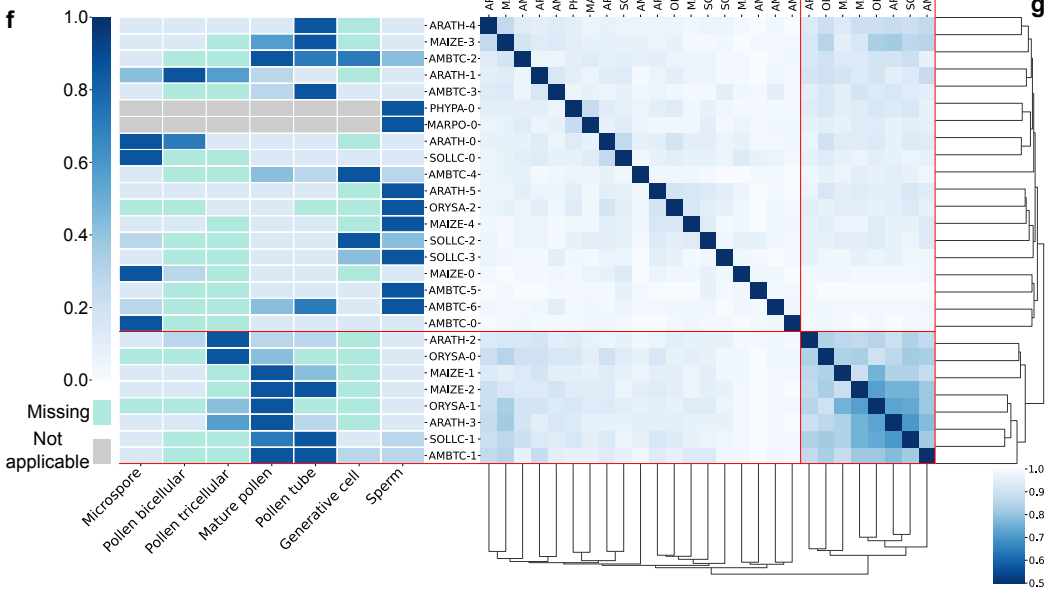
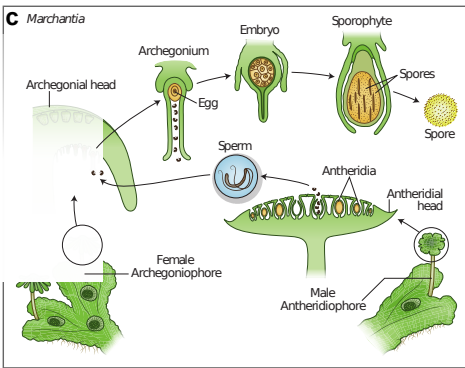
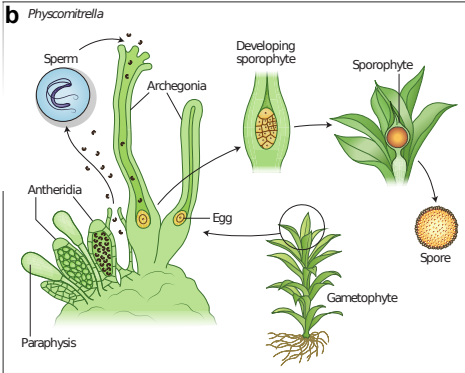
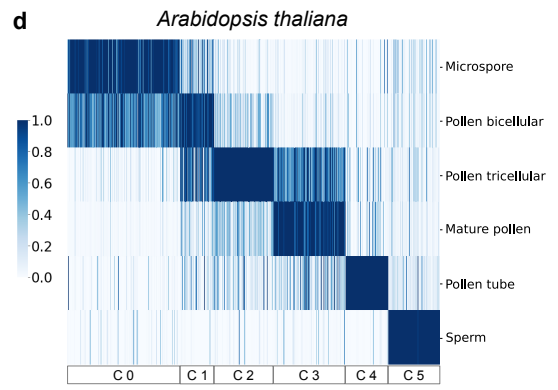
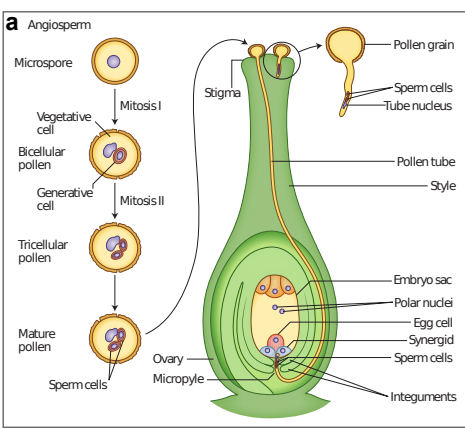
c Angiosperm

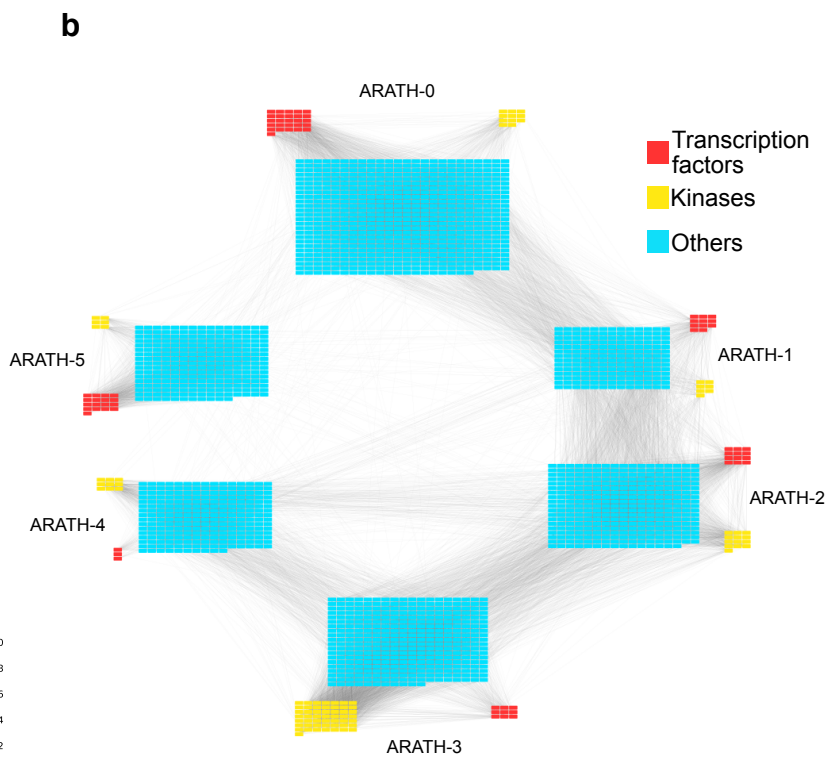
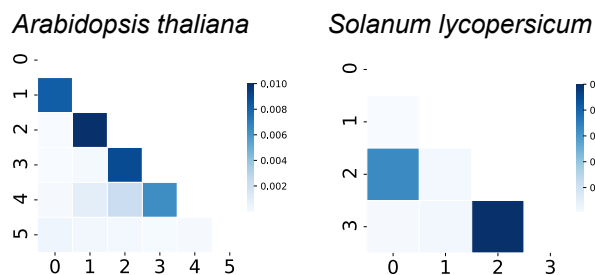
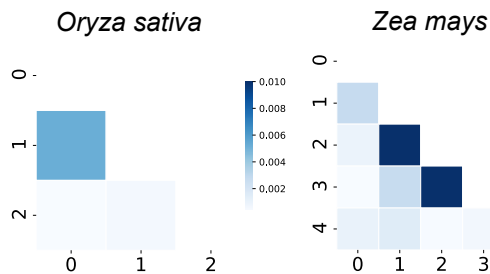
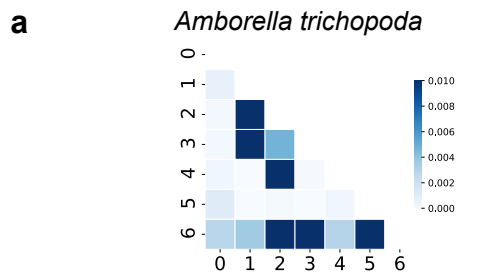




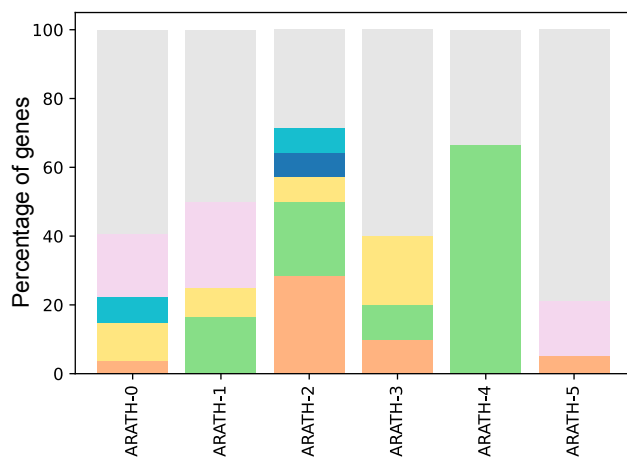








c Transcription factors



Kinases

