

1 **The evolution of the phenylpropanoid pathway entailed pronounced radiations and**
2 **divergences of enzyme families**

3
4 Sophie de Vries^{1,2*}, Janine MR Fürst-Jansen^{2*}, Iker Irisarri^{2,8*}, Amra Dhabalia Ashok², Till
5 Ischebeck^{3,4,5}, Kirstin Feussner^{3,4}, Ilka N Abreu³, Maike Petersen⁶, Ivo Feussner^{3,4,5}, Jan de
6 Vries^{2,7,8#}

7 *1 — Heinrich-Heine University Düsseldorf, Population Genetics, Universitätsstr. 1, 40225 Düsseldorf, Germany*

8 *2 — University of Goettingen, Institute for Microbiology and Genetics, Department of Applied Bioinformatics,*
9 *Goldschmidtstr. 1, 37077 Goettingen, Germany*

10 *3 — University of Goettingen, Albrecht-von-Haller-Institute for Plant Sciences, Department of Plant Biochemistry, Justus-*
11 *von-Liebig Weg 11, 37077 Goettingen, Germany*

12 *4 — University of Goettingen, Goettingen Center for Molecular Biosciences (GZMB), Goettingen Metabolomics and*
13 *Lipidomics Laboratory, Justus-von-Liebig Weg 11, 37077 Goettingen, Germany*

14 *5 — University of Goettingen, Goettingen Center for Molecular Biosciences (GZMB), Department of Plant Biochemistry,*
15 *Justus-von-Liebig Weg 11, 37077 Goettingen, Germany*

16 *6 — Philipps-Universität Marburg, Institut für Pharmazeutische Biologie und Biotechnologie, Robert-Koch-Str. 4, 35037,*
17 *Marburg, Germany*

18 *7 — University of Goettingen, Goettingen Center for Molecular Biosciences (GZMB), Department of Applied*
19 *Bioinformatics, Goldschmidtstr. 1, 37077 Goettingen, Germany*

20 *8 — University of Goettingen, Campus Institute Data Science (CIDAS), Goldschmidtstr. 1, 37077 Goettingen, Germany*

21
22 *these authors contributed equally

23 # author for correspondence:

24 Jan de Vries

25 Georg-August University Göttingen, Institute for Microbiology and Genetics, Department of Applied Bioinformatics,
26 Goldschmidtstr. 1, 37077 Göttingen, Germany, phone: +49-551-39-23755

27 devries.jan@uni-goettingen.de

28
29 **ORCID**s: Sophie de Vries 0000-0002-5267-8935, Janine MR Fürst-Jansen 0000-0002-5269-8725, Iker Irisarri 0000-0002-
30 3628-1137, Amra Dhabalia Ashok 0000-0001-5787-6941, Till Ischebeck 0000-0003-0737-3822, Ilka N Abreu 0000-0003-
31 4728-0161, Kirstin Feussner 0000-0003-1634-8258, Maike Petersen 0000-0001-7769-8556, Ivo Feussner 0000-0002-9888-
32 7003, Jan de Vries, 0000-0003-3507-5195

33
34 **AUTHOR EMAILS:** sophie.devries@uni-goettingen.de; Janine.fuerst-jansen@uni-goettingen.de; iker.irisarri@uni-
35 goettingen.de; amradhabalia.ashok@uni-goettingen.de; tischeb@uni-goettingen.de; kfeussn@uni-goettingen.de;
36 ilkanacif.deabreu@uni-goettingen.de; petersen@mailier.uni-marburg.de; ifeussn@uni-goettingen.de; devries.jan@uni-
37 goettingen.de

38
39 **KEYWORDS:** phenylpropanoid biosynthesis; plant evolution; evolution of gene families; evo-physio

40
41 **RUNNING HEAD:** Phenylpropanoid evolution: divergence & radiation

42 **ABSTRACT**

43 Land plants constantly respond to fluctuations in their environment. Part of their response is
44 the production of a diverse repertoire of specialized metabolites. One of the foremost sources
45 for metabolites relevant to environmental responses is the phenylpropanoid pathway, which
46 was long thought to be a land plant-specific adaptation shaped by selective forces in the
47 terrestrial habitat. Recent data have however revealed that streptophyte algae, the algal
48 relatives of land plants, have candidates for the genetic toolkit for phenylpropanoid
49 biosynthesis and produce phenylpropanoid-derived metabolites. Using phylogenetic and
50 sequence analyses, we here show that the enzyme families that orchestrate pivotal steps in
51 phenylpropanoid biosynthesis have independently undergone pronounced radiations and
52 divergence in multiple lineages of major groups of land plants; sister to many of these
53 radiated gene families are streptophyte algal candidates for these enzymes. These radiations
54 suggest a high evolutionary versatility in the enzyme families involved in the
55 phenylpropanoid-derived metabolism across embryophytes. We suggest that this versatility
56 likely translates into functional divergence and may explain the key to one of the defining
57 traits of embryophytes: a rich specialized metabolism.

58

59 **INTRODUCTION**

60 A diverse profile of specialized metabolites is one of the characteristics of land plants
61 (embryophytes). Almost any aspect of the biology of land plants is underpinned by
62 specialized metabolites—be it the phytohormones that are major modulators upstream in
63 various regulatory hierarchies (Scheres and van der Putten, 2017; Berens et al., 2017;
64 Blázquez et al., 2020) or pigments that give land plants their color and attune photochemical
65 properties (Jahns and Holzwarth, 2012).

66 A key aspect of the biological relevance of most specialized metabolites is their use
67 under challenging environmental conditions. Indeed, the elaboration of their specialized
68 metabolism is considered one of the drivers for the massive radiation of embryophytes on
69 land (Weng 2013). Moreover, a diversity of specialized metabolism likely played a key role
70 during the earliest steps of plants on land—allowing for the production of compounds that
71 protected land plants against the challenges of the terrestrial environment such as drought and
72 increased UV radiation (Rensing, 2017; de Vries and Archibald, 2018; Fürst-Jansen et al.,
73 2020; Jiao et al., 2020). A major pathway giving rise to a variety of specialized metabolites
74 that act in warding off environmental stressors is the biosynthesis of phenylpropanoids
75 (Dixon and Paiva, 1995; Dixon et al., 2002; Vogt, 2010).

76 The phenylpropanoid pathway is the source of precursors for thousands of metabolites
77 with multifaceted functions, and accounts for about 40% of organic carbon on earth (Vogt,
78 2010). One facet of these functions is that phenylpropanoid-derived compounds act as
79 structural polymers, foremost among which are the different types of lignin (Ralph et al.,
80 2004; Vanholme et al., 2012; Vanholme et al., 2019). Another prominent facet is that these
81 metabolites act as UV-protecting substances. While some of the best-known UV screens are
82 flavonoids, various other compounds stemming from the phenylpropanoid pathway are
83 equally potent UV protectants (Sheahan, 1996; Booij-James et al., 2000; Sytar et al. 2018;
84 Xue et al., 2020). The list of links between phenylpropanoid-derived compounds and the
85 response to environmental challenges could be continued; in fact, the response to almost any
86 abiotic stressor that plants face in the terrestrial habitat involves the action of
87 phenylpropanoid-derived compounds (for comprehensive reviews see, e.g., Dixon and Paiva,
88 1995; Vogt, 2010). Furthermore, phenylpropanoid-derived compounds are involved in the
89 defense responses against plant pathogens in many land plant lineages (Danielsson et al.
90 2011; Ponce de Leon et al., 2012; König et al., 2014; Overdijk et al., 2016; Carella et al.,
91 2019).

92 All embryophytes make use of the enzymatic routes in the phenylpropanoid pathway.
93 For example, the utilization of flavonoids under UV stress appears to be a conserved response
94 across Embryophyta (Wolf et al., 2010; Clayton et al., 2018). However, not all embryophytes
95 produce the same compounds under the same stress conditions—in contrary, the diversity of
96 compounds is immense. Major differences in the biosynthesis of phenylpropanoid-derived
97 compounds occur in distinct lineages of land plants. This includes specialized roles such as
98 the flower coloration determining anthocyanins that attract pollinators (Miller et al., 2011;
99 Sheehan et al., 2012); such a role of anthocyanins is obviously limited to flowering plants and
100 can vary even among closely related species (Saito and Harborne, 1992). That said,
101 Piatkowski et al. (2020) phylogenetically inferred that orthologs for the entire anthocyanin
102 biosynthesis pathway are already present in the ancestor of seed plants and more than half of
103 the important orthogroups were already present in the most recent common ancestor of all
104 land plants. An important recent insight into the deep evolutionary roots of flavonoid
105 biosynthesis was the discovery of auronidins—a novel class of red flavonoid pigments that
106 are synthesized in the bryophyte *Marchantia polymorpha* (Berland et al. 2019). Further, for
107 example, Renault et al. (2017a) reported on the enrichment of the *Physcomitrium patens*
108 (moss) cuticle in phenolic compounds—an enrichment that hinges on the action of a
109 cytochrome P450 enzyme that is orthologous to enzymes that act in lignin biosynthesis; the

110 production of lignin might trace its evolutionary roots back to an ancient set of enzymes
111 acting in the production of complex, phenol-enriched polymers (Renault et al., 2019).
112 Carella et al. (2019) showed that the liverwort model plant *Marchantia polymorpha* triggers
113 phenylpropanoid biosynthesis upon attack by the oomycete phytopathogen *Phytophthora*
114 *palmivora*. Similar responses towards phytopathogens are known from gymnosperms (Oliva
115 et al., 2015) and angiosperms (Dixon and Paiva, 1995; Bednarek et al., 2005; Kaur et al.,
116 2010; Chezem et al., 2017; Carella et al., 2019). Thus, all land plants use the core framework
117 of the phenylpropanoid pathway to produce—often lineage-specific—variations of
118 phenylpropanoid derivatives that aid in response to biotic and abiotic stressors.

119 The production of the chemical repertoire of land plants is often catalyzed by
120 members of large enzyme-coding gene families (Shockey et al., 2003; Nelson and Werck-
121 Reichhart, 2011; Renault et al., 2017b), and this also seems to be the case for the enzymes
122 involved in phenylpropanoid biosynthesis (Hamberger et al., 2007; Xu et al., 2009; Vogt,
123 2010). It is thus conceivable that various adaptive forces have shaped the families of enzymes
124 that act in the phenylpropanoid pathway, leading to multiple independent cases of sub- and
125 neofunctionalization (see also Rensing, 2010). An inference of the common (minimal) set of
126 enzymes that were present in the last common ancestors (LCA) of (i) streptophytes, (ii) land
127 plants and their closest streptophyte algal relatives, and (iii) land plants can thus shed light on
128 which enzymatic building blocks evolution acted upon to give rise to the elaborate chassis of
129 the phenylpropanoid pathway.

130 The phenylpropanoid pathway has long been considered to be specific to
131 Embryophyta. However, homologs of the genes coding for the enzymes that constitute the
132 embryophytic phenylpropanoid pathway can be found in extant algal relatives of land plants,
133 suggesting that they were already present in a common ancestor shared by streptophyte algae
134 and land plants (de Vries et al., 2017; Renault et al., 2019; Maeda and Fernie, 2021). Since
135 the beginning of 2020, we have genome data from all major lineages of Streptophyta—except
136 Coleochaetophyceae (Szövényi et al., 2021); only using this extended repertoire of species
137 and sequences allows us to pinpoint which subfamilies and/or which ancestral enzyme of
138 multiple subfamilies were present in the aforementioned LCAs. Compounds that, in land
139 plants, emerge from the phenylpropanoid pathway are indeed found in algae; these include
140 flavonoids and lignin-like compounds in streptophyte algae (Delwiche et al., 1989; Sørensen
141 et al., 2011; Jiao et al., 2020) and core phenylpropanoid building blocks as well as flavonoids
142 in a phylodiverse set of algae (Goiris et al. 2014). Interestingly, lignin-like compounds were
143 even found in distantly-related red macroalgae (Martone et al., 2009)—although this is likely

144 a case of convergence that builds on an unknown enzymatic framework. However, even
145 within the green lineage (Chloroplastida), the question of the deep evolutionary roots of the
146 phenylpropanoid pathway is wide open.

147 Investigations of the algal relatives of land plants have strongly benefitted from recent
148 progress in phylogenomics on plants and algae. A major outcome of these recent
149 phylogenomic analyses was that the Zygnematophyceae have been pinpointed as the class of
150 algae most closely related to land plants (Wodniok et al., 2011; Wickett et al., 2014; Leebens-
151 Mack et al., 2019). Hand in hand with these phylogenomic efforts went the generation of
152 genomic (Hori et al., 2014; Nishiyama et al., 2018; Cheng et al., 2019; Wang et al., 2020;
153 Jiao et al., 2020) and transcriptomic data on streptophyte algae (Ju et al., 2015; Rippin et al.,
154 2017; de Vries et al., 2018; de Vries et al., 2020). Additionally, critical gaps in the land plant
155 tree of life have been filled; this includes recent publications of the first genomes of
156 liverworts (Bowman et al., 2017), ferns (Li et al., 2018) and hornworts (Szövényi et al., 2015;
157 Li et al., 2020; Zhang et al., 2020). These data allow for the fine-grained tracing of the
158 evolution of key plant enzyme families across the green tree of life. Recent studies have
159 illuminated the diversity of enzymes in the routes towards flavonoids and anthocyanins as
160 well as the PAL-dependent pathway of salicylic acid biosynthesis via benzoic acid
161 (Piatkowski et al., 2020; Güngör et al., 2021, de Vries et al., 2021).

162 In this study, we infer the evolutionary history of eleven critical enzyme families
163 known to be woven into the mesh of routes from phenylpropanoids to lignin biosynthesis in
164 land plants; we have paid particular attention to the routes leading to the biosynthesis of
165 lignin. We use the new diversity of genomic and transcriptomic data from land plants as well
166 as streptophyte and chlorophyte algae to infer the origin of these large gene families. The
167 datasets were chosen in a manner that they cover the breadth of streptophyte diversity while
168 providing a balanced sampling; the latter is especially relevant in light of the high number of
169 genomes of flowering plants. We aimed to include at least one representative of each of the
170 major lineages of streptophytes in the datasets we surveyed. Our data pinpoint deep homologs
171 of candidate enzymes in streptophyte algae for L-phenylalanine ammonia-lyase (PAL), 4-
172 coumarate-CoA ligase (4CL), caffeoyl-CoA O-methyltransferase (CCoAOMT); further, for
173 streptophyte and chlorophyte algae, we pinpoint homologs for cinnamoyl-CoA reductase
174 (CCR), cinnamyl alcohol dehydrogenase (CAD), and potentially relevant monoacylglycerol
175 lipases (MAGLs). Further, we find that often the functionally characterized enzymes of the
176 core phenylpropanoid and lignin biosynthesis routes derive from lineage-specific radiations,
177 limiting the inference of function outside the model system. Nonetheless we could infer

178 which subfamilies were present in LCAs along the trajectory of streptophyte evolution, even
179 though ancestral functional inference was limited. That said, for enzyme families with deep
180 homologs, we approximated the function through domain prediction and the conservation (or
181 lack thereof) of key residues of known functional importance. We found that all enzyme
182 families underwent several lineage-specific expansions and losses as well as bursts in growth
183 of enzyme families that occurred early during the radiation of land plants. We hypothesize
184 that lineage-specific expansions in these enzyme families is linked with the diversity of
185 lineage-specific phenylpropanoid derivatives and functions that occur in the species analyzed
186 here.

187

188 **RESULTS AND DISCUSSION**

189 **The checkered occurrence of phenylalanine ammonia-lyase among streptophyte algae**

190 The conversion of the aromatic amino acid phenylalanine and/or tyrosine into cinnamate
191 and/or *p*-coumarate is the first step of the plant phenylpropanoid pathway (Figure 1). This
192 first committed step is catalyzed by PAL and the bifunctional L-phenylalanine/ L-tyrosine
193 ammonia-lyase (PTAL; Barros and Dixon, 2020). For a long time, it was thought that among
194 Chloroplastida, PAL/PTAL were limited to land plants; their gain was considered to have
195 occurred via a lateral gene transfer event that has occurred at the base of the land plant clade
196 (Emiliani et al., 2009). Recently, however, genes coding for putative PAL-like enzymes were
197 detected in streptophyte algae, such as the filamentous streptophyte alga *Klebsormidium*
198 *nitens* (de Vries et al., 2017). In light of the recent surge in available genomes from across the
199 green tree of life, we set out to explore the evolutionary history of PAL.

200 Using *AtPAL1* as a bait sequence, we screened protein data from diverse land plants
201 and all streptophyte algal genomes available. Hits among streptophyte algae fell into three
202 categories: (i) proteins of between 480 (*Klebsormidium nitens* PAL, kfl00104_0290_v1.1) and
203 527 amino acids (two homologs in *Chara braunii*; g57646_t1 and g34530_t1); (ii) short
204 proteins such as ME1156409C09523 of *Mesotaenium endlicherianum*, which is 184 amino
205 acid in length; (iii) long proteins of between 991 (*Chlorykbus atmophyticus*
206 Chrsp482S06115) and 1115 amino acids (*Klebsormidium nitens* kfl00024_0250_v1.1).
207 Proteins falling into the third category are fusions of an aromatic amino acid lyase domain
208 and a putative tRNA synthetase; homologs of such also occur in land plants (e.g.,
209 AT3G02760 and Os05g0150900). The short proteins of category ‘(ii)’ are found in some
210 Zygnematophyceae and are proteins of unknown function with a putative HAL domain.
211 Genomic artefacts leading to this result can be excluded given that these types of protein-

212 encoding genes have been recovered for independent Zygnematophyceae that are likely >500
213 million years divergent from one another. The proteins in category ‘(i)’ are those with the
214 highest identity to *bona fide* land plant PALs. This category includes the promising PAL
215 candidate kfl00104_0290_v1.1 (see de Vries et al., 2017). We therefore set out to further
216 explore these PAL-like candidates, which noteworthyly were only found in the genome of
217 *Klebsormidium nitens* and *Chara braunii* and no other streptophyte algal genome data.

218 To understand the evolutionary history of PAL in streptophytes, we computed a
219 maximum likelihood phylogeny (Figure 2). The phylogenetic analysis included the
220 aforementioned PAL homologs from diverse Streptophyta as well as bacterial and fungal
221 PALs. In agreement with previous studies (Emiliani et al., 2009; de Vries et al., 2017), the
222 fungal and bacterial PAL sequences are closely related (bacterial clade: bootstrap support
223 100; fungal clade: bootstrap support 85) to the clade of land plant PALs (bootstrap support:
224 96). Additionally, we included diverse eukaryotic and prokaryotic histidine ammonia lyases
225 (HALs) based on the set obtained from de Vries et al. (2017), which, as in this latter study,
226 form two clades with eukaryotic and prokaryotic HALs (bootstrap support 83 and 97). All
227 putative PAL-like candidate sequences from *Chara braunii* clustered with HAL sequences,
228 one (*Chara braunii* g66119_t1) with a cyanobacterium (bootstrap support 100), and two
229 (*Chara braunii* g34530_t1 and *Chara braunii* g57646_t1) sister to an entire bacterial HAL
230 clade (bootstrap support 100). Both sequences retrieved for *Klebsormidium nitens*, in
231 contrast, clustered with the PAL clades, one *K. nitens* kfl00024_0250_v1.1 with low support
232 (bootstrap level 61) and a rather long branch as sister to plant and fungal PALs, showing that
233 its placement is not fully resolved and that further analyses are required to identify its true
234 identity. The second *Klebsormidium nitens* sequence, kfl00104_0290_v1.1, clustered within
235 the clade of bacterial PALs (bootstrap support 100), of which some were already functionally
236 characterized—for example, the characterized PAL of *Nostoc punctiforme* (Moffit et al.,
237 2007). This is in agreement with the placement of this protein sequence in de Vries et al.
238 2017 and supports it as a putative PAL sequence.

239 In sum, the evolutionary origin of streptophyte PAL appears to be complex and
240 remains obscure: it may be that PAL had a distinct origin in streptophyte algae and land
241 plants, yet the pattern may also be explained by an origin via an endosymbiotic gene transfer
242 from the cyanobacterial plastid progenitor that was retained by streptophytes with a gain of
243 an extra C-terminal domain later in the evolution of embryophytes, resulting in the two
244 distinct PAL clades. It is however important to note that the 3’ region of the genomic locus
245 that codes for the shorter *K. nitens* protein kfl00104_0290_v1.1 contains sequence

246 information that resembles code for the missing C-terminal stretch; thus, the C-terminal
247 stretch might have simply been secondarily lost in *K. nitens*. Independently, the presence of
248 PALs in fungi further complicates the evolutionary scenario. Rampant gene losses during
249 eukaryotic evolution or convergent domain acquisitions, as well as horizontal gene transfer
250 from plants (as hypothesized by Emiliani et al., 2009) are other scenarios that can explain the
251 evolutionary origin of PALs in streptophytes and thus ultimately in land plants.

252

253 **Streptophyte algae have an expanded and divergent repertoire of cytochrome P450** 254 **monooxygenases with no clear C4H orthologs**

255 After the synthesis of cinnamate by PAL, two routes open up (Figure 1). One of them is the
256 conversion of cinnamate into *p*-coumarate, which is catalyzed by cinnamate **4**-hydroxylase
257 (C4H). C4H belongs to the large class of CYP450 enzymes, which are present among all
258 domains of life (Omura, 1999). Among the CYP450 enzymes C4H belongs to the CYP450
259 subfamily 73 (CYP73). In land plants, CYP450s have undergone massive duplication and
260 subfunctionalization, underpinning the specialized metabolic capabilities of embryophytes
261 (Nelson and Werck-Reichhart, 2011); for example, the CYP73 subfamily belongs to the
262 larger CYP71 clan. The specific CYP450 monooxygenases that fall into the group of C4H
263 appear to be limited to land plants: clear orthologs can be found in bryophytes and
264 tracheophytes (Emiliani et al., 2009; de Vries et al., 2017). That said, the product of the
265 reaction carried out by C4H in land plants (*p*-coumarate) has been detected via UHPLC-
266 MS/MS in phylodiverse algae (Goiris et al., 2014). Therefore, there appears to exist a route
267 towards *p*-coumarate that is either independent of C4H via direct transformation of tyrosine
268 by PTAL or carried out by a highly divergent C4H homolog. PTALs have so far however
269 been observed in monocots (Barros et al., 2016; Barros and Dixon, 2020), suggesting a
270 different CYP73 subfamily enzyme that may carry out the reaction. Owing to the recent
271 increase in genomic data available for streptophyte algae, we revisited the question of when
272 C4H-based *p*-coumarate might have emerged and explored CYP450 evolution.

273 We sampled C4H homologs from seven land plant genomes that had a BLAST bit
274 score (a normalized alignment score) of at least 200, as well as seven streptophyte algal and
275 five chlorophyte algal genomes that had a bit score of at least 100. We aligned all C4H
276 homologs and computed a maximum likelihood phylogeny (Figure S1). The well-
277 characterized C4H of *Arabidopsis* fell into a clade with full (100 out of 100) bootstrap
278 support; this clade included at least one C4H homolog from each of the other six land plants
279 genomes, corroborating the notion that all land plants have C4H orthologs, which appear

280 conserved in their function from bryophytes to tracheophytes (Russel and Conn, 1967; Urban
281 et al., 1994; Ro et al., 2001; Wohl and Petersen, 2020) and thus since the LCA of land plants.
282 However, no algal sequences fell into this clade. That said, we observed four well-supported
283 clades of streptophyte algal CYP450 enzymes (Figure S1). Investigating the genetic
284 distances, we find that some of the streptophyte algal sequences have a closer genetic
285 distance to the C4H-like clade than to sequences from land plants (including *Arabidopsis*
286 *thaliana*) from other CYP450 subfamilies (Table S1). While these sequences remain of
287 unknown function, they are candidates for the CYP450 enzyme family that catalyzes the
288 C4H-function in algae.

289

290 **A deep split of streptophyte 4CL/ACS**

291 The second route that opens up after the PAL-dependent step is the conversion of cinnamate
292 into cinnamoyl-CoA. This is carried out by the AMP-forming synthetase/ligase 4CL and
293 potentially other enzymes annotated as acyl-CoA synthetases (ACS/ACoS) (Shockey et al.,
294 2003; Figure 1). Altogether, these enzymes belong to a large family of distantly related acyl-
295 activating enzymes (AAEs), such as the long-chain acyl-CoA synthetases (LACS) and many
296 more (Shockey et al., 2003; Figure S2). Homologs with affinity to 4CL appear to occur
297 across chlorophytes and streptophytes (Labeeuw et al., 2015). At least in *Arabidopsis*, the
298 family of 4CLs has expanded and includes four canonical (“4CL”) and nine additional 4CL-
299 like (“4CLL”) members, falling into AAE clade IV and V as defined by Shockey et al.
300 (2003). We thus set out to understand what the 4CL repertoire of the last common ancestor of
301 land plants and the one shared with algae might have looked like.

302 In order to trace the radiation of 4CLs across the green tree of life, we sampled 4CL
303 homologs from genomes of nine land plants, seven streptophyte algae (plus four
304 transcriptomes of streptophyte algae) and five chlorophyte algae that had a minimum of 400
305 and a maximum of 1150 amino acids in length and showed affinity to the 4CL clade in a
306 larger phylogenetic survey (Figure S2). We recovered a large clade (bootstrap support 85)
307 that included all *bona fide* 4CL paralogs and ACOS5 of *Arabidopsis thaliana* (Figure 3);
308 ACOS5 has been previously associated with the *bona fide* 4CL clade (Shockey et al., 2003)
309 but it did show only inconsistent activity on typical substrates of 4CL (Costa et al., 2005) and
310 appears to have a very specific function in sporopollenin biosynthesis of pollen (de Azevedo
311 Souza et al., 2009). These observations agree with the presence of different (but conserved)
312 amino acids at sites that bind hydroxycinnamate in typical 4CLs (Figure 3), which might
313 suggest a different natural substrate for ACOS given that affinity is mostly determined by the

314 binding pocket size (Hu et al. 2010). We further recovered the angiosperm-specific separation
315 defined by Ehling et al. (1999) into class I and class II 4CLs. Additional lineage-specific
316 radiations occurred, for example in *Physcomitrium patens*, which fell into a clade of
317 bryophyte sequences (bootstrap support of 70) and in *Selaginella moellendorffii* (spread out
318 over the fully-supported clade of 4CLs). The common ancestor of land plants appears to have
319 possessed an ACOS5-like and one 4CL-like gene, all other 4CL paralogs in this clade likely
320 emerged later during land plant evolution. Clustering with AAE clade IV (including
321 *AtACOS5*, and *At4CL1,2,3* and 5 (bootstrap support 85)) are sequences from five
322 streptophyte algae. Each of the five streptophyte algae possesses one homolog to these five
323 types of AMP-forming ligases with 4-coumarate-CoA synthesizing activity. When we
324 predicted the tertiary structure of *Chlorokybus atmophyticus* Chrsp175S02417 and *Penium*
325 *margaritaceum* 006213.t1 via Iterative Threading ASSEMBLY Refinement (I-TASSER; Zhang
326 2008), we recovered, in both cases, firefly luciferases as best match (TM-scores 0.852 and
327 0.918; 1BA3 and 2D1S; Franks et al., 1998; Nakatsu et al., 2006). Investigating the putative
328 structure of the other streptophyte algal sequences (*Zygnema circumcarinatum*
329 DN42558_c0_g1_i1, *Spirogyra pratensis* 3442_c2_g1_i6, and *Klebsormidium nitens*
330 00016_0470_v1.1), however, always recovered *Populus tomentosa* 4CL (3A9U; Hu et al.,
331 2010) as their closest structural analog (TM scores of 0.957, 0.969, and 0.961, respectively).
332 Hence, we hypothesize that a 4CL/ACOS5-like encoding gene was present in the last
333 common ancestor of all streptophytes. Underpinning this hypothesis is that the sequence of
334 the amino acids in the binding pocket in the streptophyte algal 4CL homologs is consistent
335 with that of 4CL homologs from other land plants (including that of *Arabidopsis thaliana*;
336 Figure 3). Further, the amino acids relevant for the enzymatic function (i.e., the residues
337 KQK involved in adenylation, nucleophilic substitution and coumaroyl-AMP cleavage) are
338 also conserved across most 4CL/ACOS5 sequences, including those of the streptophyte
339 algae. Variation in these residues is already apparent in 4CLL homologs and outside of the
340 4CL/ACOS5/4CLL clade these residues show high variability (Figure 3).

341 A similar pattern was observed when we investigated the domain structure of all
342 recovered sequences. Most of the 4CL/ACOS5-like sequences contained four domains:
343 Phosphopantetheine binding ACP domain (IPR025110), AMP-binding, conserved site
344 (IPR020845), AMP-dependent synthetase-like superfamily (IPR042099), AMP-dependent
345 synthetase/ligase (IPR000873; Figure S3). There were four exceptions to this pattern. They
346 include one sequence from the water fern *Azolla filiculoides* (Azfis0013.g013344) and two
347 hornwort sequences from *Anthoceros agrestis* BONN (Sc2ySwM344.2803.3 and

348 Sc2ySwM344.2803.4), which all missed the Phosphopantetheine binding ACP domain
349 (IPR025110). The other exception was *AtACOS5*, which is the only sequence in this clade
350 that missed the conserved AMP-binding site (IPR020845). The domain pattern is similar
351 across the 4CLL-like clade, too. Yet, more sequences miss either the IPR025110 and/or the
352 IPR020845 domain. The streptophyte algal sequences within the ACOS5/4CL clade
353 contained all four domains, whereas algal sequences outside of this clade missed at least
354 one—but recovered several other domains. These additional domains are not conserved
355 within the phylogenetic subclades of these algal sequences and only exceptionally occur in
356 the *4CL/ACOS4/4CLL* (two sequences) or *Other AMP-dependent synthetase and ligase*
357 *family protein* clades (two sequences).

358 We recovered a second clade of spermatophyte sequences (bootstrap support 89)
359 representing AAE clade V enzymes, which contains 4CLL8 and several other ATP-ligases of
360 *Arabidopsis thaliana* with predicted 4CL activity, including OPCL1 (matching 4CL-like 5
361 with a 100% amino acid identity according to Uniprot). Sequences in this clade, however,
362 diverge in the amino acids that are involved in the formation of the binding pocket in the
363 canonical 4CLs (Figure 3), which might point to a different substrate preference of the
364 enzymes in this clade. In fact, OPCL1 and many of these “4CLLs” (e.g., AT5G63380,
365 AT1G20500, AT4G05160) showed higher activity on fatty acids and fatty acid-derived
366 precursors for the phytohormone jasmonic acid than cinnamate-derived compounds in an *in*
367 *vitro* substrate survey carried out by Kienow et al. (2008). It is thus questionable that the
368 enzymes of this clade act as *bona fide* 4CLs. Homologs to these sequences are found in
369 *Brachypodium distachyon* and *Picea abies*, suggesting an origin in the last common ancestor
370 of seed plants followed by two duplication events with either (a) both taking place in the last
371 common ancestor of angiosperms or (b) one early on in the LCA of seed plants and the
372 second in the LCA of angiosperms. Bootstrap support to include the *Picea abies* sequences in
373 the clade containing AT5G63380 (4CLL9) is however low (bootstrap 54). Each duplication
374 event was followed by independent lineage-specific radiations giving rise to a whole plethora
375 of possible candidates for 4CL, but also a large evolutionary potential with regard to substrate
376 specificity and flexibility. The 4CLL clade of spermatophytes is nested within a larger, lowly
377 supported clade (bootstrap 65) that included sequences from across the land plant tree of life.
378 Here, pronounced and independent expansion occurred in most of the major lineages of land
379 plants, leading to large clades of, for example, proteins of the hornwort *Anthoceros* and the
380 lycophyte *Selaginella*. As noted above, most of the *Selaginella* and some of the *Anthoceros*
381 homologs retained the conserved KQK residues required for the catalytic activity but others

382 did not (Figure 3), which suggests the presence of species-specific functions. Outside of the
383 entire 4CL-ACOS and 4CLL clade ('Streptophyte 4CL/ACOS/4CLL-likes'; bootstrap 84)
384 clustered various highly divergent ATP-dependent synthetases and ligases that exist
385 throughout the green tree of life including sequences from chlorophytes. None of these
386 synthetases and ligases retained the catalytic triad KQK.

387 Altogether, our phylogenetic data indicated that a 4CL/ACOS5-like encoding gene
388 was present at the base of Streptophyta. Domain annotation and the analysis of amino acid
389 patterns in the binding pocket and functional sites support this idea. Further, the similarity of
390 these residues between the candidates of streptophyte algal homologs for 4CLs and the
391 sequences of 4CL proteins with high activity on cinnamate derivatives as substrates (see also
392 Costa et al., 2005), indicates that 4CL activity may evolved more than 700 million years ago
393 in streptophytes.

394

395 **Patchy distribution of CCR-likes in streptophyte algae and massive independent** 396 **radiations in land plants**

397 *En route* to the production of different lignin monomers is the NADPH-dependent reduction
398 of the activated acyl-group of the phenylpropanoid backbone molecules. This first step
399 towards an aldehyde functionality is carried out by CCR (Figure 1), which falls into a larger
400 family of NADPH-dependent reductases, including dihydroflavonol reductases (DFRs) and
401 DFR-likes (DFL) (Lacombe et al., 1997; Devic et al., 1999). We previously reported the
402 presence of CCR-like protein sequences in streptophyte algae (de Vries et al. 2017). Since
403 these previous analyses, however, genome data on additional major lineages of land plants
404 and streptophyte algae have become available.

405 With these new data at hand, covering most major lineages of streptophyte algae and
406 all major lineages of land plants, it is now possible to infer the evolutionary history of CCR-
407 like and DFR-like sequences. We computed a maximum likelihood phylogeny of CCR
408 homologs with a minimum of 220 amino acids that we detected in genomes of 15 land plants,
409 seven streptophyte algae, and five chlorophytes; additionally, we included sequences found in
410 the transcriptomes of the Zygnematophyceae *Spirogyra pratensis* (de Vries et al., 2020),
411 *Zygnema circumcarinatum* (de Vries et al., 2018), and the Coleochaetophyceae *Coleochaete*
412 *orbicularis* (Ju et al., 2015; Figure 4).

413 The CCR homologs were distributed over several major clades. This included the
414 CCRL/DRL-like sequences described as TETRAKETIDE α -PYRONE REDUCTASE
415 (TKPR) by Grienberger et al. (2010), which is an important enzyme acting in the

416 production of sporopollenin. We recovered a well-supported clade of TKPR1 homologs
417 (bootstrap support of 84) and fully supported clade of TKPR2 homologs. Both clades of
418 TKPRs contained homologs from across the diversity of land plants, bolstering the idea that
419 TKPR1 and TKPR2 split early during plant evolution (Grienenberger et al., 2010)—before
420 the most recent common ancestor of land plants came about. Our domain structure analyses
421 found that TKPR1 possessed the NAD-dependent epimerase/dehydratase (IPR001509) and
422 NAD(P)-binding domain superfamily (IPR036291) domains, which appear to be present in
423 most sequences included in the phylogeny as well as the Tetraketide alpha-pyrone reductase
424 1 (IPR033267) domain (Figure S4). In contrast, TKPR2 only encoded the first two domains,
425 which is more similar to what is found in the CCR clade.

426 The *bona fide* CCRs and CCR-likes were spread out over two clades. These two
427 clades were nested in a weakly-supported monophylum (bootstrap support of 53), which was
428 sub-divided into four medium to fully supported clades. A fully supported clade of CCR-likes
429 (including AT4G30470 and AT2G23910) included sequences from across embryophytes; we
430 coined this monophylum CRL-A. *AtCCRL1* and *AtCCRL2* appear to be co-orthologs to one
431 sequence in the Brassicaceae *Capsella grandiflora* (0380s0077.1.p), thus our data suggest a
432 limited distribution of direct orthologs to CCRL1 and CCRL2. Yet the two sequences
433 together fall into a large clade, here coined CRL-B (bootstrap support 73), that contained
434 sequences from all major lineages of tracheophytes. Another medium-supported clade
435 (bootstrap support 73) included the *bona fide* CCRs, CCR1 and CCR2, of *Arabidopsis*. The
436 duplication that resulted in these two CCRs occurred earliest in the common ancestor of all
437 rosids and latest in the common ancestor of Brassicaceae, yet the CCR clade included
438 homologs from across tracheophytes. Many of these lineages appear to have expanded their
439 own repertoire from one CCR1/2 homolog that was present in the last common ancestor of
440 tracheophytes. Interestingly, a clade of divergent monocot CCRs display several
441 replacements in key amino acids involved in the binding of the substrate's phenolic ring, in
442 particular a replacement of non-polar aliphatic Ile to aromatic Tyr/Phe (Figure S5), which
443 might reflect a difference in substrate affinity. The catalytic triad SYK (Figure S5; Pan et al.
444 2014) is required for enzymatic activity and is overall conserved across CCR/DFR-likes and
445 FLDHs. Interestingly, several CRL-As possess non-conservative amino acid replacements
446 from large phenolic (Tyr/Phe) to smaller (His, Leu, Ser, Gly) amino acids, which might
447 suggest divergent substrate affinities for CRL-As. This is consistent with the domain
448 structure of many CRL-A sequences, which often lack the NAD-dependent
449 epimerase/dehydratase (IPR001509) domain, but possess additional domains such as 3-β-

450 hydroxysteroid dehydrogenase/isomerase (IPR002225) or match an additional NAD(P)-
451 binding domain (IPR016040; Figure S4). This pattern is only occasionally occurring in
452 sequences from the CCR or other CRL clades.

453 Finally, there is a third clade with a bootstrap support of 69 that included tracheophyte
454 sequences (forming a sub-clade with a bootstrap value of 90) and a single *Marchantia*
455 *polymorpha* homolog; we coined this clade CRL-C. Altogether, this suggests that the LCA of
456 all land plants had two homologs of CCRs/CRLs: one CRL-A and one CRL-B, CRL-C or
457 CCR homolog. In vascular plants, duplications have resulted in sub-clades of the CRL-
458 B/CRL-C/CCR homologs.

459 Within the larger clade that encases the DFRs, DFRLs, CCRs and CCRLs
460 (“Chloroplastida CCR/DFR-like”; bootstrap 86), one supported clade of Zygnematophyceae
461 (bootstrap support 77) and one supported clade of chlorophyte and streptophyte algae
462 (bootstrap support 76) exists. This points to a distinct DFR/DFRL/CCR/CCRL clade that
463 arose in the ancestor of Zygnematophyceae, yet its placement within the phylogeny other
464 than it belonging to the larger DFR/DFRL/CCR/CCRL clade is uncertain. The divergent
465 pattern of amino acids, which perform substrate and cofactor binding in land plants, suggest
466 that these algal homologs might vary in the substrate and enzymatic activity compared to
467 plant CCR/DFRs (Figure S5) It is, however, certain that within land plants, a pronounced
468 radiation of CCRs occurred.

469

470 **CADs are present across the green lineage**

471 The second reduction step of the activated acyl-group of the phenylpropanoid backbone and
472 one of the last steps in lignin biosynthesis is the production of phenylpropanoid-derived
473 alcohols from the corresponding aldehydes. An example is the conversion of *p*-coumaroyl
474 aldehyde into *p*-coumaryl alcohol (Kim et al., 2004; Pan et al., 2014). The required reduction
475 is catalyzed by CAD (Figure 1), which is the rate determining enzyme by which lignin is
476 produced (Gross et al., 1973; Mansell et al., 1974). In *Arabidopsis thaliana*, there are at least
477 eleven enzymes belonging to the CAD family. Enzymes of the CAD family have been
478 divided into five major groups, of which group IV was described as monocot-specific
479 (Saballos et al., 2009). In our previous studies many of the chlorophyte and streptophyte
480 potential CAD homologs, identified mostly from transcriptomes and few genomes of algae,
481 were described as CAD-like or CAD group II/III-affiliated (de Vries et al. 2017; de Vries et
482 al. 2020). Sequences clustering with those of CAD group II have been characterized as
483 sinapyl alcohol dehydrogenase (SAD) or show predicted structural similarity to SAD

484 enzymes (Guo et al. 2010; de Vries et al. 2017). Additionally, some SADs appear involved
485 not in the synthesis of lignin but defense compounds such as lignans (Suzuki and Umezawa,
486 2007; Guo et al., 2010; Saleem et al., 2010; Barakate et al., 2011), and it may thus be that
487 CAD group II is functionally versatile.

488 Here, we used the 11 canonical CAD sequences to understand the diversity in CAD
489 homologs across streptophytes. This includes also CAD group II sequences, for which
490 homologs in other species may have other substrate specificities and thus are involved in
491 different steps of the phenylpropanoid pathway (Barakat et al., 2009; Guo et al., 2010). We
492 computed a phylogeny of CAD homologs (Figure 5) detected in phylodiverse Chloroplastida.
493 While the resolution of the backbone is weak, we recovered all five CAD groups defined by
494 Saballos et al. (2009). All the CAD groups were resolved as land plant-specific clades of
495 CAD homologs with robust support. Each clade contained a varying set of major land plant
496 lineages (described below); the clades of putative streptophyte algal CAD homologs
497 contained both fewer proteins and fall in-between the five CAD-groups. Hence, this more
498 phylodiverse dataset tells a more complicated evolutionary history for CAD homologs than
499 the less-diverse data from de Vries et al. (2017).

500 Our data suggest that the common ancestor of Zygnematophyceae and land plants
501 may have possessed two CAD-like genes, which was followed by lineage-specific radiations.
502 While it appears—based on the overall topology of the tree—appealing to suggest that one
503 gene gave rise to CAD-group V and the other ancestral gene was the basis of CAD-group I to
504 IV, the low statistical support for the backbone of the phylogeny does not allow to confirm
505 such hypothesis (Figure 5). We can infer that the earliest land plants likely inherited a few (or
506 just one) CAD homolog from their algal progenitors. Most of the radiation of CADs has
507 occurred in plants dwelling on land. Of the canonical CAD group—and based on those part
508 of the topology with good bootstrap support—CAD-group V (containing *AtCAD1*) is the
509 only group present in all major land plant lineages (Figure 5). CAD-group I was likely
510 present in the LCA of tracheophytes, as it includes sequences from angiosperms,
511 gymnosperms, ferns, and the lycophyte *Selaginella moellendorffii*. CAD-groups II, III, and
512 IV include only angiosperm sequences—but note that with very weak bootstrap support (53)
513 a sequence from the gymnosperm *Gnetum montanum* associates with the group
514 CADII/III/IV; likely, expansion resulted in the ancestral gene of CAD-group II and III, which
515 diverged into CAD group III genes in angiosperms and after another expansion CAD-group
516 II originated in the LCA of dicots. This is in contrast to de Vries et al. (2017), where the
517 streptophyte algal CAD-like sequences were clustering with CAD group II/III sequences, but

518 resembles the placement of transcriptomic CAD-like sequences from *Spirogyra pratensis* and
519 *Mougotia* sp. in an already more diverse phylogenetic analysis (de Vries et al., 2020). This is
520 a clear case where including a larger diversity of streptophyte sequences to the analysis
521 enables us to better understand the complexity of the evolution of highly radiated gene
522 families. An analysis of the residues salient to CAD function showed a general conservation
523 of residues involved in zinc and NADP⁺ binding across Chloroplastida CAD-like sequences
524 (cf. Youn et al. 2006), whereas the residues in the binding pocket are generally less conserved
525 in these sequences (Figure 5). However, within a canonical CAD-groups or a CAD-like clade
526 we see a general conservation of residues in the binding pocket. The binding pockets of each
527 CAD-like clade appear different than those of the canonical CAD-groups. Domain analyses
528 suggest that some of the CAD-like sequences of streptophyte algae do not encode all of the
529 five domains present in most canonical CAD sequences, yet many of the
530 Zygnematophyceae CAD-like sequences encode all of these five domains. Similar patterns
531 emerge for other CAD-like sequences (Figure S6).

532 All CAD groups are shaped by multiple lineage-specific duplications and losses. This
533 hampers the inference of function and substrate specificity of the diverse CAD-like
534 sequences. Additionally, several lineages have originated a variety of CAD homologs that are
535 not yet designated to previous groups. In the absence of functional data, we however will not
536 give them a group designation but rather designate them as lineage-specific CAD-homologs
537 of unknown SAD or CAD function. CAD-like homologs found in streptophyte algae show
538 similar functional residues to other CAD-likes of land plants (not included in any of the five
539 clades of canonical CADs), including at the binding pockets—the pattern of residues is
540 similar to what is observed for other land plant sequences in-between the *bona fide* CAD
541 groups. The *bona fide* CAD groups showed a more homogeneous pattern of functional
542 residues.

543 Overall, both the topology of the phylogenetic tree and the conservation of key
544 residues point to (i) a deep evolutionary origin of CAD homologs and (ii) independent
545 radiations of CADs—not only in land plants but also in streptophyte algae.

546

547 **Massive independent radiations of acyltransferases and scattered candidates in** 548 **streptophyte algae**

549 A versatile group of enzymes that are important for the processes leading up to the lignins but
550 also compounds with antioxidant and antimicrobial properties are the BAHD acyltransferases
551 (named after the first enzymes characterized for this family BEAT, AHCT, HCBT, and DAT;

552 see also D’Auria, 2006). Most prominent among these are the versatile hydroxycinnamoyl-
553 CoA shikimate/quinate hydroxycinnamoyltransferases (HCT) (Eudes et al., 2016). Recently,
554 Kriegshauser et al. (2020) reported on the functional conservation of the HCT homologs
555 found in bryophytes with those of seed plants.

556 We computed a phylogeny including phylodiverse acyltransferases (Figure S7).
557 Overall, the topology of the tree corroborates the findings of Kriegshauser et al. (2020) that a
558 clear clade of HCT proteins first emerged in land plants—streptophyte algal sequences were
559 few, divergent from HCT, and scattered over the tree without clear affinity to characterized
560 acyltransferases. Without functional analyses, there is no solid foundation for predicting their
561 function—making them exciting candidates for future studies. What our data, however,
562 clearly reveal is that the last common ancestor of land plants likely had an expanded
563 repertoire of acyltransferases that further diversified during the radiation of plants on land.

564

565 **A clear clade of C3Hs is limited to land plants**

566 At several steps of the phenylpropanoid pathway, enzymes belonging to the cytochrome P450
567 family CYP98A, within the large CYP71 clan (Nelson and Werck-Reichhart, 2011), carry out
568 hydroxylations of *p*-coumarate-derived compounds (such as *p*-coumaroyl esters; for function
569 of C3H see Figure 1). In land plants, these hydroxylations are important for the production of
570 lignins, lignans, volatile phenylpropanoids, coumarins, and many more phenylpropanoid-
571 derived compounds. Previously, de Vries et al. (2017) described the detection of C3H in all
572 land plants and one putative C3H ortholog in *Klebsormidium nitens*. Now, with genomic gaps
573 in the streptophyte tree of life filled, we revisited the distribution of C3H.

574 The number C3H homologs detected in genomes of 15 land plants, 7 streptophyte
575 algae, and five chlorophytes and in the transcriptomes of *Spirogyra pratensis* (de Vries et al.,
576 2020), *Zygnema circumcarinatum* (de Vries et al., 2018), and *Coleochaete orbicularis* (Ju et
577 al., 2015) varied strongly between lineages. When sampling the sequences via BLAST (with
578 AT2G40890 as query sequence), we thus included either (a) all sequences that had a bit score
579 of at least 100, or (b) the top five hits. We aligned all sequences, cropped them to the
580 alignable region and computed a maximum likelihood phylogeny (Figure 6).

581 A clade of CYP98A included sequences from all major lineages of land plants. This
582 suggests that at least one CYP98A sequence was present in the LCA of all land plants. Based
583 on the lineages included here, it appears that from this single copy gene, radiations occurred
584 in the dicot lineages and *Amborella trichocarpa*. A single copy remained in the bryophytes,
585 lycophytes, ferns, gymnosperms, and monocots. The clade of CYP98A8 and CYP98A9 was

586 in our dataset limited to the Brassicaceae *Arabidopsis thaliana* and *Capsella grandiflora*
587 (bootstrap value 100). These two enzymes function in a route derived from the
588 phenylpropanoid pathway and are involved in the formation of N^1, N^5 -di(hydroxyferuloyl)-
589 N^{10} -sinapoylspermidine (Matsuno et al., 2009). The CYP98A8/9 clade falls into the larger
590 CYP98A clade (bootstrap support 100) together with the C3H sequences, suggesting that
591 they are the closest paralogs of C3H in *Arabidopsis thaliana* and *Capsella grandiflora*.
592 Analyses of their substrate recognition sites (SRS; cf. Rupasinghe et al. 2003) support the
593 divergent functional roles between CYP98A8/9 and canonical C3H. In particular, the first
594 two SRS (SRS1 and SRS2) show various distinct amino acid differences between the
595 CYP98A8/9 clade and the C3H clade (Figure S8). SRS1 and 2 are predicted to be involved in
596 binding of the substrate tails, which exhibit strong variation in size (Rupasinghe et al. 2003),
597 and thus may be critical for the substrate specificity of these paralogs. Indeed, SRS1 and 2 are
598 the two SRSs showing the strongest variation across the entire phylogeny, including also
599 other CYP450 subfamilies (Figure 6), corroborating this hypothesis.

600 We recovered additional land plant-specific clades of CYP450 enzymes (Figure 6)
601 such as one containing TRANSPARENT TESTA7 (TT7)-like sequences (bootstrap support
602 92); TT7 is a cytochrome P450 75B enzyme that is involved in flavonoid biosynthesis
603 (Tanaka et al., 1997; Schoenbohm et al., 2000). The BLAST search used for sampling C3H
604 homologs further recovered *AtCYP71B* and *AtCYP76C* members. In our phylogenetic
605 analysis, we inferred that *AtCYP71B34* and *AtCYP71B35* were likely born out of an
606 *Arabidopsis*-specific duplication, while CYP71B enzymes in general are present across
607 angiosperms. In contrast, *AtCYP76C1* and *AtCYP76C4* appear to have originated prior to the
608 split of *Arabidopsis thaliana* and *Capsella grandiflora*; a CYP76C4 ortholog appears to have
609 been lost in the latter plant species. The CYP76C clade may also be represented in other
610 species outside of angiosperms, because we find a sequence from *Gnetum montanum*
611 clustering with these sequences with a bootstrap support of 87. The rather long branch
612 warrants attention and would require a more CYP76C focused phylogenetic analysis, which
613 is not the point of this paper.

614 The algal sequences showed strong divergence to the C3Hs, forming only a larger
615 (fully supported) streptophyte-specific clade with all the recovered and functionally diverse
616 CYP450 enzymes (Figure 6). Domain structures of these streptophyte algal sequences is the
617 same to what is observed for the C3H homologs, but is in general conserved across the
618 phylogeny independent of the CYP450 subfamily assignment (Figure S9). Only few scattered
619 exceptions occur. Given the low support of most of the tree backbone, the role(s) of the

620 streptophyte algal homologs detected here remains elusive. Most of the streptophyte algae
621 show independent radiations of their CYP450 enzymes complicating functional predictions
622 even further. Thus, while there are interesting CYP450 candidates in streptophyte algae, a
623 clear C3H clade likely first arose early during the evolution of embryophytes.

624

625 **Monoacylglycerol lipases: multiple early radiations, independent subfunctionalization,**
626 **and the origin of CSE**

627 The conversion of caffeoyl-5-*O*-shikimate to caffeic acid may be a step along the
628 biosynthetic routes that lead to the production of G- and S-lignins in certain vascular plants
629 (Figure 1). The enzyme responsible for this step is caffeoyl-5-*O*-shikimate esterase (CSE).
630 CSE converts caffeoyl-5-*O*-shikimate to caffeic acid and was hypothesized to act together
631 with 4CL/ACOS5 to circumvent the catalysis of caffeoyl-5-*O*-shikimate to caffeoyl-CoA via
632 HCT (Vanholme et al. 2013). The latter pathway was proposed for tobacco by Hoffmann et
633 al. (2003) and confirmed for *Arabidopsis thaliana in vitro* by Vanholme and colleagues
634 (2013). Yet, based on *cse* mutants in *Arabidopsis thaliana*, they suggested that synthesis of
635 caffeoyl-CoA is more likely to occur via CSE and 4CL/ACOS5 than directly from caffeoyl
636 shikimate by HCT *in planta*. That said, in the model grasses *Brachypodium distachyon* and
637 *Zea mays* no CSE orthologs are present and crude extracts from these species show little
638 signs for the characteristic esterase activity (Ha et al., 2016). On the other hand, non-vascular
639 plants such as the model system *Physcomitrium patens* possess homologs of these enzymes
640 (Renault et al., 2017a), which suggests a secondary loss of CSE in the respective monocots.
641 CSE belongs to the family of putative monoacylglycerol lipases (MAGL). MAGLs are found
642 across eukaryotes and functional analyses in human, yeast and *Arabidopsis* have shown that
643 they possess monoacylglycerol lipase activity (Labar et al., 2010; Aschauer et al., 2016; Kim
644 et al., 2016). In contrast to other MAGLs of *Arabidopsis thaliana*, AtCSE (MAGL3) was
645 found to exhibit no hydrolytic activity on monoacylglycerols (MAGs) as a substrate (Kim et
646 al., 2016)—which applies to other enzymes of *Arabidopsis thaliana* that belong to the family
647 of MAGLs, too. Indeed, out of the 16 MAGLs that Kim and colleagues (2016) tested, only
648 MAGL6 and 8 showed high activity on MAG as substrate. Given the functional diversity in
649 MAGLs (Kim et al. 2016) and the unequal distribution of caffeoyl-5-*O*-shikimate across
650 embryophytes, functional analyses are required to fully understand how easily MAGLs can
651 lose or gain their MAGL activity. Yet, phylogenetic analyses can pinpoint the diversity of the
652 family across the green lineage.

653 Here we use phylogenetic analysis to pinpoint the distributions of the diverse MAGL
654 families, including CSE across the green lineage. In total, we recovered all 16 MAGL
655 sequences of *Arabidopsis thaliana* in the similarity search; using maximum likelihood
656 phylogenetics, we recovered clades for all the 16 MAGLs; some MAGL clades are widely
657 distributed throughout streptophytes, while others appear to have originated in embryophytes,
658 where they have again undergone lineage-specific expansions. Essentially, we recovered two
659 large clades: one restricted to streptophytes, containing homologs of MAGL2, 4 and 13; the
660 other has representation in chlorophytes as well and includes homologs of MAGL1, 3, 5, 6, 7,
661 8, 9, 10, 11, 12, 14, 15 and 16 (Figure 7).

662 Focusing on the MAGL2/4/13 clade first, we observe that MAGL13 has
663 representation in angiosperms, gymnosperms and ferns, suggesting its origin to be in the
664 LCA of tracheophytes, while MAGL2 and 4 came from a duplication event before the split
665 between *Arabidopsis* and *Capsella*. However, MAGL2/4 orthologs are present in other
666 species including *Picea abies*, suggesting that the common ancestor of seed plants possessed
667 a *MAGL2/4-like* and a *MAGL13* gene. Forming a clade with MAGL2/4/13 are lycophyte,
668 bryophyte and streptophyte algal sequences, which in general branch in an order expected
669 based on their species phylogeny (although within species duplication events have occurred).
670 This suggests that already at the base of streptophytes a *MAGL2/4/13-like* gene was present.

671 In the second large clade that includes also chlorophyte sequences, we find the clade
672 containing the CSE/MAGL3 orthologs. This clade includes sequences from both vascular and
673 non-vascular plants, pointing to an origin of CSE in the last common ancestor of land plants.
674 This adds support for a secondary loss in those monocots without a CSE ortholog. Despite the
675 origin of CSE in the common ancestor of land plants and a clear CSE ortholog in
676 *Physcomitrium patens* (3c19_14430V3.1.p), the substrate of CSE, caffeoyl-5-*O*-shikimate,
677 was not detected in crude extracts of the moss (Renault et al., 2017a). The HCT-based
678 reaction leading to caffeoyl-CoA has been confirmed *in vitro* using moss HCT (Kriegshauser
679 et al. 2021). Hence, the CSE homologs of *Physcomitrium patens* may have another function.
680 Indeed, the atypical function of *AtCSE*, together with the lack of other MAGL family
681 members, to act on MAGs (Kim et al. 2016), suggests that the functional spectrum of the
682 MAGL family is not very limited in land plants.

683 Members of the MAGL family share several conserved motifs across diverse
684 eukaryotes. One such motif (amino acid positions 132-141 in MAGL6 and 167-176 in
685 CSE/MAGL3) is situated in a region likely involved in substrate binding based on the crystal
686 structure of human MAGL (Labar et al., 2010). Within this motif a leucine in position

687 number four is found in a diverse set of 249 mammal and Sauria MAGLs investigated here
688 (Figure 7, inset) and most plant MAGLs including MAGL6 and 8. This is followed by
689 another hydrophobic amino acid (I in mammals/Sauria and V or L in most plant MAGLs). It
690 is striking that exactly these highly conserved amino acids are changed to a phenylalanine
691 and a serine in the Arabidopsis CSE and some homologs from other species. These changes
692 from two very hydrophobic amino acids to an aromatic and a hydrophilic one could be one of
693 the reasons for a change in substrate specificity from a substrate with a hydrophobic acyl
694 chain to a more hydrophilic substrate with aromatic properties. Based on this hypothesis,
695 CSEs would be restricted to some of the members of this clade that, however, stem from
696 across the diversity of land plants.

697 MAGL1 appears to have originated prior to the split of angio- and gymnosperms,
698 while specific MAGL14 and 16 orthologs likely arose after the split of asterids and rosids,
699 but a MAGL14/16 ortholog was likely present in the LCA of land plants. MAGL15, like
700 MAGL1, originated prior to the split of gymno- and angiosperms, and MAGL5 appeared to
701 come from a duplication of MAGL15 later on possibly in the ancestors of dicots.
702 Interestingly, a *MAGL5/15-like* sequence was already encoded in the genome of the ancestor
703 of streptophytes. The same is true for a *MAGL6/7/8/9/10/11/12-like* gene, which similar to
704 the *MAGL5/15-like* genes shows independent paths of radiation in streptophyte algae and
705 land plants. MAGL6,7,8,10 and 11 are only present in the here included Brassicaceae, while
706 a *MAGL6/7/8/10/11-like* gene was already present in the common ancestor of
707 spermatophytes. The same evolutionary history describes the scenario under which MAGL9
708 and 12 originated.

709 Finding the MAGL6/7/8/10/11 subclade specifically expanded in seed plants is
710 noteworthy. At least *AtMAGL8* localizes to lipid droplets (Kim et al., 2016), which are
711 structures found in various photosynthetic eukaryotes but are well-known from seeds. Thus,
712 expansion of this clade might be a read-out of spermatophyte-specific additions to the ancient
713 set of proteins relevant to LD formation and function (see de Vries and Ischebeck, 2020).

714 All in all, the MAGLs have experienced an early radiation in streptophytes. Given
715 that even the Arabidopsis MAGLs without detectable activity on MAG (see Kim et al., 2016)
716 do not form one single monophylum, it is conceivable that subfunctionalization of members
717 of the MAGL family occurred multiple times independently. This may likewise be true for all
718 independent expansions of MAGL-encoding genes observed in any other species included
719 here. The versatility in functional evolution of MAGLs makes it difficult to make robust
720 predictions of putative MAGL functions.

721

722 **COMT: convergence and complexity**

723 In angiosperms, ferulate 5-hydroxylase (F5H) and caffeate *O*-methyltransferase (COMT)
724 carry out important catalytic steps along the route from *p*-coumaroyl-CoA to S-lignin. COMT
725 catalyzes the methylation of caffeic acid or 5-hydroxyferulic acid, the product formed by
726 F5H. Like C3H and C4H, F5H belongs to the large CYP450 clan 71 (Nelson and Werck-
727 Reichhart, 2011). The function of F5H evolved at least twice in P450 enzymes, once in the
728 ancestor of angiosperms and once in the ancestor of lycophytes (Weng et al., 2008; Weng and
729 Chapple, 2010).

730 For angiosperm F5H, no clear putative orthologs were found outside of flowering
731 plants and likewise no clear orthologs were found for the lycophyte F5H (i.e. “SmF5H”),
732 which forms a separate clade from the angiosperm F5H sequences (Figure S10). This
733 corroborates previous results (de Vries et al., 2017) and is in agreement with the hypothesis
734 that F5H function evolved at least twice in the evolution of land plants (Weng et al., 2008).
735 Additionally, the average pairwise identity of the F5H homologs was quite low (18.5%)—
736 hampering robust phylogenetic analyses. We thus did not further delve into the evolution of
737 F5H. COMT however caught our attention.

738 The lycophyte *Selaginella moellendorffii* not only uses a genetically distant F5H
739 enzyme; the same appears to be true for COMT (Weng et al., 2011). This highlights a
740 promiscuity for substrate specificity and activity in P450 enzymes that is yet to be discovered
741 and mere orthology analyses can only go so far as to discover putative candidates. Using
742 phylogenetics, we explored the diversity of methyltransferases by screening for sequences
743 homologous to COMT/OMT1 of *Arabidopsis thaliana* across our phylodiverse dataset. This
744 approach identified not only clear orthologs but can also serve as a backbone to map relevant
745 mutations facilitating in functionally convergence in this group of enzymes and by that may
746 highlight possible candidates for *in vivo* and *in vitro* studies.

747 For most clades of land plant methyltransferases, based on the here recovered
748 topology, predicting a putative function was not straightforward. This applied even more so
749 to the homologs of COMT/OMT found in chlorophyte and streptophyte algae. We recovered
750 a clade of methyltransferases that included chlorophyte and streptophyte algae as well as
751 diverse land plant sequences (coined ‘Chloroplastida OMT’ in Figure 8); among these
752 clustered *Arabidopsis thaliana* proteins such as COMT, indole glucosinolate
753 methyltransferases (IGMT), and nicotinate *N*-methyltransferase (NANMT; see Li et al.,
754 2017)—hence different methyltransferases that act on a range of aromatic compounds. What

755 this means for the presence of a putative COMT in algae is obscure. However, it corroborates
756 the previously observed patchy detection of OMT1 across the green lineage based on a
757 reciprocal BLASTp searches (de Vries et al., 2017). Clear orthologs of *AtCOMT* were only
758 detected for a few angiosperms, notably not including any monocot sequence that we
759 included in our dataset. Of all methyltransferases in our dataset, only NANMT formed a
760 clade of clear orthologs that included more than one major lineage of land plants by encasing
761 sequences from angiosperms and *Picea abies* (bootstrap support 87). All other orthogroups
762 appear, like COMT, to be restricted to only a few of the included angiosperm lineages.

763 The lycophyte *Selaginella moellendorffii* has a COMT that is distantly related to
764 OMT of angiosperms. It appears to have acquired its OMT activity through convergent
765 evolution and was coined *SmCOMT* (Weng et al., 2011). In agreement with this, *SmCOMT*
766 did not cluster with the *AtOMT1* sequence in our analyses. Instead, it forms its own (weak)
767 clade with only one other sequence from *Selaginella moellendorffii* (bootstrap support 63).
768 The other *SmCOMT*-like sequences (described in Weng et al., 2011) were distributed over
769 the phylogeny and appear to be specific to *Selaginella moellendorffii*. Nonetheless, this
770 pattern highlights a certain versatility in the evolutionary history of substrate specificity of *O*-
771 methyltransferases in land plants. This appears to be only logical, noting the large lineage-
772 specific expansions in the larger clade of *O*-methyltransferases that encompasses all *O*-
773 methyltransferases from *Arabidopsis thaliana*—that is COMT, IGMTs, NANMT and *N*-
774 acetylserotonin *O*-methyltransferase (ASMT; for more on this enzyme see, e.g., Tan et al.,
775 2012; Byeon et al., 2016) (bootstrap-support 99). Within this clade fall also algal sequences
776 from chlorophytes and streptophyte algae. Their position within the clade is undetermined
777 due to low bootstrap support. These sequences appear highly divergent, many of them cluster
778 with rather long branches. Yet, some of the sequences from our previous analysis found a
779 reciprocal BLASTp hit to *AtCOMT*, including *Klebesormidium nitens* 00158_0100v1.1 that
780 clusters in a fully supported clade of *Klebesormidium* paralogs; these sequences are
781 promising candidates to explore caffeic acid *O*-methyltransferase activity. Indeed, when we
782 modelled the tertiary structure of *Klebesormidium nitens* 00158_0100v1.1 and
783 *Mesotaenium endlicherianum* ME000591S08520 using I-TASSER (Zhang, 2008), we
784 recovered *Medicago sativa* and *Lolium perenne* COMT as its closest structural analogs
785 (1KYZ; Zubieta et al., 2002; 3P9C; Louie et al., 2010; TM-scores 0.865 and 0.956,
786 respectively).

787 Like the land plant COMTs, also the algal COMT-likes appear to have undergone
788 independent radiations in this large gene family. Given the observed convergent evolution of

789 COMT activity in *Selaginella moellendorffii*, the question of whether there is COMT activity
790 across Streptophyta remains wide open.

791 To gain a first insight into whether COMT activity can be expected from other
792 streptophyte lineages, we investigated the conservation of residues relevant for the function
793 of COMT, including those that form the substrate binding pocket (Figure 8). The functional
794 residues were identified from COMT of *Lolium perenne* (Louie et al., 2010). Across our
795 phylogeny, these residues differ between the clades of canonical ASMT, NANMT, COMT
796 and IGMT, while they are conserved within them (Figure 8). The binding pocket of *ArOMT1*
797 and its orthologs consist of the amino acid pattern MSNGGG, whereas the pattern for the
798 residues important for the function of the enzyme is HDE. While HDE appears conserved
799 across the majority of sequences analyzed here, independent of the specific function of the
800 enzyme (e.g. ASMT, IGMTs and COMT all have the pattern HDE), the binding pocket is
801 highly variable among the functionally characterized enzymes. This suggests that the
802 reaction-determining residues are those that form the binding pocket and not those that are
803 catalytically important. This seems logical given that all these enzymes catalyze similar types
804 of reactions. The triple G in the binding pocket is also far more conserved across the entire
805 phylogeny, with only few exceptions occurring, while the first three residues are highly
806 variable. Indeed, the COMT-specific MSN motif is not present in the functionally
807 characterized COMT from *Selaginella moellendorffii*, rather it is MTN, which however is a
808 change between similar amino acids (Ser to Thr). Apart from the COMT orthologs and
809 *SmCOMT* sequence no other sequences from any other lineage encode the binding pocket
810 pattern M(S/T)NGGG, suggesting that none has a canonical preference for binding 5-
811 hydroxyconiferaldehyde. Yet, several homologs—including those of streptophyte algae—
812 would have the ability to catalyze the reaction based on the conservation of the residue
813 pattern HDE. What is more, all the sequences that could not be properly identified as
814 orthologs to ASMT, NANMT, COMT and IGMT (with the exception of *SmCOMT*) show no
815 similarity in their first three residues of the binding pocket to either of these enzyme families.
816 This would suggest that the enzymes from most streptophyte lineages included in this
817 analysis use different substrates than those functionally characterized in *Arabidopsis*
818 *thaliana*.

819

820 **CCoAOMT-like sequences emerged in Phragmoplastophyta**

821 Within the phenylpropanoid pathway, caffeoyl-CoA *O*-methyltransferase (CCoAOMT) and
822 most of its homologs are the enzymes that catalyze the first committed step to many of at

823 least two types of lignin (S- and G-lignin). These enzymes methylate caffeoyl-CoA and thus
824 give rise to feruloyl-CoA (Ye et al., 1994; Ye and Varner, 1995; Martz et al., 1998; Do et al.,
825 2007; Vanholme et al., 2012); in the past, it was also proposed that after the conversion to 5'-
826 hydroxy-feruloyl-CoA, CCoAOMT can methylate this compound to produce Sinapoyl-CoA
827 (Maury et al., 1999; Ferrer et al., 2005). One of the paralogs that exist in *Arabidopsis*
828 *thaliana* (tapetum-specific *O*-methyltransferase [TSM1], AT1G67990) however, shows
829 activity towards a coniferyl derivative that is formed at N^{10} by F5H starting from N^1, N^5, N^{10} -
830 tris-(hydroxyferuloyl) spermidine (Fellenberg et al., 2009); TSM1 catalyzes the production of
831 N^1, N^5 -bis-(hydroxyferuloyl)- N^{10} -synapoylspermidine (Fellenberg et al., 2009). Thus,
832 CCoAOMTs appear to be versatile in their substrate specificity and can act on different steps
833 in the phenylpropanoid pathway.

834 Here we used phylogenetics to disentangle the distribution of these enzymes across
835 the green lineage. A duplication gave rise to the genes encoding the functionally divergent
836 enzymes *AtCCoAMT* and *AtTSM1* (bootstrap support 77; Figure 9). These two
837 methyltransferases are embedded in a larger clade containing the other CCoAOMT enzymes
838 *AtCCoAOMT1*, *AtCCoAOMT7* and *AtCCoAOMT-like* (AT1G24735). The latter appears to
839 be specific to the Brassicaceae included in this dataset, while homologs of *AtCCoAOMT7*
840 occur across dicots and were detected in *Amborella trichopoda*, but were absent from the
841 included monocots. Only *AtCCoAOMT1* had a wider distribution. Its cluster (bootstrap
842 support 79) contains angiosperms, gymnosperms, lycophytes, ferns and bryophytes,
843 excluding the sequenced hornworts from the genus *Anthoceros* (Figure 9). Assuming a
844 monophyly of Bryophyta (Puttick et al., 2018), this suggests a loss of CCoAOMT1 in at least
845 the sequenced *Anthoceros* species, and that CCoAOMT1 was present in the common ancestor
846 of land plants. Lineage-specific duplications of CCoAOMT1 appear to have happened,
847 indicated by the expansions seen in tobacco, spruce, the lycophyte *Selaginella moellendorffii*
848 and the water fern *Azolla filiculoides*. The expanded repertoire of sequences in monocots and
849 *Gnetum* indicate additional lineage-specific duplications outside of the CCoAOMT clade.
850 The case of TSM1 suggests that neo-functionalization can easily occur within this type of
851 methyltransferases. We noted that the residues involved in substrate binding (Ferrer et al.
852 2005) are identical in *AtCCoAMT* and *AtTSM1* (Figure 9). A possible explanation might be
853 that the make-up of binding pocket allows for a certain versatility in substrates. Given these
854 observations, the paralogs within this and the other CCoAOMT clades cannot be assumed to
855 hold the function of CCoAOMTs. Likewise, it cannot be ruled out that their LCA may have
856 had this function. As sister to the methyltransferase clade, including the CCoAOMT

857 homologs lies a cluster of genes encoding putative candidates for streptophyte algal
858 CCoAOMTs. These were limited to representatives of the two streptophyte algal lineages
859 closest to land plants: the Coleochaetophyceae *Coleochaete scutata* and *Coleochaete*
860 *orbicularis* as well as the Zygnematophyceae *Spirogloea muscicola*. These algal sequences
861 have the same domain structure as the majority of all *S*-adenosyl-*L*-methionine (SAM)-
862 dependent methyltransferases included in the phylogeny (Figure S11). Only single sequences,
863 scattered across the phylogeny and diversity of species included here, vary in their domain
864 structure showing a loss of a domain loss, a gain of an additional domain, or both. The
865 analyses of the specific functional residues gave more insights into the streptophyte algal
866 sequences within the clade of SAM-dependent methyltransferases in the peripheral routes of
867 the phenylpropanoid pathway. These algal sequences maintain the residues involved in ion
868 and cofactor binding, but differ strongly in the substrate binding site (cf. Ferrer et al. 2005;
869 Figure 9).

870 Altogether, it appears that the family of CCoAOMT-like proteins has its origin in
871 Phragmoplastophyta. Clarifying the function of the putative CCoAOMT-like enzymes in
872 Coleochaetophyceae and Zygnematophyceae has the potential to shed light on a
873 synapomorphy with physiological relevance.

874

875 **Conclusion**

876 All genes for enzymes that act in early steps in the chassis of the phenylpropanoid pathway
877 investigated here (Figure 1) can be traced back to the LCA of land plants with the exception
878 of COMT (Figure 10); most of these can even be traced back to some ancestor that land
879 plants shared with streptophyte algae. While most of our knowledge on how these genes
880 work comes from angiosperms, this does not capture the sequence diversity in enzymes—and
881 it underpins the versatility in producing specialized metabolites.

882 Our data pinpoint that most of the enzymes have undergone massive lineage-specific
883 expansions. A lineage-specific expansion is palpable even despite the fact that sampling of
884 sequences across the Streptophyta is still strongly biased towards seed plants. These data
885 offer a framework for pinpointing those candidate genes/enzymes that are bound to shed light
886 on the evolution of key enzymatic steps—and novel ones. Such work is exemplified by
887 studies on *Selaginella* or bryophyte model systems such as *Physcomitrium patens*.
888 Characterizing enzymes that are even more divergent from what we know from angiosperms
889 should yield surprising insights and novel routes in this bountiful pathway.

890

891 MATERIAL AND METHODS

892 Dataset of protein sequences and screening for homologs

893 We downloaded protein data from: (a) genomes of fifteen land plants: *Anthoceros agrestis* as
894 well as *Anthoceros punctatus* (Li et al., 2020), *Amborella trichopoda* (*Amborella* Genome
895 Project, 2013), *Arabidopsis thaliana* (Lamesch et al., 2010), *Azolla filiculoides* (Li et al.,
896 2018), *Brachypodium distachyon* (The International Brachypodium Initiative, 2010),
897 *Capsella grandiflora* (Slotte et al., 2013), *Gnetum montanum* (Wan et al., 2018), *Marchantia*
898 *polymorpha* (Bowman et al., 2017), *Nicotiana tabacum* (Sierro et al., 2014), *Oryza sativa*
899 (Ouyang et al., 2007), *Picea abies* (Nystedt et al., 2013), *Physcomitrium patens* (Lang et al.,
900 2018), *Salvinia cucullata* (Li et al., 2018), *Selaginella moellendorffii* (Banks et al., 2011), and
901 *Theobroma cacao* (Argout et al., 2011); (b) the genomes of seven streptophyte algae:
902 *Chlorokybus atmophyticus* (Wang et al., 2020), *Chara braunii* (Nishiyama et al., 2018),
903 *Klebsormidium nitens* (Hori et al., 2014), *Mesotaenium endlicherianum* (Cheng et al., 2019),
904 *Mesostigma viride* (Wang et al., 2020), *Penium margaritaceum* (Jiao et al., 2020), *Spirogloea*
905 *musciicola* (Cheng et al., 2019); (c) the genomes of five chlorophytes: *Bathycoccus prasinos*
906 (Moreau et al., 2012), *Chlamydomonas reinhardtii* (Merchant et al., 2007), *Coccomyxa*
907 *subellipsoidea* (Blanc et al., 2012), *Micromonas pusilla*, *Micromonas* sp. (Worden et al.,
908 2009), *Ostreococcus lucimarinus* (Palenik et al., 2007), *Ulva mutabilis* (De Clerck et al.,
909 2018), *Volvox carteri* (Prochnik et al., 2010). Additionally, we included sequences found in
910 the transcriptomes of *Spirogyra pratensis* (de Vries et al., 2020), *Zygnema circumcarinatum*
911 (de Vries et al., 2018), and *Coleochaete orbicularis* (Ju et al., 2015).

912 For each of the protein families we investigated here, the representative *Arabidopsis*
913 *thaliana* protein was used as a query sequence for a BLASTp against this dataset. Initially,
914 we considered all homologs recovered at a cutoff level of 10^{-7} . However, due to the large size
915 of the protein families (i.e. high number of well-supported homologs obtained), refinement of
916 the datasets was carried out as described in the individual sections for these enzymes in the
917 Results and Discussion section.

918

919 Alignments, phylogenetic analysis, and primary sequence analysis

920 Using the homologs detected based on the above described BLASTp search for a given
921 enzyme, we generated alignments using MAFFT v7.453 (Katoh and Standley, 2013) with a
922 L-INS-I approach. Alignments were cropped, if necessary, to retain conserved domains that
923 were alignable for all homologs; alignments are provided in Supplemental Datasets S1 to
924 S11. We computed maximum likelihood phylogenies using IQ-TREE multicore version 1.5.5

925 (Nguyen et al., 2015), with 100 bootstrap replicates. To determine the best model, we used
926 ModelFinder (Kalyaanamoorthy et al., 2017) and picked the best models based on the
927 Bayesian Information Criterion. The best models were: LG+G4 (Le and Gascuel, 2008) for
928 4CL, CCR, CCoAOMT; LG+I+G4 for PAL, CAD, MAGL/CSE, COMT, and for the
929 preliminary phylogeny of 4CL; LG+F+I+G4 for C4H, F5H, and C3H; WAG+F+G4 (Whelan
930 and Goldman, 2001) for HCT.

931 Protein structure prediction was carried out using the sequences as input in the online
932 Iterative Threading ASSEMBLY Refinement (I-TASSER; Zhang, 2008; Yang et al., 2015).
933 Functional residue analyses were based on published structural analyses (Rupasinghe et al.,
934 2003; Ferrer et al., 2005, Youn et al., 2006; Hu et al., 2010; Pan et al., 2014) and alignments
935 were viewed with SeaView v.4 (Gouy et al., 2009) and plotted with ETE3 (Huerta-Cepas et
936 al., 2016).

937

938 **Protein domain predictions**

939 Protein domains for all protein sequences for the enzyme families 4CL, CCR, C3H, CAD and
940 CCoAOMT included in the phylogenies were predicted using InterProScan version 5.47-82.0
941 (Jones et al., 2014). The presence or absence of protein domains were mapped onto the
942 phylogenies of the afore mentioned gene families as presence/absence heatmaps, which were
943 visualized using iTOL v6 (Letunic and Bork, 2019).

944

945 **ACKNOWLEDGEMENTS**

946 J.M.R.F.-J. is grateful for being supported by the Ph.D. program "Microbiology and
947 Biochemistry" within the framework of the "Göttingen Graduate Center for Neurosciences,
948 Biophysics, and Molecular Biosciences" (GGNB) at the University of Goettingen; A.D.A. is
949 grateful for being supported through the International Max Planck Research School (IMPRS)
950 for Genome Science. J.d.V. thanks the European Research Council for funding under the
951 European Union's Horizon 2020 research and innovation programme (Grant Agreement No.
952 852725; ERC-StG "TerreStriAL"). M.P., I.F., and J.d.V. are grateful for support through the
953 German Research Foundation (DFG) within the framework of the Priority Programme
954 "MAdLand – Molecular Adaptation to Land: Plant Evolution to Change" (SPP 2237; VR
955 132/4-1; PE 360/37-1; FE 446/14-1)

956

957 **FIGURE LEGENDS**

958 **Figure 1. Enzymes involved in the biosynthesis of phenylpropanoid-derived compounds**
959 **investigated here.** A simplified schematic of the phenylpropanoid pathway and its routes to
960 different derivatives are shown. Boxes indicate enzyme families, which are mentioned above
961 each box and color-coded. Their coloration is the same as in Figure 10. Dotted lines indicate
962 putative steps in the pathway.

963

964 **Figure 2. A phylogenetic framework for the origin of streptophyte PAL.** PAL and HAL
965 homologs were screened for in fifteen land plant, seven streptophyte algae, and five
966 chlorophytes. Among Chloroplastida, PAL homologs were only recovered from genomes of
967 land plants and the streptophyte algae *Klebsormidium nitens* and *Chara braunii*. From all
968 detected homologs, a rooted maximum likelihood phylogeny was computed using LG+I+G4
969 as model for protein evolution (chosen according to BIC). 100 bootstrap replicates were
970 computed; only bootstrap values ≥ 50 are shown and bootstrap values of 100 are depicted by a
971 filled dot. Colored font and dots correspond to the support recovered for the higher-order
972 clades labeled on the right of the phylogenies.

973

974 **Figure 3. 4CL homologs occur across Streptophyta.** 4CL homologs were sampled from
975 protein data of nine land plant, seven streptophyte algal and five chlorophyte algal genomes.
976 Only protein sequences with a minimum length of 400 and a maximum length of 1150 amino
977 acids were included. From all detected homologs, a rooted maximum likelihood phylogeny
978 was computed using LG+G4 as model for protein evolution (chosen according to BIC). 100
979 bootstrap replicates were computed; only bootstrap values ≥ 50 are shown and bootstrap
980 values of 100 are depicted by a filled dot. Colored font and dots correspond to the support
981 recovered for the higher-order clades labeled on the right of the phylogenies. On the right we
982 show key residues for substrate binding and function of canonical 4CL as reported by Hu et
983 al. (2010).

984

985 **Figure 4. The complex evolutionary history of CCR in Chloroplastida.** CCR homologs of
986 a minimum of 220 amino acids were sampled from protein data of 15 land plant, seven
987 streptophyte algal, and five chlorophyte algal genomes; additionally, we included sequences
988 found in the transcriptomes of *Spirogyra pratensis* (de Vries et al., 2020), *Zygnema*
989 *circumcarinatum* (de Vries et al., 2018), and *Coleochaete orbicularis* (Ju et al., 2015). From
990 all detected homologs, an unrooted maximum likelihood phylogeny was computed using
991 LG+G4 as model for protein evolution (chosen according to BIC). 100 bootstrap replicates

992 were computed; only bootstrap values ≥ 50 are shown and bootstrap values of 100 are
993 depicted by a filled dot. Colored font and dots correspond to the support recovered for the
994 higher-order clades labeled on the right of the phylogenies.

995

996 **Figure 5. Phylogenetic analysis highlights CAD candidates across Chloroplastida.** CAD
997 homologs were sampled from protein data from fifteen land plant, seven streptophyte algal,
998 and five chlorophyte algal genomes as well as sequences found in the transcriptomes of
999 *Spirogyra pratensis* (de Vries et al., 2020), *Coleochaete scutata* and *Zygnema*
1000 *circumcarinatum* (de Vries et al., 2018), and *C. orbicularis* (Ju et al., 2015). From all
1001 detected homologs, an unrooted maximum likelihood phylogeny was computed using
1002 LG+I+G4 as model for protein evolution (chosen according to BIC). 100 bootstrap replicates
1003 were computed; only bootstrap values ≥ 50 are shown and bootstrap values of 100 are
1004 depicted by a filled dot. Colored font and dots correspond to the support recovered for the
1005 higher-order clades labeled on the right of the phylogenies. The five groups of CADs were
1006 named in accordance with Saballos et al. (2009). Next to the sequence labels residues from
1007 the binding pocket, NADP⁺- and Zn²⁺ binding are shown—based on Youn et al. (2016).

1008

1009 **Figure 6. A clade of C3H orthologs originated at the base of land plants.** C3H homologs
1010 were sampled from protein data of genomes of fifteen land plants, seven streptophyte algae,
1011 and five chlorophytes; additionally, sequences found in the transcriptomes of *Spirogyra*
1012 *pratensis* (de Vries et al., 2020), *Zygnema circumcarinatum* (de Vries et al., 2018), and
1013 *Coleochaete orbicularis* (Ju et al., 2015) were included. For downstream analyses, we used
1014 either (a) all sequences that had a bit score of at least 100 or (b) the top five hits. We aligned
1015 all sequences, cropped them to the alignable region and computed an unrooted maximum
1016 likelihood phylogeny was computed using LG+F+I+G4 as model for protein evolution
1017 (chosen according to BIC). 100 bootstrap replicates were computed; only bootstrap values
1018 ≥ 50 are shown and bootstrap values of 100 are depicted by a filled dot. Colored font and dots
1019 correspond to the support recovered for the higher-order clades labeled on the right of the
1020 phylogenies. Two large clades that contained only (a) *Anthoceros* and (b) chlorophyte and
1021 streptophyte algal sequences were collapsed; the full tree is shown in Figure S12.

1022

1023 **Figure 7. The occurrence of MAGLs across diverse Streptophyta and a phylogenetic**
1024 **framework for the deep evolutionary roots of CSE.** MAGL/CSE homologs were sampled
1025 from protein data from 15 land plant, seven streptophyte algal, and five chlorophyte algal

1026 genomes as well as sequences found in the transcriptomes of *Spirogyra pratensis* (de Vries et
1027 al., 2020), *Coleochaete scutata* and *Zygnema circumcarinatum* (de Vries et al., 2018), and
1028 *Coleochaete orbicularis* (Ju et al., 2015). From all detected homologs, an unrooted maximum
1029 likelihood phylogeny was computed using LG+I+G4 as model for protein evolution (chosen
1030 according to BIC). 100 bootstrap replicates were computed; only bootstrap values ≥ 50 are
1031 shown and bootstrap values of 100 are depicted by a filled dot. Colored font and dots
1032 correspond to the support recovered for the higher-order clades labeled on the right of the
1033 phylogenies. Purple font highlights those streptophyte algal sequences that share the
1034 conserved alpha helix cap domain with CSE. Logos are based on a motif (amino acids 132-
1035 141 in MAGL6 and 167-176 in CSE/MAGL3) that is situated in a region likely involved in
1036 substrate binding based on the crystal structure of human MAGL.

1037

1038 **Figure 8. Low resolution on the complex evolutionary history of COMT.** We explored
1039 the diversity of methyltransferases by screening for sequences homologous to Arabidopsis
1040 COMT/OMT1 across genome data from fifteen land plant, seven streptophyte algae, and five
1041 chlorophytes; additionally, we included sequences found in the transcriptomes of *Spirogyra*
1042 *pratensis* (de Vries et al., 2020), *Zygnema circumcarinatum* (de Vries et al., 2018), and
1043 *Coleochaete orbicularis* (Ju et al., 2015). From all detected homologs, an unrooted maximum
1044 likelihood phylogeny of 226 sequences was computed using LG+I+G4 as model for protein
1045 evolution (chosen according to BIC). 100 bootstrap replicates were computed; only bootstrap
1046 values ≥ 50 are shown and bootstrap values of 100 are depicted by a filled dot. Colored font
1047 and dots correspond to the support recovered for the higher-order clades labeled on the right
1048 of the phylogenies. Blue font highlights streptophyte algal sequences; bold font pinpoints
1049 those, that recovered land plant COMT as closest structural analogs in I-TASSER-based
1050 modeling. On the right we show residues important for substrate binding and function of
1051 canonical COMT as reported by Louie et al. (2010).

1052

1053 **Figure 9. A phylogenetic framework for the evolutionary origin of CCoAOMTs in**
1054 **Phragmoplastophyta.** CCoAOMT homologs were sampled from protein data from sixteen
1055 land plant, seven streptophyte algal, and five chlorophyte algal genomes as well as sequences
1056 found in the transcriptomes of *Spirogyra pratensis* (de Vries et al., 2020),
1057 *Coleochaete scutata* and *Zygnema circumcarinatum* (de Vries et al., 2018), and *Coleochaete*
1058 *orbicularis* (Ju et al., 2015). From all detected homologs, an unrooted maximum likelihood
1059 phylogeny of 138 sequences was computed using LG+G4 as model for protein evolution

1060 (chosen according to BIC). 100 bootstrap replicates were computed; only bootstrap values
1061 ≥ 50 are shown and bootstrap values of 100 are depicted by a filled dot. Colored font and dots
1062 correspond to the support recovered for the higher-order clades labeled on the right of the
1063 phylogenies. Purple font highlights those streptophyte algal sequences that share the
1064 conserved alpha helix cap domain with CSE. The alignment on the right shows functionally
1065 characterized sites involved in substrate, ion and co-factor recognition of CCoAOMT (Ferrer
1066 et al., 2005).

1067

1068 **Figure 10. A summary of the proposed evolutionary trajectory of key enzymes in the**
1069 **phenylpropanoid pathway across the green lineage.** At the bottom is a cladogram of the
1070 green lineage. The most recent common ancestors (MRCA) of Chloroplastida, Streptophyta,
1071 Phragmoplastophyta, Bryophyta, Embryophyta and Tracheophyta are indicated at their
1072 respective nodes. On top the cladogram is the proposed evolutionary trajectory of the enzyme
1073 families PAL, C4H, 4CL, CCR, CAD, C3H, CSE/MAGL (CSE), and CCoAOMT. The
1074 names of the enzyme families are indicated on the left of the trajectory. The enzyme (sub-
1075)families present in a specific common ancestor have been plotted onto the respective nodes
1076 of the cladogram below the evolutionary scenario of the enzyme families involved in the
1077 phenylpropanoid pathway and lignin biosynthesis. White dots indicate absence/loss of a gene
1078 family one dot indicates the presence of one representative of the gene family and several
1079 dots indicate an expansion (two or more members of the gene family) in at least one species
1080 of the represented lineages in the cladogram. Colors are chosen to distinguish different
1081 enzyme families and subfamilies. Question marks label sequences of ambiguous affiliation.

1082

1083 REFERENCES

- 1084 Amborella Genome Project (2013) The *Amborella* genome and the evolution of flowering
1085 plants. *Science*, 342, 1241089
- 1086 Argout, X., Salse, J., Aury, J.-M., Gultinan, M. J., Droc, G., Gouzy, J., ... Lanaud, C.
1087 (2011). The genome of *Theobroma cacao*. *Nature Genetics*, 43, 101–108.
- 1088 Aschauer, P., Rengachari, S., Lichtenegger, J., Schittmayer, M., Padmanabha Das, K. M.,
1089 Mayer, N., ... Oberer, M. (2016) Crystal structure of the *Saccharomyces cerevisiae*
1090 monoglyceride lipase Yju3p. *Biochimica et Biophysica Acta (BBA) - Molecular and Cell*
1091 *Biology of Lipids*, 1861, 462–470.

- 1092 Banks, J. A., Nishiyama, T., Hasebe, M., Bowman, J.L., Gribskov, N., dePamphilis, C., Albert,
1093 V. A., ... Grigoriev, I. V. (2011). The Selaginella genome identifies genetic changes
1094 associated with the evolution of vascular plants. *Science*, 332, 960-963.
- 1095 Barakat, A., Bagniewska-Zadworna, A., Choi, A. Plakkat, U., DiLoreto, D. S., Yellanki, P.,
1096 & Carlson, J. E. (2009). The cinnamyl alcohol dehydrogenase gene family in *Populus*:
1097 phylogeny, organization, and expression. *BMC Plant Biology* 9, 26.
- 1098 Barakate, A., Stephens, J., Goldie, A., Hunter, W.N., Marshall, D., Hancock, R.D., Lapierre,
1099 C., Morreel, K., Boerjan, W., & Halpin, C. (2011). Syringyl lignin is unaltered by severe
1100 sinapyl alcohol dehydrogenase suppression in tobacco. *The Plant Cell*, 23, 4492-4506.
- 1101 Barros, J., Serrani-Yarce, J.C., Chen, F., Baxter, D., Venables, B.J., & Dixon, R.A. (2016).
1102 Role of bifunctional ammonia-lyase in grass cell wall biosynthesis. *Nature Plants*, 2,
1103 16050.
- 1104 Barros, J., & Dixon, R.A. (2020). Plant Phenylalanine/Tyrosine Ammonia-lyases. *Trends in*
1105 *Plant Science*, 25, 66-79.
- 1106 Bednarek, P., Schneider, B., Svatoš, A., Oldham, N.J., Hahlbrock, K. (2005). Structural
1107 complexity, differential response to infection, and tissue specificity of indolic and
1108 phenylpropanoid secondary metabolisms in Arabidopsis roots. *Plant Physiology*, 138,
1109 1058–1070.
- 1110 Berens, M.L., Berry, H.M., Mine, A., Argueso, C.T., & Tsuda, K. (2017). Evolution of
1111 hormone signaling networks in plant defense. *Annual Review of Phytopathology*, 55,
1112 401-425.
- 1113 Berland, H., Albert, N.W., Stavland, A., Jordheim, M., McGhie, T.K., Zhou, Y., Zhang, H.,
1114 Deroles, S.C., Schwinn, K.E., Jordan, B.R., Davies, K.M., & Andersen, Ø.M. (2019).
1115 Auronidins are a previously unreported class of flavonoid pigments that challenges when
1116 anthocyanin biosynthesis evolved in plants. *Proceedings of the National Academy of*
1117 *Sciences USA*. 116, 20232-20239.
- 1118 Blanc, G., Agarkova, I., Grimwood, J., Kuo, A., Brueggeman, A., Dunigan, D. D., ... Van
1119 Etten, J. L. (2012). The genome of the polar eukaryotic microalga *Coccomyxa*
1120 *subellipsoidea* reveals traits of cold adaptation. *Genome Biology*, 13, R39.
- 1121 Blázquez, M.A., Nelson, D.C., & Weijers, D. (2020). Evolution of plant hormone response
1122 pathways. *Annual Review of Plant Biology*, 71, 327-353.
- 1123 Booij-James, I.S., Dube, S.K., Jansen, M.A.K., Edelman, M., & Mattoo, A.K. (2000).
1124 Ultraviolet-B radiation impacts light-mediated turnover of the photosystem II reaction

- 1125 center heterodimer in *Arabidopsis* mutants altered in phenolic metabolism. *Plant*
1126 *Physiology*, 124, 1275–1284
- 1127 Bowman, J. L., Kohchi, T., Yamato, K.T., Jenkins, J., Shu, S., Ishizaki, K., ... Schmutz, J.
1128 (2017). Insights into land plant evolution garnered from the *Marchantia polymorpha*
1129 genome. *Cell*, 171, 287-304.
- 1130 Byeon, Y., Lee, H. J., Lee, H. Y., & Back, K. (2016). Cloning and functional characterization
1131 of the *Arabidopsis* *N*-acetylserotonin *O*-methyltransferase responsible for melatonin
1132 synthesis. *Journal of Pineal Research*, 60, 65-73.
- 1133 Carella, P., Gogleva, A., Hoey, D.J., Bridgen, A.J., Stolze, S.C., Nakagami, H., & Schornack,
1134 S. (2019). Conserved biochemical defenses underpin host responses to oomycete
1135 infection in an early-divergent land plant lineage. *Current Biology*, 29, 2282–2294.e5.
- 1136 Cheng, S., Xian, W., Fu, Y., Marin, B., Keller, J., Wu, T., ... Melkonian, M. (2019) Genomes
1137 of subaerial Zygnematophyceae provide insights into land plant evolution. *Cell*, 179,
1138 1057-1067.e14.
- 1139 Chezem, W.R., Memon, A., Li, F.-S., Weng, J.-K., Clay, N.K. (2017). SG2-type R2R3-MYB
1140 transcription factor MYB15 controls defense-induced lignification and basal immunity in
1141 *Arabidopsis*. *The Plant Cell*, 29, 1907-1926.
- 1142 Clayton, W. A., Albert, N. W., Thrimawithana, A. H., McGhie, T. K., Deroles, S.
1143 C., Schwinn, K. E., ... Davies, K. M. (2018). UVR8-mediated induction of flavonoid
1144 biosynthesis for UVB tolerance is conserved between the liverwort *Marchantia*
1145 *polymorpha* and flowering plants. *Plant Journal*, 96, 503–517.
- 1146 Costa, M.A., Bedgar, D.L., Moinuddin, S.G.A., Kim, K.-W., Cardenas, C.L., Cochrane, F.C.,
1147 ... Lewis, N. G. (2005). Characterization in vitro and in vivo of the putative multigene 4-
1148 coumarate:CoA ligase network in *Arabidopsis*: syringyl lignin and sinapate/sinapyl
1149 alcohol derivative formation. *Phytochemistry*, 66: 2072–2091.
- 1150 Danielsson, M., Lundén, K., Elfstrand, M., Hu, J., Zhao, T., Arnerup, J., Ihrmark, K.,
1151 Swedjemark, G., Borg-Karlson, A.K., & Stenlid, J. (2011). Chemical and transcriptional
1152 responses of Norway spruce genotypes with different susceptibility to *Heterobasidion*
1153 spp. infection. *BMC Plant Biology*, 11, 154.
- 1154 D'Auria, J.C. (2006). Acyltransferases in plants: a good time to be BAHD. *Current Opinion*
1155 *in Plant Biology*, 9, 331-340.
- 1156 De Clerck, O., Kao, S.-M., Bogaert, K. A., Blomme, J., Foflonker, F., Kwantes, M. ...
1157 Bothwell, J. H. (2018) Insights into the evolution of multicellularity from the sea lettuce
1158 genome. *Current Biology*, 28, 2921-2933.e2925.

- 1159 Delwiche, C. F., Graham, L. E., Thomson, N. (1989). Lignin-like compounds and
1160 sporopollenin in *Coleochaete*, an algal model for land plant ancestry. *Science*, 245, 399–
1161 401.
- 1162 Devic, M., Guillemot, J., Debeaujon, I., Bechtold, N., Bensaude, E., Koornneef, M.,
1163 Pelletier, G. & Delseny, M. (1999). The BANYULS gene encodes a DFR-like protein
1164 and is a marker of early seed coat development. *The Plant Journal*, 19, 387-398.
- 1165 de Azevedo Souza, C., Kim, S.S., Koch, S., Kienow, L., Schneider, K., McKim, S.M.,
1166 Haughn, G.W., Kombrink, E., & Douglas, C.J. (2009). A novel fatty acyl-CoA
1167 synthetase is required for pollen development and sporopollenin biosynthesis in
1168 *Arabidopsis*. *Plant Cell* 21, 507-525.
- 1169 de Vries, J., Curtis, B.A., Gould, S.B. & Archibald, J.M. (2018). Embryophyte stress
1170 signaling evolved in the algal progenitors of land plants. *Proceedings of the National*
1171 *Academy of Sciences USA*, 115, E3471–E3480.
- 1172 de Vries, J., de Vries, S., Slamovits, C.H., Rose, L.E. & Archibald, J.M. (2017) How
1173 embryophytic is the biosynthesis of phenylpropanoids and their derivatives in
1174 streptophyte algae? *Plant Cell Physiol* 58, 934–945.
- 1175 de Vries, J., & Archibald, J.M. (2018) Plant evolution: Landmarks on the path to terrestrial
1176 life. *New Phytologist* 217, 1428–1434.
- 1177 de Vries, J. de Vries, S., Curtis, B. A., Zhou, H., Penny, S., Feussner, K., Pinto, D. M.,
1178 Steinert, M., Cohen, A. M., von Schwartzberg, K., & Archibald, J. M. (2020) Heat
1179 stress response in the closest algal relatives of land plants reveals conserved stress
1180 signalling circuits. *The Plant Journal*, 103, 1025-1048.
- 1181 de Vries, J., & Ischebeck, T. (2020) Ties between stress and lipid droplets pre-date seeds.
1182 *Trends Plant Sci.* 12, 1203-1214.
- 1183 de Vries, S., Herrfurth, C., Li, F.-W., Feussner, I., & de Vries, J. (2021) An ancient route
1184 towards salicylic acid and its implications for the perpetual *Trichormus–Azolla*
1185 symbiosis. *bioRxiv* preprint doi: <https://doi.org/10.1101/2021.03.12.435107>
- 1186 Dixon, R. A., & Paiva, N. L. (1995) Stress-induced phenylpropanoid metabolism. *Plant Cell*,
1187 7, 1085-1097.
- 1188 Dixon, R.A., Achnine, L., Kota, P., Liu, C.-J., Srinivasa Reddy, M.S., & Wang, L. (2002).
1189 The phenylpropanoid pathway and plant defence—a genomics perspective. *Molecular*
1190 *Plant Pathology* 3, 371-390.
- 1191 Do, C.-T., Pollet, B., Thévenin, J., Sibout, R., Denoue, D., Barrière, Y., Lapierre, C., &
1192 Jouanin, L. (2007) Both caffeoyl coenzyme A 3-O-methyltransferase 1 and caffeic

- 1193 acid *O*-methyltransferase 1 are involved in redundant functions for lignin, flavonoids and
1194 sinapoyl malate biosynthesis in *Arabidopsis*. *Planta*, 226, 1117–1129.
- 1195 Ehltng, J., Büttner, D., Wang, Q., Douglas, C.J., Somssich, I.E., & Kombrink, E. (1999).
1196 Three 4-coumarate:coenzyme A ligases in *Arabidopsis thaliana* represent two
1197 evolutionarily divergent classes in angiosperms. *The Plant Journal*, 19, 9-20.
- 1198 Emiliani, G., Fondi, M., Fani, R. & Gribaldo, S. (2009). A horizontal gene transfer at the
1199 origin of phenylpropanoid metabolism: a key adaptation of plants to land. *Biology*
1200 *Direct*, 4, 7.
- 1201 Eudes, A., Pereira, J.H., Yogiswara, S., Wang, G., Benites, V.T., Baidoo, E.E.K., Lee, T.S.,
1202 Adams, P.D., Keasling, J.D., & Loqué, D. (2016) Exploiting the substrate promiscuity of
1203 hydroxycinnamoyl-CoA: shikimate hydroxycinnamoyl transferase to reduce lignin. *Plant*
1204 *Cell Physiol*, 57, 568-579.
- 1205 Fellenberg, C., Milkowski, C., Hause, B., Lange, P., Böttcher, C., Schmidt, J., & Vogt,
1206 T. (2008). Tapetum-specific location of a cation-dependent *O*-methyltransferase
1207 in *Arabidopsis thaliana*. *The Plant Journal*, 56, 132– 145.
- 1208 Franks, N.P., Jenkins, A., Conti, E., Lieb, W.R., & Brick, P. (1998). Structural basis for the
1209 inhibition of firefly luciferase by a general anesthetic. *Biophysical Journal*, 75, 2205-
1210 2211.
- 1211 Ferrer, J.L., Zubieta, C., Dixon, R.A., & Noel, J.P. (2005). Crystal structures of alfalfa
1212 caffeoyl coenzyme A 3-*O*-methyltransferase. *Plant Physiology*, 137, 1009–1017.
- 1213 Fürst-Jansen, J.M.R., de Vries, S. & de Vries, J. (2020). Evo-physio: on stress responses and
1214 the earliest land plants. *Journal of Experimental Botany*, 11, 3254–3269.
- 1215 Goiris, K., Muylaert, K., Voorspoels, S., Noten, B., De Paepe, D., Baart, G. J. E., & De
1216 Cooman, L. (2014). Detection of flavonoids in microalgae from different evolutionary
1217 lineages. *Journal of Phycology*, 50, 483–492.
- 1218 Gouy, M., Guindon, S., & Gascuel, O. (2010) SeaView Version 4: A multiplatform graphical
1219 user interface for sequence alignment and phylogenetic tree building. *Molecular Biology*
1220 *and Evolution*, 27, 221–224.
- 1221 Güngör, E., Brouwer P., Dijkhuizen, L. W., Shaffar, D. C., Nierop, K. G. J., de Vos, R. C. H.,
1222 Toraño, J. S., van der Meer, I. N., & Schliepman, H. (2021) *Azolla* ferns testify: seed
1223 plants and ferns share a common ancestor for leucoanthocyanidin reductase enzymes.
1224 *New Phytologist*, 229, 1118-1132.

- 1225 Guo, D.-M., Ran, J.-H., & Wang, X.-Q. (2010). Evolution of the Cinnamyl/Sinapyl Alcohol
1226 Dehydrogenase (CAD/SAD) Gene Family: The Emergence of Real Lignin is Associated
1227 with the Origin of Bona Fide CAD. *Journal of Molecular Evolution*, 71, 202-218.
- 1228 Grienenberger, E., Kim, S. S., Lallemand, B., Geoffroy, P., Heintz, D., de Azevedo Souza,
1229 C., Heitz, T., Douglas, C. J., Legrand, M. (2010) Analysis of TETRAKETIDE α -
1230 PYRONE REDUCTASE function in *Arabidopsis thaliana* reveals a previously unknown,
1231 but conserved, biochemical pathway in sporopollenin monomer biosynthesis. *The Plant*
1232 *Cell* 22, 4067-4083.
- 1233 Gross, G. G., Stöckigt, J., Mansell, R. L., & Zenk, M. H. (1973). Three novel enzymes
1234 involved in the reduction of ferulic acid to coniferyl alcohol in higher plants: ferulate:
1235 CoA ligase, feruloyl-CoA reductase and coniferyl alcohol oxidoreductase. *FEBS Lett.* 31,
1236 283–286
- 1237 Ha, C.M., Escamilla-Trevino, L., Yarce, J.C.S., Kim, H., Ralph, J., Chen, F., & Dixon, R.A.
1238 (2016). An essential role of caffeoyl shikimate esterase in monolignol biosynthesis in
1239 *Medicago truncatula*. *The Plant Journal*, 86, 363-375.
- 1240 Hamberger, B., Ellis, M., Friedmann, M., de Azevedo Souza, C., Barbazuk, B., & Douglas,
1241 C.J. (2007). Genome-wide analyses of phenylpropanoid-related genes in *Populus*
1242 *trichocarpa*, *Arabidopsis thaliana*, and *Oryza sativa*: the *Populus* lignin toolbox and
1243 conservation and diversification of angiosperm gene families. *Canadian Journal of*
1244 *Botany*, 85, 1182–1201.
- 1245 Hoffmann, L., Maury, S., Martz, F., Geoffroy, P., & Legrand, M. (2003). Purification,
1246 cloning, and properties of an acyltransferase controlling shikimate and quinate ester
1247 intermediates in phenylpropanoid metabolism. *The Journal of Biological Chemistry*,
1248 278, 95–103.
- 1249 Hori, K., Maruyama, F., Fujisawa, T., Togashi, T., Yamamoto, N., Seo, M., ... Ohta, H.
1250 (2014). *Klebsormidium flaccidum* genome reveals primary factors for plant terrestrial
1251 adaptation. *Nature Communications*, 5, 3978.
- 1252 Hu, Y., Gai, Y., Yin, L., Wang, X., Feng, C., Feng, L., Li, D., Jiang, X.N., & Wang, D.C.
1253 (2010). Crystal structures of a *Populus tomentosa* 4-coumarate:CoA ligase shed light on
1254 its enzymatic mechanisms. *The Plant Cell*, 22, 3093-3104.
- 1255 Huerta-Cepas, J., Serra, F., & Bork, P. ETE 3: Reconstruction, analysis, and visualization of
1256 phylogenomic data. *Molecular Biology and Evolution*, 33, 1635–1638.
- 1257 Jahns, P., & Holzwarth, A.R. (2012). The role of the xanthophyll cycle and of lutein in
1258 photoprotection of photosystem II. *BBA – Bioenergetics*, 1817, 182–193.

- 1259 Jiao, C., Sørensen, I., Sun, X., Sun, H., Behar, H., Alseikh, S., ... Rose, J. K. C. (2020). The
1260 *Penium margaritaceum* genome: hallmarks of the origins of land plants. *Cell*, 181,
1261 P1097-1111.E12
- 1262 Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., ... Hunter, S. (2014).
1263 InterProScan 5: genome-scale protein function classification. *Bioinformatics*, 30, 1236-
1264 1240.
- 1265 Ju, C., Van de Poel, B., Cooper, E. D., Thierer, J. H., Gibbons, T. R., Delwiche, C. F., &
1266 Chang, C. (2015). Conservation of ethylene as a plant hormone over 450 million years of
1267 evolution. *Nature Plants*, 1, 14004.
- 1268 Kaur, H., Heinzl, N., Schöttner, M., Baldwin, I.T., & Gális, I. (2010). R2R3-NaMYB8
1269 regulates the accumulation of phenylpropanoid-polyamine conjugates, which are
1270 essential for local and systemic defense against insect herbivores in *Nicotiana attenuata*.
1271 *Plant Physiology*, 152, 1731-1747.
- 1272 Kim, S.-J., Kim, M.-R., Bedgar, D.L., Moinuddin, S.G.A., Cardenas, C.L., Davin, L.B.,
1273 Kang, C., & Lewis, N.G. (2004). Functional reclassification of the putative cinnamyl
1274 alcohol dehydrogenase multigene family in *Arabidopsis*. *Proceedings of the National*
1275 *Academy of Sciences USA*, 101, 1455-1460.
- 1276 Kim, R. J., Kim, H.J., Shim, D., & Suh, M.C. (2016). Molecular and biochemical
1277 characterizations of the monoacylglycerol lipase gene family of *Arabidopsis thaliana*.
1278 *The Plant Journal*, 85, 758-771.
- 1279 König, S., Feussner, K., Kaefer, A., Landesfeind, M., Thurow, C., Karlovsky, P., Gatz, C.,
1280 Polle, A., & Feussner, I. (2014). Soluble phenylpropanoids are involved in the defense
1281 response of *Arabidopsis* against *Verticillium longisporum*. *New Phytologist*, 202, 823-
1282 837.
- 1283 Kriegshauser, L., Knosp, S., Grienberger, E., Tatsumi, K., Gütle, D. D., Sørensen, I., ...
1284 Renault, H. (2021) Function of the HYDROXYCINNAMOYL-CoA:SHIKIMATE
1285 HYDROXYCINNAMOYL TRANSFERASE is evolutionarily conserved in
1286 embryophytes. *The Plant Cell*, in press, preprint on bioRxiv, 2020.09.16.300285.
- 1287 Labar, G., Bauvois, C., Borel, F., Ferrer, J. □L., Wouters, J., & Lambert, D.M. (2010), Crystal
1288 Structure of the Human Monoacylglycerol Lipase, a Key Actor in Endocannabinoid
1289 Signaling. *ChemBioChem*, 11, 218-227.
- 1290 Labeeuw, L., Martone, P.T., Boucher, Y., & Case, R.J. (2015) Ancient origin of the
1291 biosynthesis of lignin precursors. *Biology Direct* 10, 23.

- 1292 Lacombe, E., Hawkins, S., Doorselaere, J.V., Piquemal, J., Goffner, D., Poeydomenge, O.,
1293 Boudet, A-M., & Grima-Pettenati, J. G. (1997). Cinnamoyl CoA reductase, the first
1294 committed enzyme of the lignin branch biosynthetic pathway: cloning, expression and
1295 phylogenetic relationships. *The Plant Journal* 11, 429–441.
- 1296 Lang, D., Ullrich, K. K., Murat, F., Fuchs, J., Jenkins, J., Haas, F. B., ... Rensing, S. A.
1297 (2018) The *Physcomitrella patens* chromosome-scale assembly reveals moss genome
1298 structure and evolution. *The Plant Journal*, 93, 515–533.
- 1299 Lamesch, P., Berardini, T. Z., Li, D., Swarbreck, D., Wilks, C., Sasidharan, R., ... Huala, E.
1300 (2011). The Arabidopsis Information Resource (TAIR): improved gene annotation and
1301 new tools. *Nucleic Acids Research*, 40, D1202-D1210.
- 1302 Le S.Q, & Gascuel O. (2008). An Improved General Amino Acid Replacement Matrix. *Mol.*
1303 *Biol. Evol.* 25, 1307-1320.
- 1304 Leebens-Mack, J.H., Barker, M.S., Carpenter, E.J., Deyholos, M. K., Gitzendanner, M. A.,
1305 Graham, S. W., ... Wong, G. K.-S. (2019) One thousand plant transcriptomes and the
1306 phylogenomics of green plants. *Nature*, 574, 679–685.
- 1307 Letunic, I., & Bork, P. (2019). Interactive Tree Of Life (iTOL) v4: recent updates and new
1308 developments. *Nucleic Acids Research*, 47, W256-W259.
- 1309 Li, W., Zhang, F., Wu, R., Jia, R., Li, G., Guo, Y., Liu, C., & Wang, W. (2017). A novel N-
1310 Methyltransferase in Arabidopsis appears to feed a conserved pathway for nicotine
1311 detoxification among land plants and is associated with lignin biosynthesis. *Plant*
1312 *Physiology*, 174, 1492-1504.
- 1313 Li, F.-W., Nishiyama, T., Waller, M., Frangedakis, E., Keller, J., Li, Z. ... Szövényi, P.
1314 (2020) Anthoceros genomes illuminate the origin of land plants and the unique biology
1315 of hornworts. *Nature Plants*, 6, 259-272.
- 1316 Li, F.-W., Brouwer, P., Carretero-Paulet, L., Cheng, S., de Vries, J., Delaux, P.-M., ... Pryer,
1317 K. M. (2018) Fern genomes elucidate land plant evolution and cyanobacterial symbioses.
1318 *Nature Plants*, 4, 460-472.
- 1319 Louie, G. V., Bowman, M. E., Tu, Y., Mouradov, A., Spangenberg, G., & Noel,
1320 J.P. (2010). Structure–function analyses of a caffeic acid O-methyltransferase from
1321 perennial ryegrass reveal the molecular basis for substrate preference. *The Plant Cell*, 22,
1322 4114– 4127.
- 1323 Maeda, H. A., & Fernie, A. R. (2021) Evolutionary history of plant metabolism. *Annual*
1324 *Review of Plant Biology*, 72, 10.1146/annurev-arplant-080620-031054

- 1325 Mansell, R. L. G., Gross, G. G., Stöckigt, J., Franke, H., & Zenk, M. H. (1974) Purification
1326 and properties of cinnamyl alcohol dehydrogenase from higher plants involved in lignin
1327 biosynthesis. *Phytochemistry*, 13, 2427–2435.
- 1328 Matsuno, M., Compagnon, V., Schoch, G.A., Schmitt, M., Debayl, D., Bassard, J.-E. ...
1329 Werck-Reichhart, D. (2009). Evolution of a novel phenolic pathway for pollen
1330 development. *Science* 325, 1688–1692.
- 1331 Martone, P. T., Estevez, J. M., Lu, F., Ruel, K., Denny, M. W., Somerville, C., & Ralph, J.
1332 (2009). Discovery of lignin in seaweed reveals convergent evolution of cell-wall
1333 architecture. *Current Biology*, 19, 169-175.
- 1334 Martz, F., Maury, S., Pincon, G., & Legrand, M. (1998). cDNA cloning, substrate specificity
1335 and expression study of tobacco caffeoyl- CoA 3-O-methyltransferase, a lignin
1336 biosynthetic enzyme. *Plant Molecular Biology*, 36, 427–437.
- 1337 Maury, S., Geoffroy, P., & Legrand, M. (1999). Tobacco O-methyltransferases involved in
1338 phenylpropanoid metabolism: the different caffeoyl- coenzyme A/5-hydroxyferuloyl-
1339 coenzyme A 3/5-O-methyltransferase and caffeic acid/5-hydroxyferulic acid 3/5-O-
1340 methyltransferase classes have distinct substrate specificities and expression patterns.
1341 *Plant Physiology*, 121, 215–224.
- 1342 Merchant, S. S., Prochnik, S. E., Vallon, O., Harris, E. H., Karpowicz, S. J., Witman, G. B.,
1343 ... Grossman, A. R. (2007). The *Chlamydomonas* genome reveals the evolution of key
1344 animal and plant functions. *Science* 318, 245-250.
- 1345 Miller, M., Owens, S.J., & Rørslett, B. (2011). Plants and colour: flowers and pollination.
1346 *Optics & Laser Technology*, 43, 282-294.
- 1347 Moffitt, M.C., Louie, G.V., Bowman, M.E., Pence, J., Noel, J.P., & Moore, B.S. (2007).
1348 Discovery of two cyanobacterial phenylalanine ammonia lyases: kinetic and structural
1349 characterization. *Biochemistry*, 46, 1004–1012.
- 1350 Moreau, H., Verhelst, B., Couloux, A., Derelle, E., Rombauts, S., Grimsley, N., ...
1351 Vandepoele, K. (2012). Gene functionalities and genome structure in *Bathycoccus*
1352 *prasinus* reflect cellular specializations at the base of the green lineage. *Genome Biology*,
1353 13, R74.
- 1354 Nakatsu, T., Ichiyama, S., Hiratake, J., Saldanha, A., Kobashi, N., Sakata, K., & Kato, H.
1355 (2006). Structural basis for the spectral difference in luciferase bioluminescence. *Nature*,
1356 440, 372-376.
- 1357 Nelson, D., & Werck-Reichhart, D. (2011). A P450-centric view of plant evolution. *The*
1358 *Plant Journal*, 66, 194-211.

- 1359 Nishiyama, T., Sakayama, H., de Vries, J., Buschmann, H., Saint-Marcoux, D., Ullrich, K.
1360 K., ... Rensing, S. A. (2018). The *Chara* genome: secondary complexity and
1361 implications for plant terrestrialization. *Cell*, 174, 448–464.
- 1362 Nystedt, B., Street, N. R., Wetterbom, A., Zuccolo, A., Lin, Y.-C., Scofield, D. G., ...
1363 Jansson, S. (2013). The Norway spruce genome sequence and conifer genome evolution.
1364 *Nature*, 497, 579-584.
- 1365 Oliva, J., Rommel, S., Fossdal, C.G., Hietala, A.M., Nemesio-Gorriz, M., Solheim, H., &
1366 Elfstrand, M. (2015). Transcriptional responses of Norway spruce (*Picea abies*) inner
1367 sapwood against *Heterobasidion parviporum*. *Tree Physiology*, 35, 1007–1015.
- 1368 Omura, T. (1999). Forty years of cytochrome P450. *Biochemical and Biophysical Research*
1369 *Communications*, 266, 690-698.
- 1370 Ouyang, S., Zhu, W., Hamilton, J., Lin, H., Campbell, M., Childs, K., ... Buell, C. R. (2007).
1371 The TIGR Rice Genome Annotation Resource: improvements and new features. *Nucleic*
1372 *Acids Research*, 35, D883–D887.
- 1373 Overdijk, E.J.R., de Keijzer, J., de Groot, D., Schoina, C., Bouwmeester, K., Ketelaar, T., &
1374 Govers, F. (2016). Interaction between the moss *Physcomitrella patens* and
1375 *Phytophthora*: a novel pathosystem for live \square cell imaging of subcellular defence. *Journal*
1376 *of Microscopy*, 263, 171–180.
- 1377 Palenik, B., Grimwood, J., Aerts, A., Rouzé, P., Salamov, A., Putman, N., ... Grigoriev, I. V.
1378 (2007). The tiny eukaryote *Ostreococcus* provides genomic insights into the paradox of
1379 plankton speciation. *Proceedings of the National Academy USA*. 104, 7705–7710.
- 1380 Pan, H., Zhou, R., Louie, G.V., Mühlemann, J. K., Bomati, E. K., Bowman, M. E., ... Wang,
1381 X. (2014). Structural Studies of Cinnamoyl-CoA Reductase and Cinnamyl-Alcohol
1382 Dehydrogenase, key enzymes of monolignol biosynthesis. *The Plant Cell* 26, 3709.
- 1383 Piatkowski, B. T., Imwattana, K., Tripp, E. A., Weston, D. J., Healey, A., Schmutz, J., &
1384 Shaw, A. J. (2020). Phylogenomics reveals convergent evolution of red-violet coloration
1385 in land plants and the origins of the anthocyanin biosynthetic pathway. *Molecular*
1386 *Phylogenetics and Evolution*, 151, 106904.
- 1387 Ponce De León, I., Schmelz, E. A., Gaggero, C., Castro, A., Álvarez, A., & Montesano, M.
1388 (2012). *Physcomitrella patens* activates reinforcement of the cell wall, programmed cell
1389 death and accumulation of evolutionary conserved defence signals, such as salicylic acid
1390 and 12-oxo-phytodienoic acid, but not jasmonic acid, upon *Botrytis cinerea* infection.
1391 *Molecular Plant Pathology*, 13, 960-74.

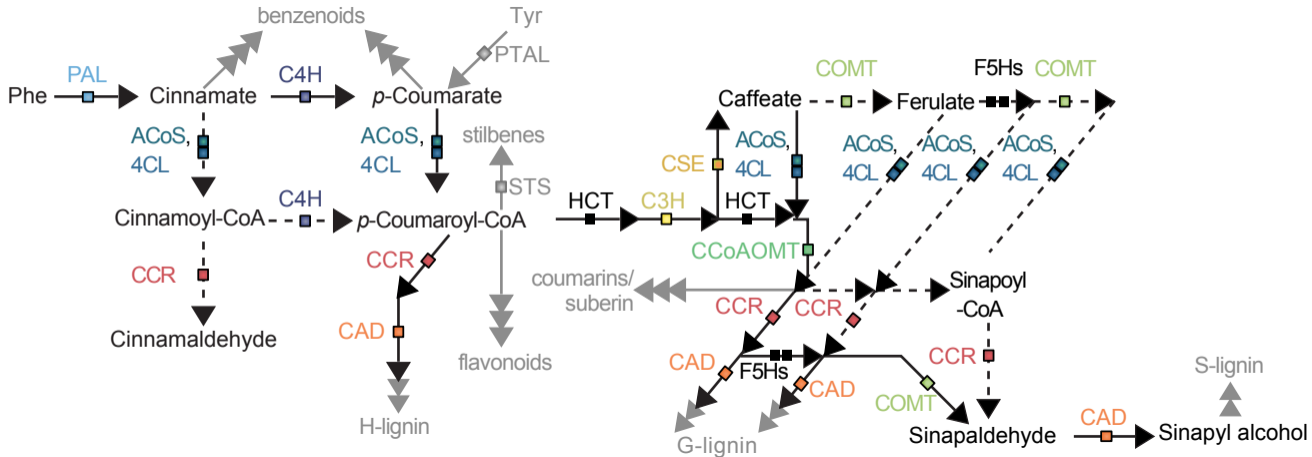
- 1392 Prochnik, S.E., Umen, J., Nedelcu, A. M., Hallmann, A., Miller, S. M., Nishii, I., ... Rokhsar,
1393 D. S. (2010). Genomic analysis of organismal complexity in the multicellular green alga
1394 *Volvox carteri*. *Science*, 329, 223–226
- 1395 Puttick, M.N., Morris, J.L., Williams, T.A., Cox, C.J., Edwards, D., Kenrick, P., Pressel, S.,
1396 Wellman, C.H., Schneider, H., Pisani, D., & Donoghue, P.C.J. (2018). The
1397 interrelationships of land plants and the nature of the ancestral embryophyte. *Current*
1398 *Biology*, 28, 733-745.
- 1399 Ralph, J., Lundquist, K., Brunow, G., Lu, F., Kim, H., Schatz, P. F., Marita, J. M., Hatfield,
1400 R. D., Ralph, S. A., Christensen, J. H., & Boerjan, W. (2004). Lignins: natural polymers
1401 from oxidative coupling of 4-hydroxyphenylpropanoids. *Phytochemistry Reviews*, 3, 29–
1402 60.
- 1403 Rupasinghe S., Baudry J. & Schuler M.A. (2003) Common active site architecture and
1404 binding strategy of four phenylpropanoid P450s from *Arabidopsis thaliana* as revealed
1405 by molecular modeling. *Protein Engineering Design and Selection*, 16, 721–731.
- 1406 Renault, H., Alber, A., Horst, N.A., Basilio Lopes, A., Fich, E. A., Kriegshausen, L., ...
1407 Werck-Reichhart, D. (2017a). A phenol-enriched cuticle is ancestral to lignin evolution
1408 in land plants. *Nature Communications*, 8, 14713.
- 1409 Renault, H., De Marothy, M., Jonasson, G., Lara, P., Nelson, D. R., Nilsson, I., André, F.,
1410 von Heijne, G., & Werck-Reichhart, D. (2017b). Gene duplication leads to altered
1411 membrane topology of a cytochrome P450 enzyme in seed plants. *Molecular Biology*
1412 *and Evolution*, 34, 2041-2056.
- 1413 Renault, H., Werck-Reichhart, D., & Weng, J-K. (2019). Harnessing lignin evolution for
1414 biotechnological applications. *Current Opinion in Biotechnology*, 56, 105-111.
- 1415 Rensing, S.A. (2014). Gene duplication as a driver of plant morphogenetic evolution. *Current*
1416 *Opinion in Plant Biology*, 17, 43-48.
- 1417 Rensing, S.A. (2018) Great moments in evolution: the conquest of land by plants. *Current*
1418 *Opinion in Plant Biology*, 42, 49–54.
- 1419 Rippin, M., Becker, B. and Holzinger, A. (2017) Enhanced desiccation tolerance in mature
1420 cultures of the streptophytic green alga *Zygnema circumcarinatum* revealed by
1421 transcriptomics. *Plant & Cell Physiology*, 58, 2067–2084.
- 1422 Rippin, M., Pichrtová, M., Arc, E., Kranner, I., Becker, B. and Holzinger, A. (2019)
1423 Metatranscriptomic and metabolite profiling reveals vertical heterogeneity within a
1424 *Zygnema* green algal mat from Svalbard (High Arctic). *Environmental Microbiology*, 21,
1425 4283-4299.

- 1426 Ro, D. K., Mah, N., Ellis, B. E., & Douglas, C. J. (2001). Functional characterization and
1427 subcellular localization of poplar (*Populus trichocarpa* x *Populus deltoides*) cinnamate
1428 4-hydroxylase. *Plant Physiology*, 126, 317-329.
- 1429 Russell, D.W., Conn, E.E. (1967). The cinnamic acid 4-hydroxylase from pea seedlings.
1430 *Archives of Biochemistry and Biophysics* 122, 256-258.
- 1431 Saballos, A., Ejeta, G., Sanchez, E., Kang, C., & Vermerris, W. (2009) A genomewide
1432 analysis of the cinnamyl alcohol dehydrogenase family in sorghum [*Sorghum bicolor*
1433 (L.) Moench] identifies SbCAD2 as the brown midrib6 gene. *Genetics*, 181, 783–795.
- 1434 Saito, N., & Harborne, J.B. (1992). Correlations between anthocyanin type, pollinator and
1435 flower colour in the labiatae. *Phytochemistry*, 31, 3009-3015.
- 1436 Scheres, B., & van der Putten, W.H. (2017) The plant perceptron connects environment to
1437 development. *Nature* 543, 337-345.
- 1438 Schoenbohm, C., Martens, S., Eder, C., Forkmann, G., & Weisshaar, B. (2000) Identification
1439 of the *Arabidopsis thaliana* flavonoid 3'-hydroxylase gene and functional expression of
1440 the encoded P450 enzyme. *Biological Chemistry*, 381, 749-753.
- 1441 Sheahan, J.J. (1996) Sinapate esters provide greater UV-B attenuation than flavonoids in
1442 *Arabidopsis thaliana* (Brassicaceae). *American Journal of Botany*, 83, 679-686.
- 1443 Sheehan, H., Hermann, K., & Kuhlemeier, C. (2012). Color and Scent: How Single Genes
1444 Influence Pollinator Attraction. *Cold Spring Harbor Symposia on Quantitative Biology*,
1445 77, 117-133.
- 1446 Shockey, J.M., Fulda, M.S., & Browse, J. (2003). Arabidopsis contains a large superfamily of
1447 Acyl-Activating Enzymes. Phylogenetic and biochemical analysis reveals a new class of
1448 Acyl-Coenzyme A Synthetases. *Plant Physiology*, 132, 1065-1076.
- 1449 Sierro, N. Battey, J., Ouadi, S., Bakaher, N., Bovet, L., Willig, A., Goepfert, S., Peitsch, M.
1450 C., & Ivanov, N. V. (2014). The tobacco genome sequence and its comparison with those
1451 of tomato and potato. *Nature Communications*, 5, 3833.
- 1452 Slotte, T., Hazzouri, K. M., Ågren, J. A., Koenig, D., Maumus, F., Guo, Y.-L., ... Wrigth, S.
1453 I. (2013). The *Capsella rubella* genome and the genomic consequences of rapid mating
1454 system evolution. *Nature Genetics*, 45, 831–835
- 1455 Sørensen, I., Pettolino, F.A., Bacic, A., Ralph, J., Lu, F., O'Neill, M.A., Fei, Z., Rose, J.K.C.,
1456 Domozych, D.S., Willats, W.G.T. (2011). The charophycean green algae provide
1457 insights into the early origins of plant cell walls. *The Plant Journal*, 68, 201–211.
- 1458 Suzuki, S., & Umezawa, T. (2007). Biosynthesis of lignans and norlignans. *Journal of Wood*
1459 *Science* 53, 273-284.

- 1460 Sytar, O., Zivcak, M., Bruckova, K., Brestic, M., Hemmerich, I., Rauh, C., Simko, I. 2018.
1461 Shift in accumulation of flavonoids and phenolic acids in lettuce attributable to changes
1462 in ultraviolet radiation and temperature. *Scientia Horticulturae*, 239, 193-204.
- 1463 Szövényi, P., Frangedakis, E., Ricca, M., Quandt, D., Wicke, S., Langdale, J.A. (2015).
1464 Establishment of *Anthoceros agrestis* as a model species for studying the biology of
1465 hornworts. *BMC Plant Biology*, 15, 1-7.
- 1466 Szövényi, P., Gunadi, A., Li, F.-W. (2021). Charting the genomic landscape of seed-free
1467 plants. *Nature Plants*, <https://doi.org/10.1038/s41477-021-00888-z>
- 1468 Tan, D.-X., Hardeland, R., Manchester, L.C., Korkmaz, A., Ma, S., Rosales-Corral, S.,
1469 Reiter, R.J. (2012). Functional roles of melatonin in plants, and perspectives in
1470 nutritional and agricultural science. *Journal of Experimental Botany*, 63, 577-597.
- 1471 Tanaka, A., Shigemitsu, T., Yokota, Y., Shika, N. (1997). A new *Arabidopsis* mutant induced
1472 synthesis with spotted pigmentation. *Genes & Genetic Systems*, 72, 141–148.
- 1473 The International *Brachypodium* Initiative. (2010). Genome sequencing and analysis of the
1474 model grass *Brachypodium distachyon*. *Nature*, 463, 763–768.
- 1475 Urban, P., Werck-Reichhart, D., Teutsch, H.G., Durst, F., Regnier, S., Kazmeier, M., &
1476 Pompon, D. (1994). Characterization of recombinant plant cinnamate 4-hydroxylase
1477 produced in yeast. Kinetic and spectral properties of the major plant P450 of the
1478 phenylpropanoid pathway. *European Journal of Biochemistry*, 222, 843–850.
- 1479 Vanholme, R., Storme, V., Vanholme, B., Sundin, S., Christensen, J.H., Goemine, G.,
1480 Halpin, C., Rohde, A., Morreel, K., Boerjan, W. (2012) A systems biology view of
1481 responses to lignin biosynthesis perturbations in *Arabidopsis*. *The Plant Cell*, 24, 3506-
1482 3529.
- 1483 Vanholme, R., Cesarino, I., Rataj, K., Xiao, Y., Sundin, L., Goeminne, G., ... Boerjan, W.
1484 (2013). Caffeoyl shikimate esterase (CSE) is an enzyme in the ligninbiosynthetic
1485 pathway in *Arabidopsis*. *Science*, 341, 1103–1106.
- 1486 Vanholme, R., De Meester, B., Ralph, J., & Boerjan, W. (2019). Lignin biosynthesis and its
1487 integration into metabolism. *Current Opinion Biotechnology*, 56, 230-239.
- 1488 Vogt, T. (2010). Phenylpropanoid biosynthesis. *Molecular Plant*, 3, 2–20.
- 1489 Wan, T. et al. (2018). A genome for gnetophytes and early evolution of seed plants. *Nature*
1490 *Plants*, 4, 82–89.
- 1491 Wang, S., Li, L., Li, H., Sahu, S. K., Wang, H., Xu, Y., ... Liu, X. (2020) Genomes of early-
1492 diverging streptophyte algae shed light on plant terrestrialization. *Nature Plants*, 6, 95-
1493 106.

- 1494 Weng, J.-K., Li, X., Stout, J., & Chapple, C. (2008) Independent origins of syringyl lignin in
1495 vascular plants. *Proceedings of the National Academy of Sciences USA*, 105, 7887–7892.
- 1496 Weng, J.-K., & Chapple, C. (2010) The origin and evolution of lignin biosynthesis. *New*
1497 *Phytologist*, 187, 273-285.
- 1498 Weng, J.-K., Akiyama, T., Ralph, J., & Chapple, C. (2011) Independent recruitment of an O-
1499 methyltransferase for syringyl lignin biosynthesis in *Selaginella moellendorffii*. *The*
1500 *Plant Cell*, 23, 2708–2724.
- 1501 Weng, J. K. (2013) The evolutionary paths towards complexity: a metabolic perspective. *New*
1502 *Phytologist*, 201, 1141-1149.
- 1503 Whelan, S., & Goldman, N. (2001). A general empirical model of protein evolution derived
1504 from multiple protein families using a maximum-likelihood approach. *Molecular Biology*
1505 *and Evolution*, 18, 691-699.
- 1506 Wickett, N.J., Mirarab, S., Nguyen, N., Warnow, T., Carpenter, E., Matasci, N., ... Leebens-
1507 Mack, J. (2014) Phylotranscriptomic analysis of the origin and early diversification of
1508 land plants. *Proceedings of the National Academy of Sciences USA*, 111, E4859–E4868.
- 1509 Wodniok, S., Brinkmann, H., Glöckner, G., Heidel, A. J., Philippe, H., Melkonian, M., &
1510 Becker, B. (2011) Origin of land plants: do conjugating green algae hold the key? *BMC*
1511 *Evol. Biol.* 11, 104.
- 1512 Wohl, J., & Petersen, M. (2020). Functional expression and characterization of cinnamic acid
1513 4-hydroxylase from the hornwort *Anthoceros agrestis* in *Physcomitrella patens*. *Plant*
1514 *Cell Reports*, 39, 597-607.
- 1515 Wolf, L., Rizzini, L., Stracke, R., Ulm, R., & Rensing, S. A. (2010) The molecular and
1516 physiological responses of *Physcomitrella patens* to ultraviolet-B radiation. *Plant*
1517 *Physiol.* 153, 1123-1134.
- 1518 Worden, A. Z., Lee, J.-H., Mock, T., Rouzé, P., Simmons, M. P., Aerts, A. L., ... Grigoriev,
1519 I. V. (2009). Green evolution and dynamic adaptations revealed by genomes of the
1520 marine picoeukaryotes *Micromonas*. *Science*, 324, 268–272.
- 1521 Xu, Z., Zhang, D., Hu, J., Zhou, X., Ye, X., Reichel, K. L., ... Yuan, J. S. (2009).
1522 Comparative genome analysis of lignin biosynthesis gene families across the plant
1523 kingdom. *BMC Bioinformatics* 10, S3.
- 1524 Xue, J.-S., Zhang, B., Zhan, H., Lv, Y.-L., Jia, X.-L., Wang, T., ..., Yang, Z.-N. (2020).
1525 Phenylpropanoid derivatives are essential components of sporopollenin in vascular
1526 plants. *Molecular Plant* 13, 1644-1653.

- 1527 Yang, J., Yan, R., Roy, A., Xu, D., Poisson, J., & Zhang, Y. (2015). The I-TASSER Suite:
1528 protein structure and function prediction. *Nature Methods* 12, 7-8.
- 1529 Ye, Z.-H., Kneusel, R.E., Matern, U., & Varner, J. E. (1994). An alternative methylation
1530 pathway in lignin biosynthesis in *Zinnia*. *The Plant Cell*, 6, 1427–1439.
- 1531 Ye, Z.-H., & Varner, J.E. (1995). Differential expression of two *O*-methyltransferases in
1532 lignin biosynthesis in *Zinnia elegans*. *Plant Physiology*, 108, 459–467.
- 1533 Youn, B., Camacho, R., Moinuddin, S.G.A., Lee, C., Davin, L.B., Lewis, N.G., ..., (2006)
1534 Crystal structures and catalytic mechanism of the *Arabidopsis* cinnamyl alcohol
1535 dehydrogenases AtCAD5 and AtCAD4. *Organic and Biomolecular Chemistry*, 4, 1687–
1536 1697.
- 1537 Zhang, Y. (2008). I-TASSER server for protein 3D structure prediction. *BMC Bioinformatics*,
1538 9, 40.
- 1539 Zhang, J., Fu, X.-X., Li, R.-Q., Zhao, X., Liu, Y., Li, M.-H., ... Chen, Z.-D. (2020). The
1540 hornwort genome and early land plant evolution. *Nature Plants*, 6, 107-118.
- 1541 Zubieta, C., Kota, P., Ferrer, J.L., Dixon, R.A., & Noel, J.P. (2002). Structural basis for the
1542 modulation of lignin monomer methylation by caffeic acid/5-hydroxyferulic acid 3/5-O-
1543 methyltransferase. *The Plant Cell*, 14, 1265-1277.



0.7

Angiosperms

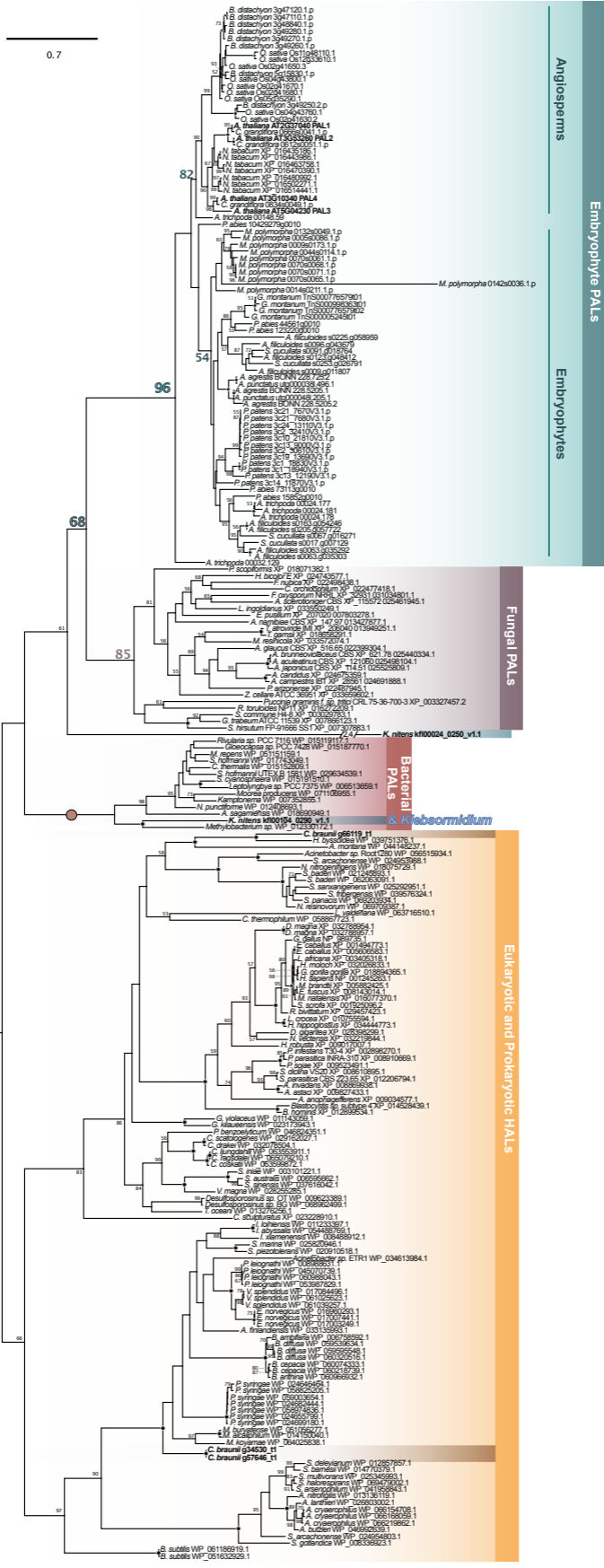
Embryophyte PALs

Embryophytes

Fungal PALs

Bacterial & Kiebsormidium

Eukaryotic and Prokaryotic HALs

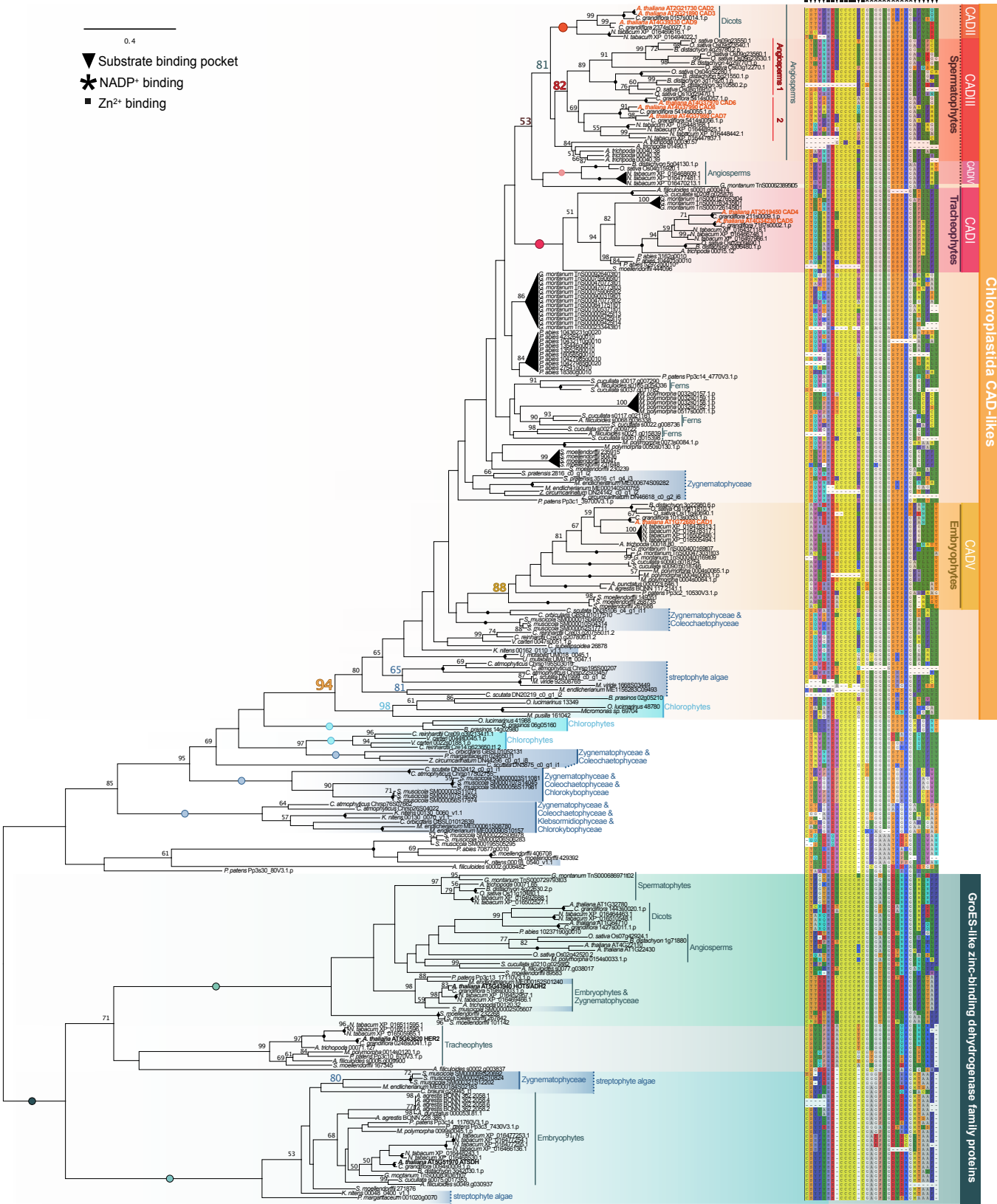


B. subtilis WP_06118691.1
B. subtilis WP_051632929.1

C. braunii g34530_t1
C. braunii g57646_t1

S. deleyianum WP_012857857.1
S. barnesi WP_014770379.1
S. multivorans WP_025345993.1
S. halodispersans WP_08473902.1
S. arsenophilum WP_041953843.1
A. nitrospira WP_013136119.1
A. lamarii WP_026903002.1
A. cyanoerophilum WP_06154708.1
A. cyanoerophilum WP_066168059.1
A. cyanoerophilum WP_066219862.1
A. butleri WP_046952631.1
S. arachonense WP_024954803.1
S. gottwaldica WP_008336923.1

- ▼ Substrate binding pocket
- * NADP⁺ binding
- Zn²⁺ binding



CADII
Spermatophytes
CADV
Tracheophytes
CADU
Chloroplastidata CAD-likees

CADV
Embryophytes

Chlorophytes
Streptophyte algae
Zygnemathophyceae & Coleochaetophyceae

Embryophytes & Zygnemathophyceae
Dicots
Angiosperms
Embryophytes
Streptophyte algae

Angiosperms 1 2

Angiosperms

Zygnemathophyceae

Zygnemathophyceae & Coleochaetophyceae

streptophyte algae

Chlorophytes

Chlorophytes

Zygnemathophyceae & Coleochaetophyceae

Zygnemathophyceae & Coleochaetophyceae & Chlorokybophyceae

Zygnemathophyceae & Coleochaetophyceae & Klebsorminophyceae & Chlorokybophyceae

embryophytes & Zygnemathophyceae

embryophytes

embryophytes & Zygnemathophyceae

embryophytes

embryophytes & Zygnemathophyceae

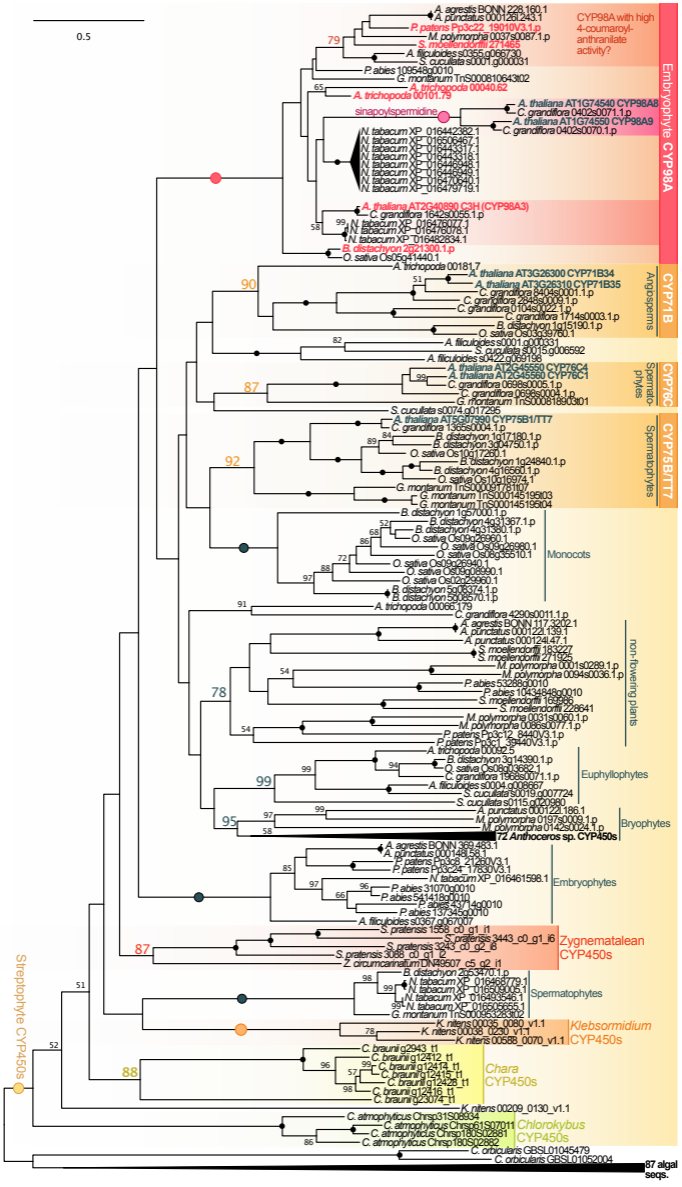
embryophytes

embryophytes & Zygnemathophyceae

embryophytes

embryophytes & Zygnemathophyceae

embryophytes



0.8

249 mammal & Sauria MAGLs
GMYLISPLVtAs

Streptophyte algal MAGLs

Streptophyte MAGL6/7/8/9/10/11/12

MAGL9
MAGL12
MAGL9/12
MAGL7
MAGL8
MAGL6
MAGL10
MAGL11

GAYLAPNCKIAs

Streptophyte algal MAGLs

Streptophyte MAGL5/15

Angiosperms
MAGL15
MAGL5
Embryophytes

GAYLAPNCKIAs

Streptophyte algal MAGLs

Embryophyte MAGL1/3

CSE/MAGL3
MAGL1

GATESAPLEXIP

Embryophyte MAGL14/16

MAGL14
MAGL16
GAYLAPNCKIAs

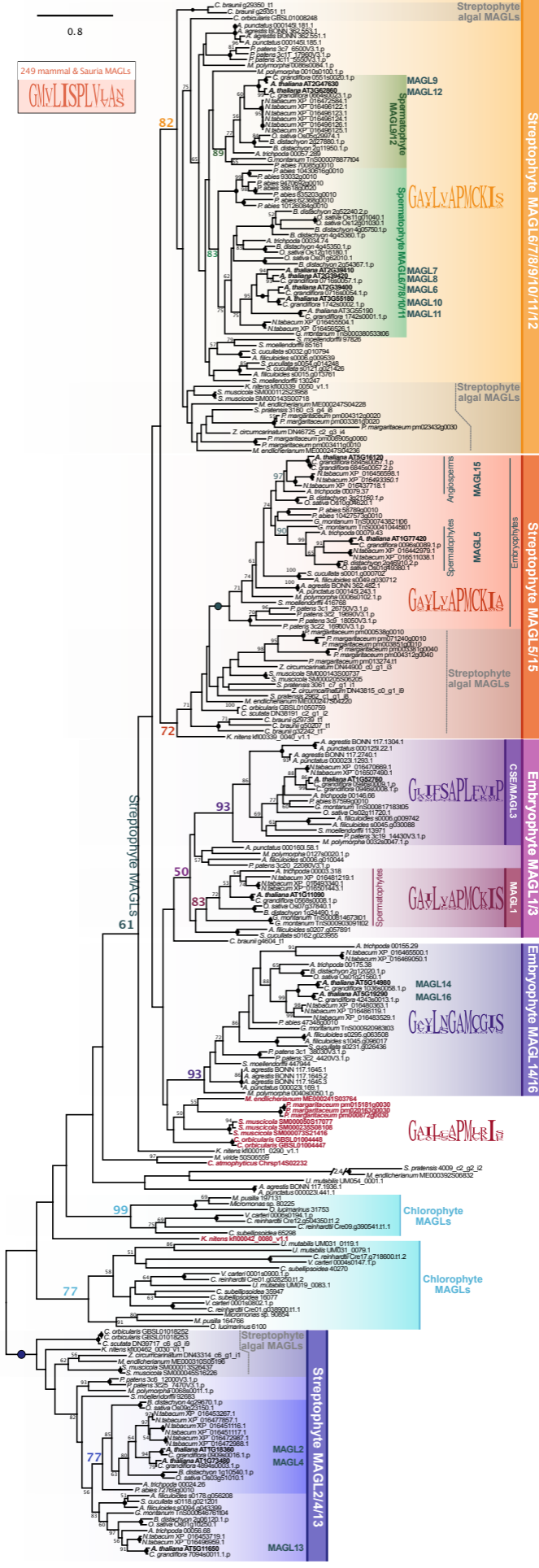
GAYLAPMERIAs

Chlorophyte MAGLs

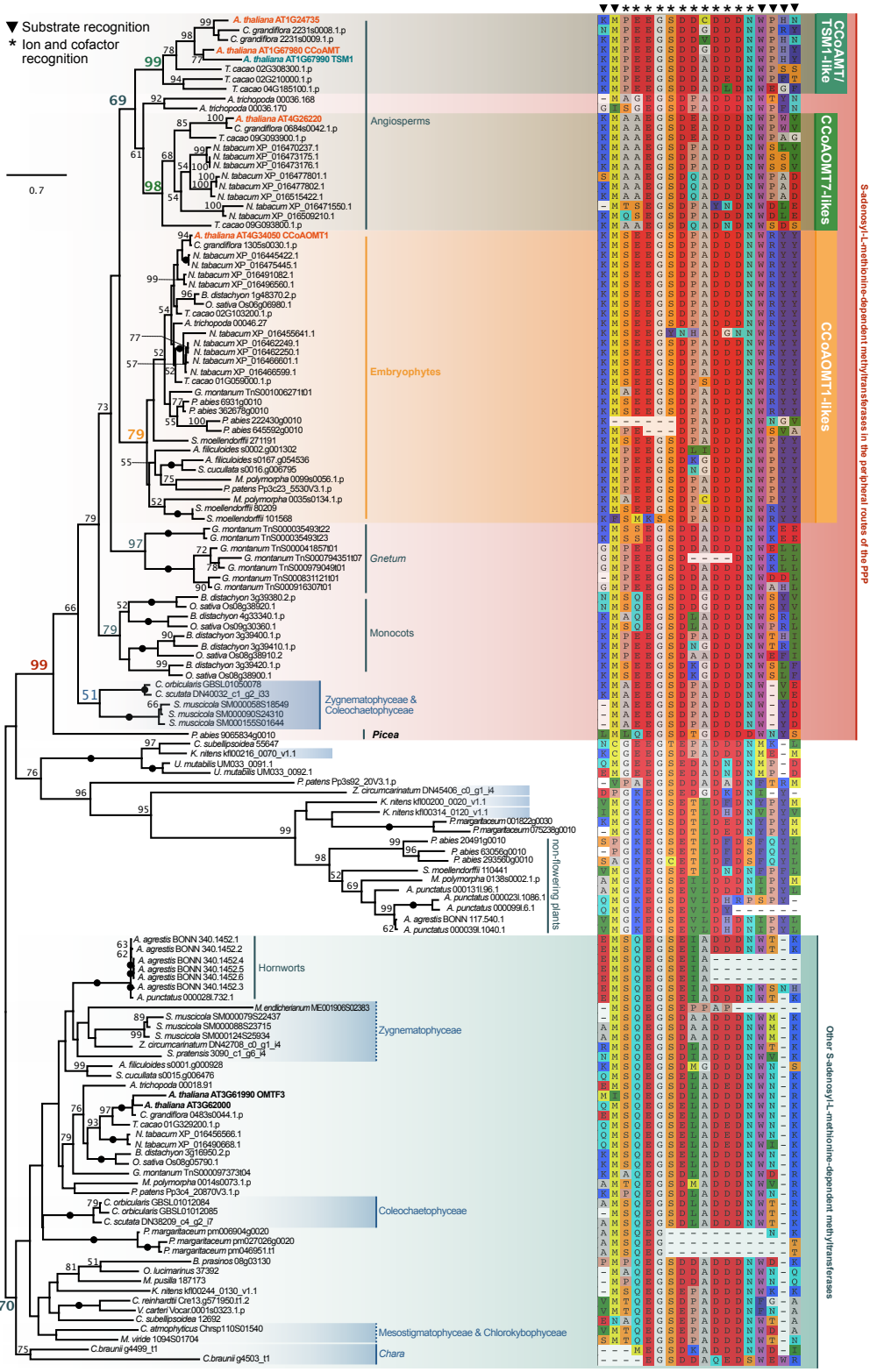
Chlorophyte MAGLs

Streptophyte algal MAGLs

Streptophyte MAGL2/4/3
MAGL2
MAGL4
MAGL13



▼ Substrate recognition
* Ion and cofactor
recognition



Sadenosyl-L-methionine-dependent methyltransferases in the peripheral routes of the PPP

Other Sadenosyl-L-methionine-dependent methyltransferases

0.7

▼

