

1 **Beyond RuBisCO: Convergent molecular evolution of multiple chloroplast**  
2 **genes in C<sub>4</sub> plants**

3  
4  
5 Claudio Casola<sup>1,2</sup>, Jingjia Li<sup>1</sup>

6 <sup>1</sup> Department of Ecology and Conservation Biology, Texas A&M University, College Station,  
7 TX, USA

8 <sup>2</sup> Interdisciplinary Graduate Program in Ecology and Evolutionary Biology, Texas A&M  
9 University, College Station, TX, USA

10  
11  
12 Corresponding author:

13 [claudio.casola@ag.tamu.edu](mailto:claudio.casola@ag.tamu.edu)  
14  
15

16 **Abstract**

17 **Background.** The recurrent evolution of the C<sub>4</sub> photosynthetic pathway in angiosperms  
18 represents one of the most extraordinary examples of convergent evolution of a complex trait.  
19 Comparative genomic analyses have unveiled some of the molecular changes associated with the  
20 C<sub>4</sub> pathway. For instance, several key enzymes involved in the transition from C<sub>3</sub> to C<sub>4</sub>  
21 photosynthesis have been found to share convergent amino acid replacements along C<sub>4</sub> lineages.  
22 However, the extent of convergent replacements potentially associated with the emergence of C<sub>4</sub>  
23 plants remains to be fully assessed. Here, we introduced a robust empirical approach to test  
24 molecular convergence along a phylogeny including multiple C<sub>3</sub> and C<sub>4</sub> taxa. By analyzing  
25 proteins encoded by chloroplast genes, we tested if convergent replacements occurred more  
26 frequently than expected in C<sub>4</sub> lineages compared to C<sub>3</sub> lineages. Furthermore, we sought to  
27 determine if convergent evolution occurred in multiple chloroplast proteins beside the well-  
28 known case of the large RuBisCO subunit encoded by the chloroplast gene *rbcL*.

29 **Methods.** Our study was based on the comparative analysis of 43 C<sub>4</sub> and 21 C<sub>3</sub> grass species  
30 belonging to the PACMAD clade, a focal taxonomic group in many investigations of C<sub>4</sub>  
31 evolution. We first used protein sequences of 67 orthologous chloroplast genes to build an  
32 accurate phylogeny of these species. Then, we inferred amino acid replacements along 13 C<sub>4</sub>  
33 lineages and 9 C<sub>3</sub> lineages using reconstructed protein sequences of their ancestral branches,  
34 corresponding to the most recent common ancestor of each lineage. Pairwise comparisons  
35 between ancestral branches allowed us to identify both convergent and divergent amino acid  
36 replacements between C<sub>4</sub>-C<sub>4</sub>, C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub> lineages.

37 **Results.** The reconstructed phylogenetic tree of 64 PACMAD grasses was characterized by  
38 strong supports in all nodes used for analyses of convergence. We identified 217 convergent  
39 replacements and 201 divergent replacements in 45/67 chloroplast proteins in both C<sub>4</sub> and C<sub>3</sub>  
40 ancestral branches. Pairs of C<sub>4</sub>-C<sub>4</sub> ancestral branches showed higher levels of convergent  
41 replacements than C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub> pairs. Furthermore, we found that more proteins shared  
42 unique convergent replacements in C<sub>4</sub> lineages, with both RbcL and RpoC1 (the RNA  
43 polymerase beta' subunit 1) showing a significantly higher convergent/divergent replacements  
44 ratio in C<sub>4</sub> branches. Notably, significantly more C<sub>4</sub>-C<sub>4</sub> pairs of ancestral branches showed  
45 higher numbers of convergent vs. divergent replacements than C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub> pairs. Our results  
46 demonstrated that, in the PACMAD clade, C<sub>4</sub> grasses experienced higher levels of molecular  
47 convergence than C<sub>3</sub> species across multiple chloroplast genes. These findings have important  
48 implications for both our understanding of the evolution of photosynthesis and the goal of  
49 engineering improved crop varieties that integrates components of the C<sub>4</sub> pathway.

50

51

## 52 **Introduction**

53           Convergent evolution represents the independent acquisition of similar phenotypic traits  
54 in phylogenetically distant organisms. Understanding the genomic changes underlying the  
55 recurrent emergence of phenotypes is a major goal of molecular evolution. The rapidly  
56 increasing taxonomic breadth of genomic resources combined with the development of rigorous  
57 frameworks to comparatively investigate molecular changes has accelerated the pace of  
58 discovery in this area. For instance, substitutions in coding regions of conserved genes have been  
59 implicated in phenotypic changes responsible for adaptation of marine mammals to an aquatic  
60 lifestyle (Foote et al., 2015; Zhou et al., 2015). Other examples of convergent phenotypes whose  
61 molecular underpinnings have been investigated include adaptations in snake and agamid lizard  
62 mitochondria (Castoe et al., 2009), echolocation in mammals (Parker et al., 2013; Thomas and  
63 Hahn, 2015; Zou and Zhang, 2015; Storz, 2016), and hemoglobin function in birds (Natarajan et  
64 al., 2016).

65  
66           Several traits are also known to have convergently evolved in land plants (e.g., Li et al.,  
67 2018; Lü et al., 2018; Preite et al., 2019). One of the most notable examples is represented by the  
68 repeated evolution of the C<sub>4</sub> photosynthetic pathway in flowering plants. The C<sub>4</sub> pathway is a  
69 complex functional adaptation that allows for better photosynthesis efficiency under certain  
70 environmental conditions, such as dry and warm climates, high light intensity, low CO<sub>2</sub>  
71 concentration, and limited availability of nutrients (Knapp and Medina, 1999; Long, 1999). The  
72 C<sub>4</sub> pathway involves cytological, anatomical and metabolic modifications thought to have  
73 evolved multiple times independently in various lineages from the C<sub>3</sub> type (Kellogg, 1999; Sage,  
74 2004; Sage et al., 2011). According to phylogenetic, anatomical and biochemical evidence, the  
75 few slightly different variants of the C<sub>4</sub> photosynthesis type originated more than 60 times in  
76 angiosperms (Sage et al., 2012; Heyduk et al., 2019). In grasses (family Poaceae) alone, the C<sub>4</sub>  
77 pathway has evolved independently ~20 times (Grass Phylogeny Working Group II, 2012).

78  
79           Transitions from C<sub>3</sub> to C<sub>4</sub> plants resulted from genetic changes that include  
80 nonsynonymous substitutions, gene duplications and gene expression alterations (Christin et al.,  
81 2007; Christin et al., 2013a; Christin et al., 2015; Goolsby et al., 2018; Heyduk et al., 2019). It  
82 has been suggested that the evolution of the C<sub>4</sub> pathways proceeded throughout a series of  
83 evolutionary steps wherein the Kranz leaf anatomy typical of this pathway originated first,  
84 followed by changes in the expression patterns of key genes and finally by adaptive  
85 modifications of protein sequences that often result in the convergent emergence of the same  
86 amino acid replacements across C<sub>4</sub> lineages (Christin et al., 2013b; Sage et al., 2012; Williams et  
87 al., 2013).

88  
89           Evidence of convergent changes in proteins associated with photosynthetic processes has  
90 steadily accumulated since genomic data from multiple C<sub>4</sub> lineages have become available in the  
91 past couple of decades. Most of these studies have focused on the ribulose-1,5-bisphosphate

92 carboxylase/oxygenase (RuBisCO), a large multimeric enzyme that catalyzes the carboxylation  
93 of ribulose-1,5-bisphosphate (RuBP), allowing plants to fix atmospheric carbon (Andersson and  
94 Backlund, 2008). RuBisCO also initiates oxygenation of RuBP, which leads to a more limited  
95 production of energy and to loss of carbon in the process of photorespiration (Andersson and  
96 Backlund, 2008; Maurino and Peterhansel, 2010). RuBisCO's limited ability to discriminate  
97 between CO<sub>2</sub> and O<sub>2</sub> has been attributed to the much higher CO<sub>2</sub> to O<sub>2</sub> atmospheric partial  
98 pressure until ~400 million years ago (Sage, 1999, 2004; Sage et al., 2012).  
99

100 Previous studies have revealed multiple convergent amino acid replacements in the large  
101 RuBisCO subunit in C<sub>4</sub> lineages, encoded by the chloroplast gene *rbcL* (Kapralov and Filatov,  
102 2007; Christin et al., 2008; Kapralov et al., 2011; Kapralov et al., 2012; Piot et al., 2018). Some  
103 of these convergent replacements have been associated to positive selection of the corresponding  
104 codons in C<sub>4</sub> monocot and eudicot lineages (Kapralov and Filatov, 2007; Christin et al., 2008;  
105 Kapralov et al., 2012; Piot et al., 2018). Notably, biochemical analyses have demonstrated that  
106 some recurrent amino acid changes in the large RuBisCO subunit of C<sub>4</sub> plants critically alter the  
107 kinetics of RuBisCO, resulting in an accelerated rate of CO<sub>2</sub> fixation at the beginning of the  
108 Calvin-Benson cycle (Studer et al., 2014). Convergent amino acid changes have also been  
109 described in enzymes that are encoded by nuclear genes and play a primary role in the C<sub>4</sub>  
110 pathway, including the phosphoenolpyruvate carboxylase PEPC (Christin et al., 2007; Besnard et  
111 al., 2009), the NADP-malic enzymes NADP-me (Christin et al., 2009a), the  
112 phosphoenolpyruvate carboxykinase PEPCK (Christin et al., 2009b) and the small RuBisCO  
113 subunit (Kapralov et al., 2011).  
114

115 Given the number of biochemical, physiological and anatomical traits that were affected  
116 in each evolutionary transition from C<sub>3</sub> to C<sub>4</sub> photosynthesis (Heyduk et al. 2019), it is likely that  
117 many genes experienced analogous selective pressures across taxa that include C<sub>4</sub> plants. This  
118 could have led to the widespread occurrence of convergent amino acid replacements among a  
119 significant fraction of proteins encoded by genes involved in photosynthesis processes. A recent,  
120 important work has produced the first analysis of convergent replacements across multiple  
121 proteins involved in the metabolism of C<sub>4</sub> and crassulacean acid metabolism (CAM) among  
122 species belonging to the portulugo clade (Caryophyllales). Goolsby and colleagues (2018)  
123 compared evolutionary patterns in 19 gene families with critical roles in metabolic pathways of  
124 both C<sub>4</sub> and CAM plants, also known as carbon-concentration mechanisms (CCMs) genes, and  
125 in 64 non-CCM gene families. They found convergent replacements in proteins from C<sub>4</sub> and  
126 CAM lineages, as well as higher levels of convergent replacements in CCM vs. non-CCM gene  
127 families (Goolsby et al., 2018). Additionally, several amino acid replacements that are prevalent  
128 among C<sub>4</sub> and CAM taxa compared to C<sub>3</sub> lineages were identified in this study (Goolsby et al.,  
129 2018).  
130

131           Altogether, the results of this and other studies demonstrated that convergent molecular  
132 evolution occurred across multiple genes in both C<sub>4</sub> and CAM groups. However, a rigorous  
133 framework to assess the full extent of molecular convergence in C<sub>3</sub> to CCMs transitions has yet  
134 to be presented. For example, analyses of convergent evolution should include null hypotheses  
135 that assume no differences between taxa with and without convergence. In the case of CCMs  
136 evolution, a plausible null hypothesis consists in statistically equivalent numbers of convergent  
137 replacements between C<sub>4</sub> (or CAM) lineages and C<sub>3</sub> lineages.

138  
139           Additionally, nonadaptive replacements should be used to normalize convergent  
140 replacements, in order to account for variation in the rates of nonsynonymous substitutions  
141 across lineages. This approach has been successfully applied in studies of molecular convergent  
142 evolution in vertebrates by assessing both convergent replacements and protein sequence  
143 changes that result in different amino acids, or *divergent replacements* (Castoe et al., 2009;  
144 Thomas and Hahn, 2015; Zou and Zhang, 2015). Furthermore, testing hypotheses about the  
145 extent of convergent molecular evolution remains particularly challenging for many nuclear  
146 genes, because of the prevalence of duplicated copies, particularly in plants (Christin et al., 2007;  
147 Goolsby et al., 2018). Single-copy nuclear or organelle genes allow to more easily recognize  
148 convergent changes and overcome possible confounding compensatory effects due to the  
149 presence of paralogous copies.

150  
151           Given these premises, we sought to test if convergent amino acid changes occur more  
152 frequently in proteins encoded by chloroplast genes in a taxon that includes multiple well-  
153 characterized lineages of C<sub>4</sub> and C<sub>3</sub> grasses. Chloroplast proteins represent an ideal set of targets  
154 to study the role of convergent evolution in C<sub>3</sub> to C<sub>4</sub> transitions for a variety of reasons. First,  
155 most chloroplast proteins are involved in biochemical and biophysical processes that are critical  
156 to photosynthesis. For instance, out of ~75 functionally annotated protein-coding genes in the  
157 maize chloroplast genome, 45 genes are implicated in photosynthesis-related processes,  
158 including *rbcL*, 17 genes coding for subunits of the photosystems I and II (PS I and PS II), 12  
159 genes coding for subunits of the NADH dehydrogenase complex, 6 genes coding for chloroplast  
160 ATPase subunits, 4 genes coding for cytochrome b<sub>6</sub>f complex subunits, and a few more genes  
161 implicated in the assembly of other protein complexes (Maier et al., 1995). Second,  
162 nonannotated orthologous copies of chloroplast genes can be readily identified across plants  
163 through sequence homology searches, taking advantage of the thousands of complete chloroplast  
164 genome sequences currently available for green plants. Third, comparative studies of convergent  
165 evolution in C<sub>4</sub> photosynthesis are facilitated by detailed reconstruction of phylogenetic  
166 relationships within groups with both C<sub>4</sub> and C<sub>3</sub> lineages. Fourth, signatures of positive selection  
167 have been found in multiple chloroplast genes in taxa that contain both C<sub>3</sub> and C<sub>4</sub> plants,  
168 although only the genes *rbcL* and *psaJ*, which encodes a small subunit of the Photosystem I  
169 complex, showed evidence of adaptive changes exclusively in C<sub>4</sub> lineages (Christin et al., 2008;

170 Goolsby et al., 2018; Piot et al., 2018). Finally, most chloroplast genes occur as single copy loci,  
171 as opposed to the multiple paralogs typically present for plant genes encoded in the nucleus.

172

173 In this study, we analyzed 67 chloroplast genes from 64 grass species, including 43 C<sub>4</sub>  
174 and 19 C<sub>3</sub> species belonging to the PACMAD clade, named after six of its most representative  
175 subfamilies: Panicoideae, Arundinoideae, Chloridoideae, Micrairoideae, Aristidoideae and  
176 Danthonioideae. Using published phylogenetic information, we identified thirteen independent  
177 C<sub>3</sub> to C<sub>4</sub> transitions in this group of species. We applied a series of tests based on convergent vs.  
178 divergent amino acid replacements and determined that convergent molecular evolution occurred  
179 at a higher rate in chloroplast genes of C<sub>4</sub> lineages compared to C<sub>3</sub> lineages, a pattern that  
180 remained largely unchanged after excluding the RbcL protein from the convergence analyses.  
181 Our findings suggest that the evolutionary trajectories of multiple chloroplast genes have been  
182 remarkably affected during the emergence of the C<sub>4</sub> adaptation in the PACMAD clade, a result  
183 that has significant implications for our understanding of C<sub>4</sub> photosynthesis evolution and  
184 organelle-nucleus interactions, and for the identification of molecular changes that might be  
185 critical to the successful development of engineered C<sub>3</sub> crops that incorporate carbon-  
186 concentration mechanisms.

187

188

## 189 **Methods**

190

### 191 ***Data source and filtering***

192 We queried NCBI GenBank (Sayers et al., 2019) for complete chloroplast genome  
193 sequences of grass species that were included in phylogenetic analyses by the Grass Phylogeny  
194 Working Group II (2012) and downloaded the corresponding coding sequences. Each species  
195 was assigned to either C<sub>3</sub> or C<sub>4</sub> type following the results of the Grass Phylogeny Working  
196 Group II (2012). Additionally, we downloaded the coding chloroplast sequences for  
197 *Dichanthelium acuminatum*, *Thyridolepis xerophila*, *Sartidia dewinteri* and *Sartidia perrieri* (C<sub>3</sub>  
198 species) (Brown and Smith, 1972; Smith and Brown, 1973; Hattersley and Stone, 1986;  
199 Hattersley et al., 1986; Besnard et al., 2014). We used the standalone blastn ver.  
200 2.2.29+(Camacho et al., 2009) with the Expect value (E) cutoff of 1e-10 to determine putative  
201 sequence orthology with coding sequences of the *Zea mays* chloroplast genes (Maier et al.,  
202 1995). Single copy putative orthologs that were present in more than 95% of the species were  
203 retained for further analysis (Table S1).

204

### 205 ***Multiple sequence alignment***

206 We aligned the individual sequences using TranslatorX ver. 1.1 (Abascal et al., 2010),  
207 and further adjusted the alignments manually using BioEdit ver. 7.0.9.0 (Hall, 1999). Stop  
208 codons and sites that could not be aligned unambiguously were removed.

209



210 ***Phylogeny reconstruction***

211 We concatenated the individual sequence alignments and extracted third codon position  
212 sites for phylogeny reconstruction. We ran PartitionFinder ver. 1.1.1 (Lanfear et al., 2012) to  
213 identify the best partitioning scheme (partitioning by gene) for the downstream analysis using  
214 both Akaike information criterion (AIC) (Akaike, 1973) and Bayesian information criterion  
215 (BIC) (Schwarz, 1978). We then used maximum likelihood framework as implemented in  
216 RAxML ver. 8.2.10 (Stamatakis, 2014) to reconstruct the phylogeny. Branch support was  
217 estimated using 1,000 bootstrap replicates. *Oryza sativa* and *Brachypodium distachyon* from the  
218 BOP (Bambusoideae, Oryzoideae and Pooideae) clade were used as outgroup, whereas all  
219 ingroup species belonged to the PACMAD clade. We used FigTree ver. 1.4.0 (Rambaut, 2012)  
220 to rearrange and visualize the phylogeny, and the figures were edited further to improve  
221 readability and to indicate C<sub>4</sub>/C<sub>3</sub> classification.

222

223 ***Ancestral state reconstruction***

224 We reconstructed ancestral states at each phylogenetic node for each individual gene  
225 using the program codeml from the software package PAML ver. 4.9a (Yang, 2007) and the  
226 basic codon substitution model (model = 0, NSsites = 0).

227

228 ***Inference of convergent and divergent replacements***

229 We extracted the reconstructed ancestral states from the codeml output. The  
230 corresponding amino acid sequences were then compared to investigate individual site changes  
231 along selected branches in the reconstructed phylogenetic tree in the context of emergence of the  
232 C<sub>4</sub> trait. For each group of species descendant from a single C<sub>4</sub> ancestor, we chose the branch  
233 between the most recent C<sub>3</sub> ancestor and the most ancestral C<sub>4</sub> node, i.e., the branch along which  
234 the C<sub>4</sub> adaptation presumably emerged (referred to as “C<sub>4</sub> ancestral branch” throughout this  
235 article, see Figs. 1 and 2). For C<sub>3</sub> species, we chose the most ancestral branch that did not share  
236 ancestry with any C<sub>4</sub> lineage (“C<sub>3</sub> ancestral branch”, see Figs. 1 and 2). In either case, if only a  
237 single species was available in a given lineage, that terminal branch was used. The outgroup  
238 species (*O. sativa* and *B. distachyon*) were not included in this analysis (Fig. 1).

239

240 We searched for amino acid changes that occurred along pairs of ancestral branches.  
241 Replacements in both branches that resulted in the same state at a given site in the two  
242 descendants were considered convergent, regardless of whether the corresponding ancestral  
243 states of ancestral were the same or different (Castoe et al., 2009). Likewise, two replacements  
244 were considered divergent if states at the descendant orthologous sites were different, regardless  
245 of the corresponding ancestral states (Castoe et al., 2009). Although two orthologous sites, by  
246 definition, descend from one ancestral site, the actual state transitions, as well as their number,  
247 between ancestral and descendant states along a given branch are not known because the states  
248 are reconstructed only at discrete time steps (i.e., at selected nodes) and represent only those

249 specific evolutionary time stamps. Therefore, potential intermediate stages, including a transient  
250 convergent phase, would remain undetected.

251  
252 We identified putative convergent and divergent amino acid changes in each gene  
253 product individually. We summarized those data within each of the three categories: (1) two C<sub>4</sub>  
254 ancestral branches (C<sub>4</sub>-C<sub>4</sub>), (2) C<sub>3</sub> ancestral branch and C<sub>4</sub> ancestral branch (C<sub>3</sub>-C<sub>4</sub>), and (3) two  
255 C<sub>3</sub> ancestral branches (C<sub>3</sub>-C<sub>3</sub>). the Boschloo's statistical exact unconditional test (Boschloo,  
256 1970) was performed was performed to test the significance of the convergent replacement  
257 excess when comparing two of the three photosynthesis type pairs using the SciPy library ver.  
258 1.7.1 in python3 (Virtanen et al. 2020).

#### 259 **Data availability**

261 Raw data, including alignments, fasta sequences, and phylogenetic analyses data, are  
262 available through the following Figshare repository:  
263 [https://figshare.com/articles/dataset/Convergence-chloroplast-genes-C4-Casola-Li-](https://figshare.com/articles/dataset/Convergence-chloroplast-genes-C4-Casola-Li-2021/15180690)  
264 [2021/15180690](https://figshare.com/articles/dataset/Convergence-chloroplast-genes-C4-Casola-Li-2021/15180690).

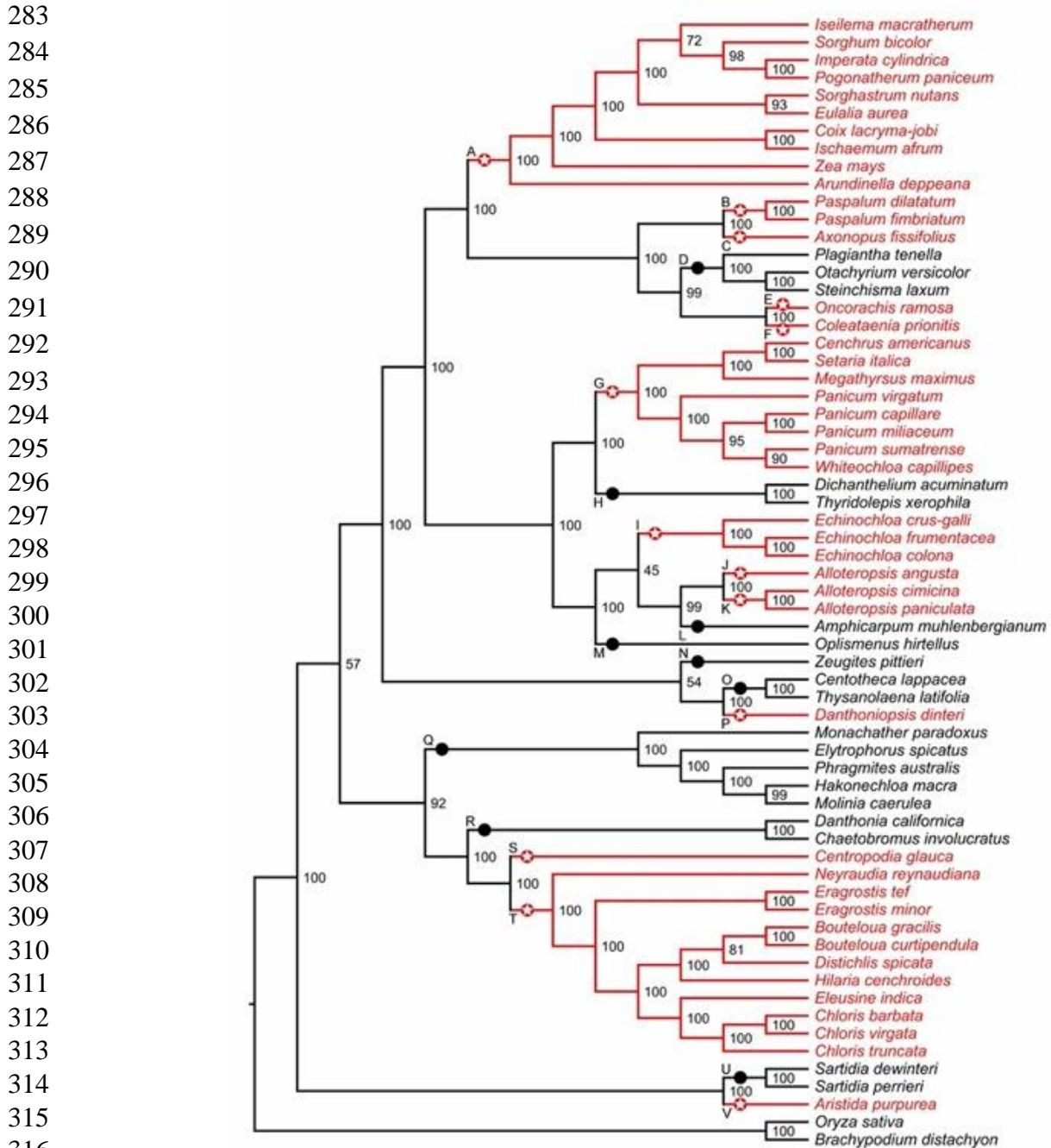
#### 265 266 267 **Results**

##### 268 ***Phylogeny reconstructions***

270 We examined 63 grass chloroplast genomes to identify gene orthologs for *Zea mays*  
271 chloroplast genes and extracted the corresponding coding and protein sequences. The resulting  
272 dataset included up to 67 DNA/protein sequences in 64 grass species that were retained for  
273 further analysis (Table S1). One to four sequences were absent in thirteen species. Out of 64  
274 species, 43 were classified as C<sub>4</sub> and 21 (including two outgroup species) as C<sub>3</sub>. The reconstructed  
275 phylogeny is well supported, except for three branches with low to moderate bootstrap values,  
276 and it is consistent for both AIC and BIC (Fig. 1 and Figs. S1-S3). We identified thirteen C<sub>4</sub>  
277 ancestral branches that represent putative C<sub>3</sub> to C<sub>4</sub> transitions, and nine C<sub>3</sub> ancestral branches  
278 (Fig. 1).

279  
280  
281  
282





317 **Figure 1. Phylogenetic relationships among 64 C<sub>4</sub> and C<sub>3</sub> grass species.**

318 The phylogeny tree was obtained using RAxML (GTR+Γ model) based on the third codon position sites  
319 in 67 chloroplast genes. The partitioning scheme was selected according to Akaike information criterion  
320 (AIC). C<sub>4</sub> and C<sub>3</sub> ancestral branches are shown in red and black, respectively. Red stars and black circles  
321 (labels A-V) indicate C<sub>4</sub> and C<sub>3</sub> ancestral branches, respectively. Numbers represent bootstrap support.

322  
323  
324  
325

326  
327

### 328 ***Convergent and divergent amino acid replacements across chloroplast proteins***

329 We assessed the level of molecular convergence in C<sub>3</sub> to C<sub>4</sub> transitions by quantifying  
330 convergent and divergent amino acid replacements across the PACMAD phylogeny by  
331 performing pairwise comparisons of reconstructed ancestral sequences (Figs. 2 and 3, Table S2;  
332 see Methods).

333  
334  
335

336  
337

338  
339  
340  
341

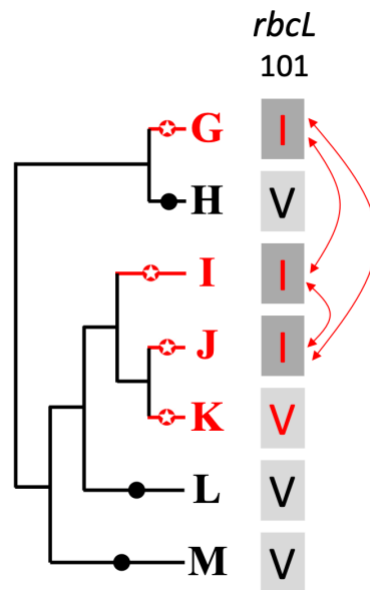
342  
343

344  
345

346  
347

348  
349

350  
351



352  
353

**Figure 2. Example of ancestral C<sub>4</sub> and C<sub>3</sub> ancestral branches and convergent changes in C<sub>4</sub> ancestral branches.**

354  
355

Pairwise comparisons were carried out between the ancestral branches (in this example, G to M; see also Fig. 1). The C<sub>4</sub> branches G, I and J shared a convergent V→I amino acid replacement (red arrows) at position 101 of the large RuBisCO subunit encoded by the gene *rbcL*. The ancestral amino acid and the convergently changed amino acid are highlighted in light and dark gray, respectively. C<sub>4</sub> and C<sub>3</sub> ancestral branches are shown in red and black, respectively.

357  
358  
359  
360

361  
362

A total of 217 sites showed at least one convergent replacement: 104 in C<sub>4</sub>-C<sub>4</sub>, 120 in C<sub>3</sub>-C<sub>4</sub> and 34 in C<sub>3</sub>-C<sub>3</sub> pairs. A further 201 sites exhibited one or more divergent replacements: 96 in C<sub>4</sub>-C<sub>4</sub>, 121 in C<sub>3</sub>-C<sub>4</sub>, and 39 in C<sub>3</sub>-C<sub>3</sub> pairs (Table 1). The difference in convergent-divergent site distributions between the three photosynthesis types was not statistically significant ( $P \geq 0.05$ , Boschloo's test; Table 1).

363  
364  
365  
366  
367

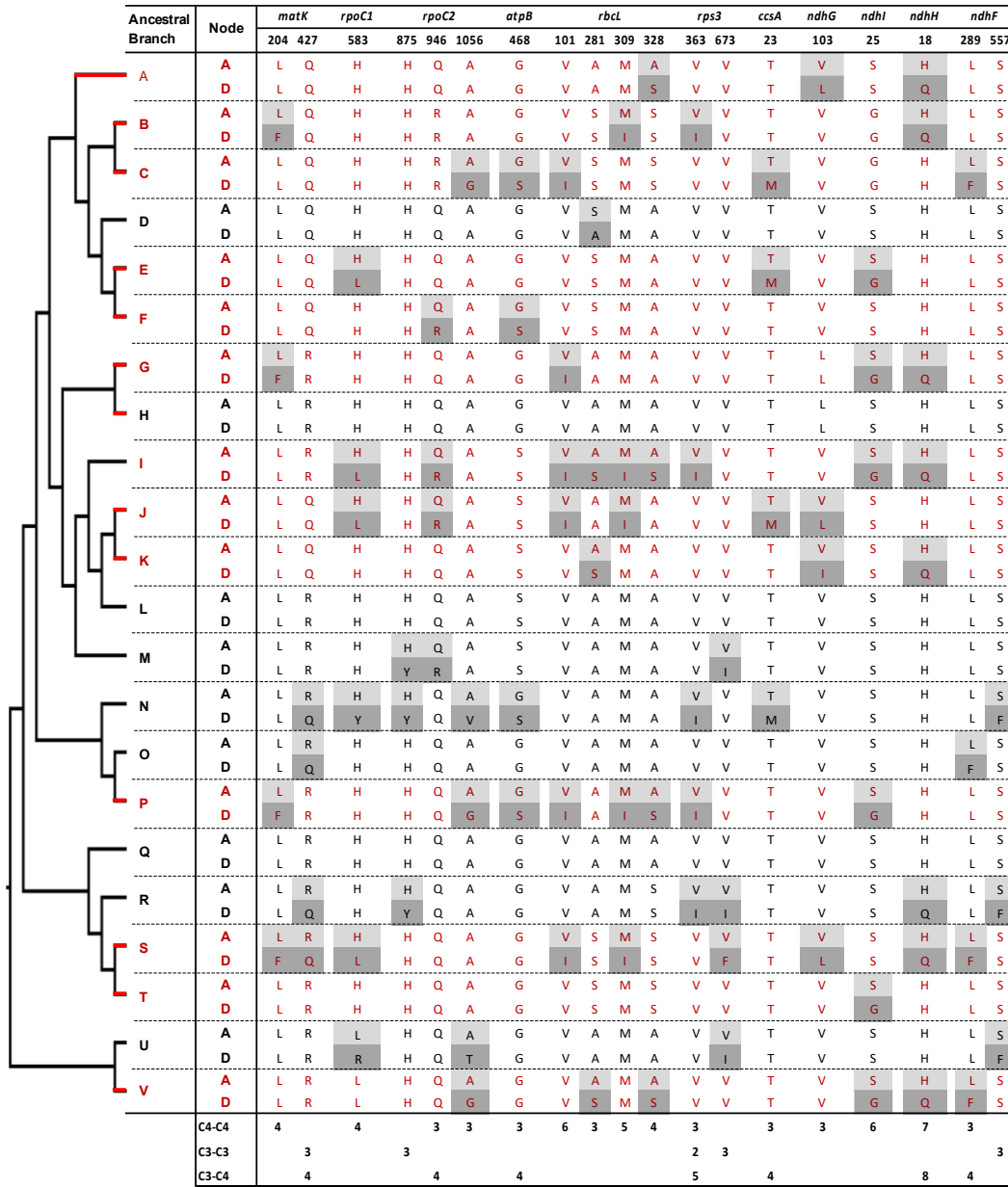
368  
369  
370  
  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390

**Table 1. Numbers of amino acid sites and genes with convergent and divergent replacements in ancestral branch comparisons.**

	C4-C4			C3-C4			C3-C3		
	Con	Div	Ratio	Con	Div	Ratio	Con	Div	Ratio
Sites	104	96	1.08	120	121	0.99	34	39	0.87
Sites*	80	64	1.25	82	69	1.19	17	16	1.06
Genes	24	23	1.04	26	32	0.81	13	17	0.76
Genes*	24	20	1.2	25	29	0.86	9	10	0.9

Comparisons were made between pairs of C4-C4, C3-C3 and C3-C4 branches. Numbers of replacements unique to a given category (\*), and the corresponding ratios *Con:Div (Ratio)*. Differences between the C3-C3 and C4-C4 categories are not statistically significant ( $P \geq 0.05$ , Boschloo's test). Con: convergent. Div: divergent.

Among the C4 ancestral branches, several individual sites showed high contrast in the number of branches involved in convergent and divergent replacements (Fig. 3, Tables S2 and S3). For example, seven C4 branches (54%) shared the H18Q replacement in the product of *ndhH*, with no divergent replacements. Six, five, and four C4 branches (46%, 38%, and 31%) showed convergent replacements at three sites in the RbcL protein (V101I, M309I, and A328S, respectively). Furthermore, six C4 branches shared the S25G replacement in the product of *ndhI* and four L204F changes in the protein encoded by *matK*. In all these cases, there were no other convergent or divergent replacements in C3-C3 or C3-C4 branch comparisons, except for one H18Q change in NdhH in a C3-C3 branch. Two sites with convergent replacements in the proteins encoded by *ndhF* (L557F) and *rpoC2* (H875Y) were found uniquely in C3-C3 pairs, and only one site in the protein Rps3 showed convergence independently in C4-C4 and C3-C3 pairs (Fig. 3).



426 **Figure 3. Amino acid replacements shared by at least three C<sub>4</sub> or C<sub>3</sub> reference branches.**  
 427 Ancestral (A) and derived (D) amino acids at replacement sites are shown. Site numbers correspond to the  
 428 *Zea mays* orthologous sequence annotation. Red and black letters and branches represent C<sub>4</sub> and C<sub>3</sub>  
 429 ancestral branches, respectively (see also Figs. 1 and 2).

430  
 431  
 432  
 433  
 434  
 435

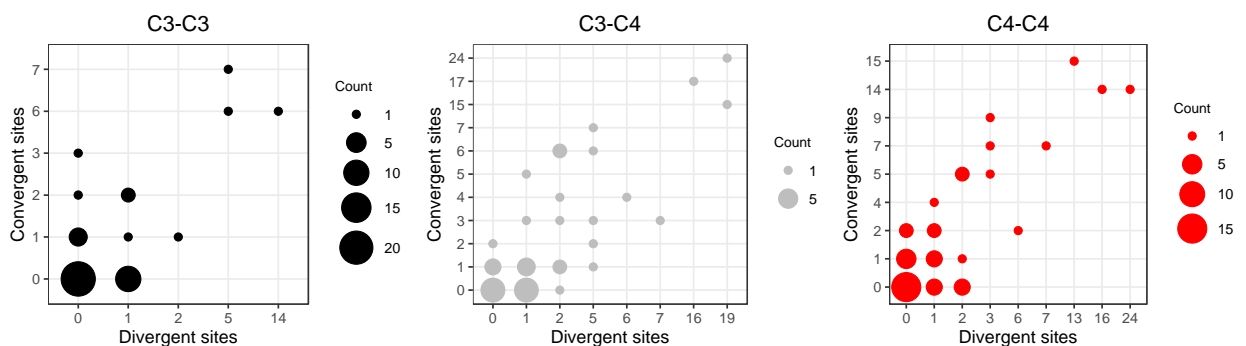
436 We then searched for convergent replacements that occurred along more than two C<sub>4</sub>  
437 branches at sites that remained otherwise conserved in C<sub>3</sub> and C<sub>4</sub> lineages, arguing that such  
438 changes could result from selective pressure rather than drift. We identified fourteen C<sub>4</sub>-specific  
439 convergent sites in proteins from 8 genes: *atpB*, *ccsA*, *matK*, *ndhF*, *ndhH*, *ndhI*, *rbcL* and *rpoC2*  
440 (Table S3). Five of these sites were found in RbcL, whereas two sites were identified in both  
441 NdhF and NdhI.

442  
443

#### 444 ***Molecular convergence in individual chloroplast proteins***

445 Convergent and divergent amino acid replacements were detected in the products of 45  
446 chloroplast genes, thirteen of which had at least one site with four or more replacements (Fig. 4,  
447 Table 1 and Table S2). Twenty-four genes had convergent changes in C<sub>4</sub>-C<sub>4</sub>, 26 in C<sub>3</sub>-C<sub>4</sub>, and 13  
448 in C<sub>3</sub>-C<sub>3</sub> types of pairs (Table 1). Although the convergent/divergent replacement ratio was  
449 higher in C<sub>4</sub>-C<sub>4</sub> pairs than C<sub>3</sub>-C<sub>4</sub> and C<sub>3</sub>-C<sub>3</sub> pairs, the differences between the three  
450 photosynthesis types was not statistically significant ( $P \geq 0.05$ , Boschloo's test; Table 1). The  
451 lack of replacements was the single most common state for chloroplast proteins across  
452 photosynthesis types; however, in C<sub>4</sub>-C<sub>4</sub> there were more genes with a higher number convergent  
453 vs. divergent replacements (Fig. 4 and Table S4).

454  
455

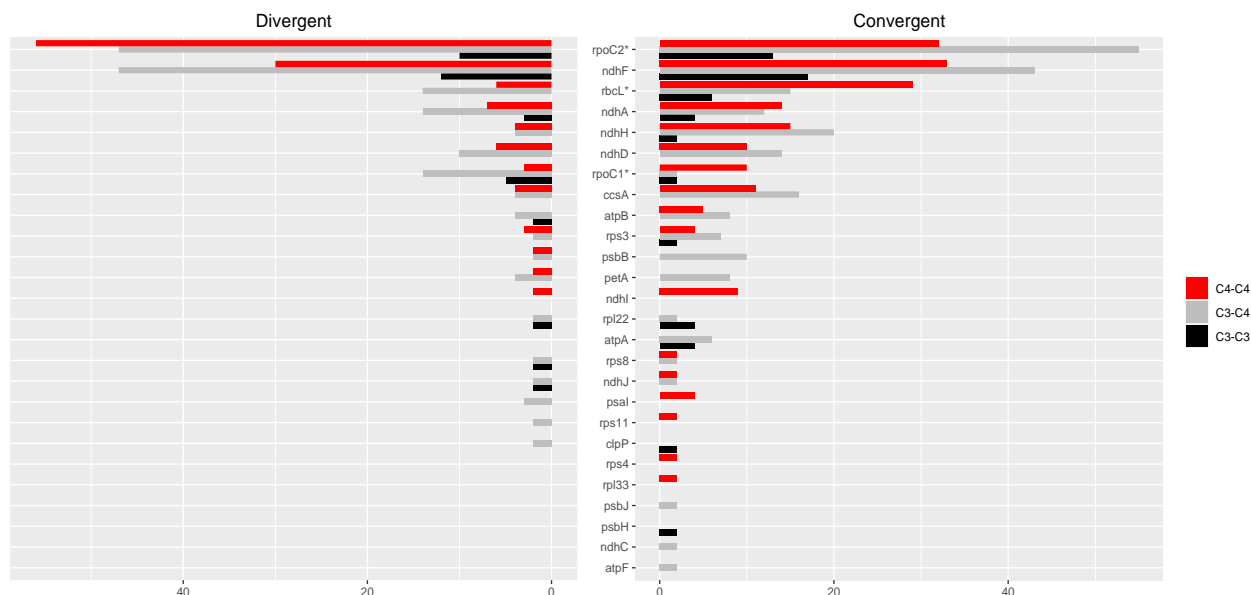


456  
457 **Figure 4. Distribution of convergent and divergent amino acid replacements in pairs of ancestral**  
458 **branches.**

459  
460

461 Overall, 26 proteins showed a higher number of convergent vs. replacement sites, of  
462 which 16, 13 and 10 were found in C<sub>4</sub>-C<sub>4</sub>, C<sub>3</sub>-C<sub>4</sub> and C<sub>3</sub>-C<sub>3</sub> pairs, respectively (Fig. 5 and Table  
463 S4). We found statistically significant differences in the number of replacements between C<sub>4</sub>-C<sub>4</sub>  
464 and C<sub>3</sub>-C<sub>4</sub> pairs, but not C<sub>3</sub>-C<sub>3</sub> pairs, in the products of the genes *rbcL*, *rpoC1* and *rpoC2* ( $P <$   
465  $0.05$ , Boschloo's test; Table S4). In RbcL and RpoC1, C<sub>4</sub>-C<sub>4</sub> pairs shared much higher number of  
466 convergent replacements, whereas the opposite was true in RpoC2. RpoC1 was also the only  
467 protein showing more convergent than divergent replacements in C<sub>4</sub>-C<sub>4</sub> pairs compared to C<sub>3</sub>-C<sub>3</sub>  
468 and C<sub>3</sub>-C<sub>4</sub> pairs. In C<sub>4</sub>-C<sub>4</sub> pairs, RpoC1 shared 4 convergent and 1 divergent replacement,

469 compared to 1 and 2 in C<sub>3</sub>-C<sub>3</sub> pairs and 1 and 5 in C<sub>3</sub>-C<sub>4</sub> pairs, respectively. Additionally, the  
470 proteins NdhG, NdhI, PsaI, RpoA, Rps4 and Rps11 exhibited convergent replacements only in  
471 C<sub>4</sub>-C<sub>4</sub> pairs (Table S4). When considering the number of affected sites rather than the number of  
472 replacements, no genes showed a significantly different pattern between photosynthesis types ( $P$   
473  $\geq 0.05$ , Boschloo's test; Table S4).  
474  
475



476  
477 **Figure 5. Amino acid replacements in chloroplast proteins with more convergent than divergent**  
478 **changes in at least one photosynthesis type.**

479 Twenty-six chloroplast proteins with more convergent than divergent changes in C<sub>4</sub>-C<sub>4</sub>, C<sub>3</sub>-C<sub>4</sub> and/or C<sub>3</sub>-  
480 C<sub>3</sub> pairs. Asterisks indicate proteins with significantly different replacements between C<sub>4</sub>-C<sub>4</sub> and C<sub>3</sub>-C<sub>4</sub>  
481 pairs.  
482  
483

484 The proteins encoded by *matK*, *rpoC2* and *ndhF* shared much higher numbers of both  
485 convergent and divergent replacements than other chloroplast proteins across all photosynthesis  
486 type comparisons (Table S4). Both *matK* and *ndhF* are known to be rapidly evolving and have  
487 been consistently used in low taxonomic level phylogenetic studies in flowering plants (Barthet  
488 and Hilu, 2008; Patterson and Givnish, 2002). The gene *rpoC2* has also been recently described  
489 as a useful phylogenetic marker in angiosperms (Walker et al., 2019).  
490

#### 491 ***Molecular convergence across ancestral branches***

492 The comparison of ancestral branch pairs with convergent and divergent replacements  
493 revealed remarkable differences between photosynthesis types. Overall, C<sub>4</sub>-C<sub>4</sub> pairs of ancestral  
494 branches showed a distribution skewed toward more convergent and divergent replacements than  
495 the two other categories (Fig. 6). There were significantly fewer pairs of C<sub>4</sub>-C<sub>4</sub> ancestral  
496 branches with no replacements and with no convergent replacements than C<sub>3</sub>-C<sub>4</sub> and C<sub>3</sub>-C<sub>3</sub> pairs



497 ( $P < 0.05$ , Boschloo's test; Table 2). Conversely, significantly more C<sub>4</sub>-C<sub>4</sub> pairs shared more  
 498 convergent than divergent replacements, and at least two convergent changes compared to C<sub>3</sub>-C<sub>4</sub>  
 499 and C<sub>3</sub>-C<sub>3</sub> pairs ( $P < 0.05$ , Boschloo's test; Table 2).

500 No significant difference was observed between pairs of C<sub>3</sub>-C<sub>4</sub> and pairs of C<sub>3</sub>-C<sub>3</sub>. We  
 501 found identical patterns when the same analyses were performed after excluding all replacements  
 502 in the RbcL protein, except for the lack of a significant difference between C<sub>4</sub>-C<sub>4</sub> and C<sub>3</sub>-C<sub>3</sub> in  
 503 the proportion of pairs with divergent replacements and pairs with more convergent than  
 504 divergent changes (Table S6).

505  
 506

507 **Table 2. Number of ancestral branches with convergent and divergent replacements.**

	C <sub>4</sub> -C <sub>4</sub>	C <sub>3</sub> -C <sub>4</sub>	C <sub>3</sub> -C <sub>3</sub>
No replacements	6 (.08)	30 (.26)	12 (.33)
No Con	12 (.15)	48 (.41)	16 (.44)
w/Con	66 (.85)	69 (.59)	20 (.56)
w/Div	63 (.81)	67 (.57)	18 (.50)
Con>Div	40 (.51)	36 (.31)	10 (.28)
Con>1	49 (.63)	39 (.33)	8 (.22)

508 Comparisons were made between pairs of C<sub>4</sub>-C<sub>4</sub>, C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub> branches. Proportions of pairs of ancestral  
 509 branches over all branches by category are shown in parenthesis. The total number of pairs of ancestral  
 510 branches are 78, 36 and 117 for C<sub>4</sub>-C<sub>4</sub>, C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub> comparisons, respectively. All comparisons between C<sub>4</sub>-C<sub>4</sub> pairs and both  
 511 C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub> pairs were statistically significantly different ( $P < 0.05$ , Boschloo's test). No comparison between  
 512 C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub> pairs was statistically significant ( $P \geq 0.05$ , Boschloo's test). Con: convergent. Div: divergent.  
 513 Con>Div: pairs of branches with more convergent than divergent replacements. Con>1: pairs of branches with more  
 514 than one convergent replacement.

515

516

517

518

519

520

521

522

523

524

525

526

527

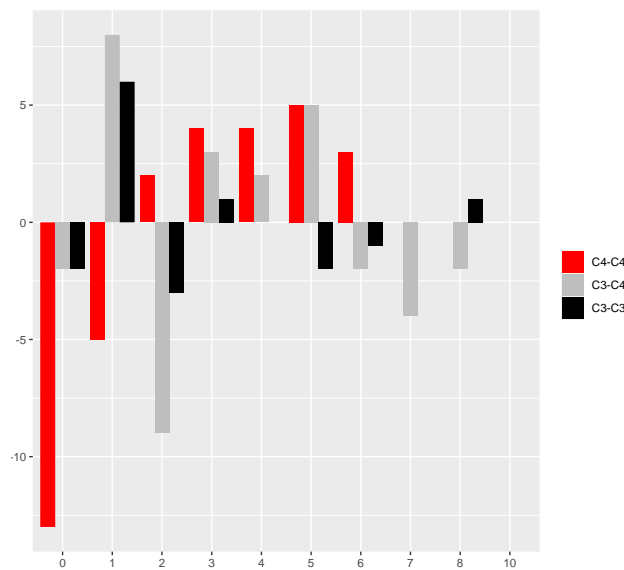
528

529

530

531

532



533 **Figure 6. Pairs of ancestral branches by convergent and divergent replacements.**

534 Difference in the number of pairs of ancestral branches for convergent and divergent categories (0-8 and  
 535 10 replacements).

536

537 ***Distribution of amino acid replacements across PACMAD lineages***

538 Convergent and divergent replacements were preferentially found in specific pairs of  
 539 ancestral branches. In C<sub>4</sub> pairs, convergent sites were most abundant between *Danthoniopsis*  
 540 *dinteri* and *Aristida purpurea* (ten sites, branches P and V in Fig. 1), whereas divergent sites  
 541 were most common between *Centropodia glauca* and *Aristida purpurea* (ten sites, branches S  
 542 and V in Fig. 1). In pairwise C<sub>3</sub> branch comparisons, most convergent sites were identified  
 543 between both *Zeugites pittieri* and Danthonieae (branches N and R in Fig. 1) and Danthonieae  
 544 and *Sartidia* spp. (branches R and U in Fig. 1), whereas the most divergent site-rich pair was  
 545 formed by *Zeugites pittieri* and *Sartidia* spp. (eight sites, branches N and U in Fig. 1; Table S5).

546

547 ***Molecular convergence in the RuBisCO large subunit***

548 We further inspected the evolution of the RuBisCO large subunit across the PACMAD  
 549 clade. A total of 4 out of 9 RbcL amino acids with convergent changes in C<sub>4</sub> ancestral  
 550 branches—V101I, A281S, M309I and A328S—have been identified in previous studies on  
 551 PACMAD grasses (Christin et al., 2008; Piot et al., 2018) as sites that experienced adaptive  
 552 evolution in C<sub>4</sub> species (Table 3).

553

554

555 **Table 3. Summary of RbcL amino acid sites with signatures of convergent evolution or positive**  
 556 **selection.**

Codon	Ancestral AA	Convergent Change/p.s.s.	#Convergent a.b.
<b>10</b>	S	G	2
<b>93</b>	E	D	2
<b>94</b>	A	P	2
<b>101</b> <sup>*†</sup>	V	I	6
142 <sup>*†</sup>	P	Several	na
<b>143</b>	T	A	3
145 <sup>*†</sup>	S	A/V	na
258 <sup>*</sup>	R	K	na
270 <sup>*</sup>	L	I	na
<b>281</b> <sup>*†</sup>	A	S	3
282 <sup>†</sup>	H	Several	na
<b>309</b> <sup>*†</sup>	M	I	5
<b>328</b> <sup>*†</sup>	A	S	4
<b>461</b> <sup>*</sup>	V	I	2
468 <sup>†</sup>	E	D	na
471 <sup>†</sup>	E	Several	na
476 <sup>†</sup>	I	L/V	na

557 Ancestral AA: ancestral amino acid. Convergent change/p.s.s.: derived amino acid in multiple C<sub>4</sub> ancestral branches  
 558 and positively selected sites from previous studies. #Convergent a.b.: number of ancestral branches with convergent  
 559 changes. Boldface: sites with convergent changes identified in this study. Asterisk: positively selected sites in  
 560 PACMAD C<sub>4</sub> lineages from Christin et al. (2008). Dagger: positively selected sites in PACMAD C<sub>4</sub> lineages from  
 561 Piot et al. (2018).

562  
563 A further site, T143A, was found to evolve under positive selection in C<sub>3</sub> to C<sub>4</sub> transitions  
564 in monocots (Studer et al., 2014). Interestingly, an adaptive S143A replacement has also been  
565 detected in the gymnosperm *Podocarpus* (Sen et al., 2011). Three more sites with convergent  
566 replacements—at positions 93, 94 and 461—correspond to amino acids that were reported to  
567 evolve under positive selection in different groups of seed plants by Kapralov and Filatov  
568 (2007). Thus, all of the *rbcL* codons that appear to have evolved convergently among the  
569 PACMAD C<sub>4</sub> lineages we have examined are also known to have experienced adaptive evolution  
570 in seed plants, but not all of them have been shown to evolve adaptively in C<sub>4</sub> grasses.

571  
572  
573

## 574 Discussion

575 The recurrent emergence of carbon-concentration mechanisms (CCMs) across multiple  
576 angiosperm clades in the past 35 million years represents one of the most striking examples of  
577 convergent evolution of a complex phenotypic trait. Several investigations have shown that the  
578 phenotypic parallelism across C<sub>4</sub> lineages is to some extent mirrored by convergent changes in  
579 the sequence of proteins with key metabolic roles in the biochemistry of C<sub>4</sub> photosynthesis, both  
580 in monocots and eudicots (Christin et al., 2007; Besnard et al., 2009, Christin et al., 2009a,  
581 Christin et al., 2009b, Kapralov et al., 2011, Goolsby et al., 2018). Furthermore, biochemical  
582 analyses have determined that some of these changes reflect adaptive shifts, as in the case of the  
583 increased availability of CO<sub>2</sub> at the RuBisCO site (Studer et al., 2014). Further evidence of  
584 changes in the selective pressure associated to the C<sub>3</sub> to C<sub>4</sub> transitions have emerged from the  
585 detection of several positively selected sites in multiple genes associated with photosynthetic  
586 processes (Christin et al., 2008; Studer et al., 2014; Goolsby et al., 2018; Piot et al., 2018). These  
587 and other discoveries have paved the way to a more nuanced understanding of the molecular  
588 basis of phenotypic convergence in CCM plants and may accelerate the development of crop  
589 varieties with augmented resistance to high temperature and low water availability.

590  
591 For these aims to be fully realized, a robust framework to assess the extent and  
592 phenotypic impact of convergent molecular changes is necessary. Along the lines of strategies  
593 applied in vertebrates research (Castoe et al., 2009, Thomas and Hahn, 2015), we presented here  
594 the results of a novel methodological approach to the study of molecular convergence in C<sub>4</sub>  
595 grasses. We investigated patterns of convergent and divergent amino acid changes in nearly 70  
596 chloroplast proteins across multiple C<sub>4</sub> and C<sub>3</sub> lineages in the PACMAD clade, with the goal of  
597 testing a specific hypothesis: is the evolution of chloroplast proteins showing stronger signatures  
598 of convergent amino acid replacements in C<sub>4</sub> lineages compared to C<sub>3</sub> lineages? This analysis  
599 also allowed us to establish if proteins other than enzymes involved in the CCM biochemistry  
600 underwent parallel amino acid changes in C<sub>4</sub> lineages. Our reasoning is that many proteins

601 expressed in the chloroplast could have experienced similar selective pressure across multiple C<sub>3</sub>  
602 to C<sub>4</sub> transitions and might have accumulated convergence replacements as a result.

603

604 We based our analysis on the identification of amino acid replacements shared by pairs of  
605 ancestral C<sub>4</sub> branches, defined here as branches corresponding to C<sub>3</sub> to C<sub>4</sub> transitions in the  
606 PACMAD phylogeny. We compared these changes to those identified in ancestral C<sub>3</sub> branches,  
607 namely all C<sub>3</sub> lineages that include only C<sub>3</sub> species (Figs. 1 and 2), and to changes found between  
608 ancestral C<sub>3</sub> and C<sub>4</sub> branches. For each of the three possible pairs of photosynthesis types C<sub>4</sub>-C<sub>4</sub>,  
609 C<sub>3</sub>-C<sub>4</sub> and C<sub>3</sub>-C<sub>3</sub>, we determined the number of amino acid sites, genes and pairs of ancestral  
610 branches with convergent replacements. We detected signatures of convergent evolution in all  
611 types of datasets. First, we identified many individual replacements that emerged repeatedly and  
612 uniquely in C<sub>4</sub> ancestral branches, particularly in the proteins RbcL, NdhH, NdhI and MatK. We  
613 also observed C<sub>3</sub>-specific convergent replacements in NdhF and RpoC2, and a case of multiple  
614 C<sub>4</sub> and C<sub>3</sub> convergent changes in Rps3. Additionally, we identified 8 chloroplast genes with one  
615 or more C<sub>4</sub>-specific convergent sites. Second, we found evidence of significantly higher rates of  
616 convergent replacements in C<sub>4</sub> lineages in both RbcL and RpoC1, and several convergent  
617 replacements that occurred exclusively in C<sub>4</sub>-C<sub>4</sub> pairs in proteins encoded by *ndhG*, *ndhI*, *psaI*,  
618 *rpoA*, *rps4* and *rps11*. These genes are involved in a variety of biological processes in the  
619 chloroplast, from the cyclic electron transport in (*ndhG* and *ndhI*) and the stabilization of (*psaI*)  
620 the photosystem I, to transcription (*rpoA* and *rpoC1*), translation (*rps4* and *rps11*) and CO<sub>2</sub>  
621 fixation (*rbcL*). Third, we identified statistically significant differences in pairs of C<sub>4</sub> branches  
622 with convergent replacements (Table 2). Crucially, we observed more pairs with higher  
623 convergent than divergent replacements in C<sub>4</sub>-C<sub>4</sub> compared to both C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub>, even after  
624 removing replacements identified in the RuBisCO large subunit, RbcL.

625

626 Altogether, these findings suggest that multiple biochemical processes occurring in the  
627 chloroplast might have experienced recurrent adaptive changes associated with the emergence of  
628 C<sub>4</sub> photosynthesis. Notably, some of these proteins are not directly involved in the light-  
629 dependent or light-independent reactions of the photosynthesis, implying that processes such as  
630 the regulation of gene expression and protein synthesis in the chloroplast are also experiencing  
631 significant selective pressures during the transition from C<sub>3</sub> to C<sub>4</sub> plants. These results should  
632 motivate further studies to determine the prevalence of convergent amino acid replacements due  
633 transitions to CCMs among the thousands of proteins encoded by nuclear genes but expressed in  
634 the chloroplast (Jarvis and López-Juez, 2013). Although such analyses are currently hindered by  
635 the limited number of sequenced nuclear genomes in taxa with multiple C<sub>3</sub> and C<sub>4</sub> lineages,  
636 including the PACMAD clade, genome-wide investigations of convergent replacements will be  
637 possible in the near future given the current pace of DNA sequencing in plants.

638

639 A further important conclusion drawn from these results is that convergent replacements  
are not uncommon between C<sub>3</sub>-C<sub>3</sub> and C<sub>3</sub>-C<sub>4</sub> lineages. This is possibly due to some

640 environmental factors affecting the evolution of chloroplast genes that are shared across grass  
641 lineages regardless of their photosynthesis type.

642

643 The analysis of individual convergent replacements in the RuBisCO large subunit both  
644 confirmed previous findings and highlighted novel potentially adaptive changes among  
645 PACMAD species. Importantly, these novel convergent replacements are known to evolve under  
646 positive selection in non-PACMAD seed plants. This underscores the potential of our approach  
647 to identify novel changes with functional significance in the transition to CCMs in grasses, as  
648 opposed to standard statistical tests of positive selection. Alternatively, some RbcL sites could  
649 experience convergence across a variety of seed plants because of selective pressure other than  
650 those associated with C<sub>3</sub> to C<sub>4</sub> transitions.

651

652 Overall, our results are robust to several possible confounding factors. First, we analyzed  
653 branches that are strongly supported in our phylogeny reconstruction. The phylogenetic tree built  
654 using the 67 chloroplast genes is well supported, with the exception of three branches with fairly  
655 low bootstrap support. However, all three branches are short and have minimal impact upon our  
656 conclusions regarding C<sub>4</sub> evolution (Fig. 1 and Figs. S1-S3). Moreover, the tree is largely  
657 consistent with a comprehensive recent study of 250 grasses based on complete plastome data  
658 (Saarela et al., 2018). Second, by focusing only on ancestral branches and ignoring amino acid  
659 replacements that may have occurred after the divergence of species within a given C<sub>4</sub> clade, our  
660 strategy provided a conservative estimate of the number of convergent changes that could have  
661 occurred during the evolution of PACMAD grasses. Third, we eliminated genes with possible  
662 paralogous copies, which could have introduced false positive replacements.

663

664 We recognize some potential caveats in our approach. By relying on a relatively small  
665 sample of PACMAD species, our statistical power to detect signatures of convergent evolution  
666 was limited. Increasing the number of ancestral C<sub>4</sub> and C<sub>3</sub> lineages should provide a broader  
667 representation of convergent replacements in C<sub>4</sub> clades. Furthermore, we applied a strict  
668 definition of convergence that ignores changes to amino acids with similar chemical properties.  
669 We think that a conservative approach was necessary given that amino acids with similar  
670 chemical properties might have a very different functional effect on protein activity given their  
671 size and tridimensional interactions with nearby residues. Third, we assumed that all the  
672 observed convergent replacements were the result of convergent phenotypic changes, which fall  
673 under the general category of homoplasy (Avice and Robinson, 2008). However, some of these  
674 replacements could instead represent hemiplasy, or character state changes due to introgression  
675 between different C<sub>4</sub> lineages, incomplete lineage sorting (ILS) of ancestral alleles or horizontal  
676 gene transfer (Avice and Robinson, 2008). Recombination between chloroplast genomes, which  
677 is required for introgression to occur, has been documented but appears to be rare (Carbonell-  
678 Caballero et al., 2015, Greiner and Sobanski, 2015, Sancho et al., 2018). Introgression or  
679 horizontal gene transfer between congeneric species has been associated to the acquisition of part

680 of the C<sub>4</sub> biochemical pathway in the PACMAD genus *Alloteropsis* (Christin et al., 2012;  
681 Olofsson et al., 2016). However, these transfers were limited to a few nuclear genes. Moreover,  
682 only a very few cases of horizontal transfer between chloroplast genomes have been reported in  
683 plants (Stegemann et al., 2012). Therefore, the contribution of hemiplasy to the observed pattern  
684 of convergent replacements in C<sub>4</sub> lineages is likely to be minimal. Finally, we treated C<sub>4</sub> species  
685 regardless of their photosynthesis subtype (NAPD-ME, NAD-ME and PEPCK), which is known  
686 to vary among PACMAD subfamilies (Taylor et al., 2010). We argue that our results are  
687 conservative with regard to this aspect because convergent replacements should be expected to  
688 occur more often between C<sub>4</sub> groups sharing the same photosynthesis subtype.

689

690

## 691 **Conclusions**

692

693 In this study, we showed that molecular convergent evolution in the form of recurrent amino acid  
694 replacements affected multiple chloroplast proteins in C<sub>4</sub> lineages of the PACMAD clade of  
695 grasses. This finding significantly broadened the number of genes known to have evolved  
696 convergently in C<sub>4</sub> species. We observed for the first time that genes not directly involved in  
697 photosynthesis-related processes experienced convergent changes, suggesting that future efforts  
698 should rely whenever possible on genome-wide analyses of amino acid changes rather than focus  
699 primarily on candidate key metabolic genes, similarly to previous investigations on gene  
700 expression patterns in C<sub>4</sub> and CAM plants. Our methodological approach based on the  
701 comparison of convergent and divergent replacements among photosynthesis types underscores  
702 the importance of a more rigorous hypothesis-based testing of convergent evolution signatures in  
703 C<sub>4</sub> plant evolution. Our results should inform more nuanced approaches to introduce CCM-like  
704 processes in C<sub>3</sub> crops.

705

706

## 707 **Acknowledgements**

708 The project was supported by the National Institute of Food and Agriculture, U.S. Department of  
709 Agriculture, under award number TEX0-1-9599, the Texas A&M AgriLife Research, and the  
710 Texas A&M Forest Service.

711

712

## 713 **References**

- 714 Abascal, F., R. Zardoya, and M. J. Telford. 2010. TranslatorX: multiple alignment of nucleotide  
715 sequences guided by amino acid translations. *Nucleic Acids Res* 38: W7–W13.
- 716 Akaike, H. 1973. Information theory and an extension of the maximum likelihood principle.  
717 Proceedings of the Second International Symposium on Information Theory: 267–281.
- 718 Andersson, I., and A. Backlund. 2008. Structure and function of Rubisco. *Plant Physiol Biochem*  
719 46: 275-291.



- 720 Avise, J. C., and T. J. Robinson. 2008. Hemiplasy: a new term in the lexicon of phylogenetics.  
721 *SystBiol* 57: 503–507.
- 722 Barthet, K M.M., W. Hilu. 2008. Evaluating evolutionary constraint on the rapidly evolving gene  
723 matK using protein composition. *J Mol Evol* 66: 85-97.
- 724 Besnard, G., M. Muasya, F. Russier, E. H. Roalson, N. Salamin, P.-A. Christin. 2009.  
725 Phylogenomics of C<sub>4</sub> Photosynthesis in Sedges (Cyperaceae): Multiple Appearances and  
726 Genetic Convergence. *Mol Biol Evol* 26: 1909–1919.
- 727 Besnard, G., P.-A. Christin, P.-J. G. Malé, E. Lhuillier, C. Lauzeral, E. Coissac, and M. S.  
728 Vorontsova. 2014. From museums to genomics: old herbarium specimens shed light on a  
729 C<sub>3</sub> to C<sub>4</sub> transition. *J Exp Bot* 65: 6711–6721.
- 730 Boschloo, R. D. 1970. Raised conditional level of significance for the 2 × 2-table when testing  
731 the equality of two probabilities. *Stat Neerl* 24: 1-9.
- 732 Brown, W. V., and B. N. Smith. 1972. Grass evolution, the Kranz syndrome, <sup>13</sup>C/<sup>12</sup>C ratios, and  
733 continental drift. *Nature* 239: 345–346.
- 734 Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, and T. L. Madden.  
735 2009. BLAST+: architecture and applications. *BMC Bioinform* 10: 421.
- 736 Carbonell-Caballero, J., R. Alonso, V. Ibañez, J. Terol, M. Talon, J. Dopazo. 2015. A  
737 phylogenetic analysis of 34 chloroplast genomes elucidates the relationships between  
738 wild and domestic species within the genus Citrus. *Mol Biol Evol* 32: 2015–2035.
- 739 Castoe, T. A., A. P. J. d. Koning, H.-M. Kim, W. Gu, B. P. Noonan, G. Naylor, Z. J. Jiang, et al.  
740 2009. Evidence for an ancient adaptive episode of convergent molecular evolution. *Proc*  
741 *Natl Acad Sci USA* 106: 8986–8991.
- 742 Christin, P.-A., M. Arakaki, C. P. Osborne, and E. J. Edwards. 2015. Genetic enablers  
743 underlying the clustered evolutionary origins of C<sub>4</sub> photosynthesis in angiosperms. *Mol*  
744 *Biol Evol* 32: 846–858.
- 745 Christin, P.-A., S. F. Boxall, R. Gregory, E. J. Edwards, J. Hartwell, and C. P. Osborne. 2013a.  
746 Parallel recruitment of multiple genes into C<sub>4</sub> photosynthesis. *Genome Biol Evol* 5: 2174–  
747 2187.
- 748 Christin, P.-A., E. J. Edwards, G. Besnard, S. F. Boxall, R. Gregory, E. A. Kellogg, J. Hartwell,  
749 C.P. Osborne. 2012. Adaptive evolution of C(4) photosynthesis through recurrent lateral  
750 gene transfer. *Curr Biol* 22: 445-9.
- 751 Christin, P.-A., C. P. Osborne, D. S. Chatelet, J. T. Columbus, G. Besnard, T. R. Hodkinson, L.  
752 M. Garrison, et al. 2013b. Anatomical enablers and the evolu- tion of C<sub>4</sub> photosynthesis  
753 in grasses. *Proc Natl Acad Sci USA* 110: 1381–1386.
- 754 Christin P.-A., B. Petitpierre, N. Salamin, L. Büchi, G. Besnard. 2009b. Evolution of C(4)  
755 phosphoenolpyruvate carboxykinase in grasses, from genotype to phenotype. *Mol*  
756 *Biol Evol* 26: 357-65.
- 757 Christin P.-A., E. Samaritani, B. Petitpierre B, N. Salamin, G. Besnard. 2009a. Evolutionary  
758 insights on C<sub>4</sub> photosynthetic subtypes in grasses from genomics and  
759 phylogenetics. *Genome Biol Evol* 1:221-30.

- 760 Christin, P.-A., N. Salamin, A. M. Muasya, E. H. Roalson, F. Russier, and G. Besnard. 2008.  
761 Evolutionary switch and genetic convergence on *rbcL* following the evolution of C<sub>4</sub>  
762 photosynthesis. *Mol Biol Evol* 25: 2361–2368.
- 763 Christin, P.-A., N. Salamin, V. Savolainen, M. R. Duvall, and G. Besnard. 2007. C<sub>4</sub>  
764 photosynthesis evolved in grasses via parallel adaptive genetic changes. *Curr Biol* 17:  
765 1241-1247.
- 766 Foote, A. D., Y. Liu, G. W. C. Thomas, T. Vinař, J. Alföldi, J. Deng, S. Dugan, et al. 2015.  
767 Convergent evolution of the genomes of marine mammals. *Nat Genet* 47: 272–275.
- 768 Goolsby, E. W., A. J. Moore, L. P. Hancock, J. M. De Vos, and E. J. Edwards. 2018. Molecular  
769 evolution of key metabolic genes during transitions to C<sub>4</sub> and CAM photosynthesis. *Am J*  
770 *Bot* 105: 602–613.
- 771 Greiner S., J. Sobanski, R. Bock. 2015. Why are most organelle genomes transmitted  
772 maternally? *Bioessays* 37: 80-94.
- 773 Grass Phylogeny Working Group II. 2012. New grass phylogeny resolves deep evolutionary  
774 relationships and discovers C<sub>4</sub> origins. *New Phytol* 193: 304–312.
- 775 Hall, T. A. 1999. BioEdit: a user-friendly biological sequence alignment editor and analysis  
776 program for Windows 95/98/NT. *Nucl Acids Symp Series* 41: 95-98.
- 777 Hattersley, P. W., and N. E. Stone. 1986. Photosynthetic enzyme activities in the C<sub>3</sub>-C<sub>4</sub>  
778 intermediate *Neurachne minor* S. T. Blake (Poaceae). *Aust J Plant Physiol* 13: 399-408.
- 779 Hattersley, P. W., S.-C. Wong, S. Perry, and Z. Roksandic. 1986. Comparative ultrastructure and  
780 gas exchange characteristics of the C<sub>3</sub>-C<sub>4</sub> intermediate *Neurachne minor* S. T. Blake  
781 (Poaceae). *Plant Cell Environ* 9: 217-233.
- 782 Heyduk, K., J. J. Moreno-Villena, I. S. Gilman, P.-A. Christin, and E. J. Edwards. 2019. The  
783 genetics of convergent evolution: insights from plant photosynthesis. *Nat Rev Genet* 20:  
784 485–493.
- 785 Jarvis, P., E. López-Juez. 2013. Biogenesis and homeostasis of chloroplasts and other  
786 plastids. *Nat Rev Mol Cell Biol* 14: 787-802.
- 787 Kapralov, M. V., and D. A. Filatov. 2007. Widespread positive selection in the photosynthetic  
788 Rubisco enzyme. *BMC Evol Biol* 7: 73.
- 789 Kapralov, M. V., J. A. C. Smith, and D. A. Filatov. 2012. Rubisco evolution in C<sub>4</sub> eudicots: An  
790 analysis of Amaranthaceae *sensu lato*. *PLoS One* 7: e52974.
- 791 Kapralov, M. V., D. S. Kubien, I. Andersson, and D. A. Filatov. 2011. Changes in rubisco  
792 kinetics during the evolution of C<sub>4</sub> photosynthesis in *Flaveria* (Asteraceae) are associated  
793 with positive selection on genes encoding the enzyme. *Mol Biol Evol* 28: 1491–1503.
- 794 Kellogg, E. A. 1999. Phylogenetic Aspects of the Evolution of C<sub>4</sub> Photosynthesis. In R. F. Sage  
795 and R. K. Monson [eds.], C<sub>4</sub> Plant Biology, 411-444. Academic Press, San Diego,  
796 California, USA.
- 797 Knapp, A. K., and E. Medina. 1999. Success of C<sub>4</sub> Photosynthesis in the Field: Lessons from  
798 Communities Dominated by C<sub>4</sub> Plants. In R. F. Sage and R. K. Monson [eds.], C<sub>4</sub> Plant  
799 Biology, 251-283. Academic Press, San Diego, California, USA.

- 800 Lanfear, R., B. Calcott, S. Y. W. Ho, and S. Guindon. 2012. PartitionFinder: Combined selection  
801 of partitioning schemes and substitution models for phylogenetic analyses. *Mol Biol Evol*  
802 29: 1695–1701.
- 803 Li, B., C. Förster, C. A. M. Robert, T. Züst, L. Hu, R. A. R. Machado, J.-D. Berset, et al. 2018.  
804 Convergent evolution of a metabolic switch between aphid and caterpillar resistance in  
805 cereals. *Sci Adv* 4: eaat6797.
- 806 Long, S. P. 1999. Environmental Responses. In R. F. Sage and R. K. Monson [eds.], *C<sub>4</sub> Plant*  
807 *Biology*, 215-249. Academic Press, San Diego, California, USA.
- 808 Lü, P., S. Yu, N. Zhu, Y.-R. Chen, B. Zhou, Y. Pan, D. Tzeng, et al. 2018. Genome encode  
809 analyses reveal the basis of convergent evolution of fleshy fruit ripening. *Nat Plants* 4:  
810 784–791.
- 811 Maier R. M., K. Neckermann, G. L. Igloi, H. Kössel. 1995. Complete sequence of the maize  
812 chloroplast genome: gene content, hotspots of divergence and fine tuning of  
813 genetic information by transcript editing. *J Mol Biol* 251: 614-28.
- 814 Maurino, V. G., and C. Peterhansel. 2010. Photorespiration: current status and approaches for  
815 metabolic engineering. *Curr Opin Plant Biol* 13: 248-255.
- 816 Natarajan, C., F. G. Hoffmann, R. E. Weber, A. Fago, C. C. Witt, and J. F. Storz. 2016.  
817 Predictable convergence in hemoglobin function has unpredictable molecular  
818 underpinnings. *Science* 354: 336-339.
- 819 Olofsson, J.K., M. Bianconi, G. Besnard, L. T. Dunning, M. R. Lundgren, H. Holota, M. S.  
820 Vorontsova, O. Hidalgo, I. L. Leitch, P. Nosil, C. P. Osborne, P.-A. Christin. 2016.  
821 Genome biogeography reveals the intraspecific spread of adaptive mutations for a  
822 complex trait. *Mol Ecol* 25: 6107-6123.
- 823 Patterson, T. B., T. J. Givnish. 2002. Phylogeny, concerted convergence, and phylogenetic niche  
824 conservatism in the core Liliales: insights from *rbcL* and *ndhF* sequence data. *Evolution*  
825 56: 233–252.
- 826 Parker, J., G. Tsagkogeorga, J. A. Cotton, Y. Liu, P. Provero, E. Stupka, and S. J. Rossiter. 2013.  
827 Genome-wide signatures of convergent evolution in echolocating mammals. *Nature* 502:  
828 228–231.
- 829 Piot, A., J. Hackel, P.-A. Christin, and G. Besnard. 2018. One-third of the plastid genes evolved  
830 under positive selection in PACMAD grasses. *Planta* 247: 255–266.
- 831 Preite, V., C. Sailer, L. Syllwasschy, S. Bray, H. Ahmadi, U. Krämer, and L. Yant. 2019.  
832 Convergent evolution in *Arabidopsis halleri* and *Arabidopsis arenosa* on calamine  
833 metalliferous soils. *Philos Trans R Soc Lond B Biol Sci* 374: 20180243.
- 834 Rambaut, A. 2012. FigTree. Tree Figure Drawing Tool, version 1.4.0. website:  
835 <http://tree.bio.ed.ac.uk/software/figtree/>.
- 836 Rosenblum, E. B., C. E. Parent, and E. E. Brandt. 2014. The molecular basis of phenotypic  
837 convergence. *Annu Rev Ecol Evol Syst* 45: 203-226.

- 838 Saarela, J. M., S. V. Burke, W. P. Wysocki, M. D. Barrett, L. G. Clark, J. M. Craine, P. M.  
839 Peterson, et al. 2018. A 250 plastome phylogeny of the grass family (Poaceae):  
840 topological support under different data partitions. *PeerJ* 6: e4299.
- 841 Sage, R. F. 1999. Why C<sub>4</sub> Photosynthesis? In R. F. Sage and R. K. Monson [eds.], *C<sub>4</sub> Plant*  
842 *Biology*, 3-16. Academic Press, San Diego, California, USA.
- 843 Sage, R. F. 2004. The evolution of C<sub>4</sub> photosynthesis. *New Phytol* 161: 341–370.
- 844 Sage, R. F., P.-A. Christin, and E. J. Edwards. 2011. The C<sub>4</sub> plant lineages of planet Earth. *J Exp*  
845 *Bot* 62: 3155–3169.
- 846 Sage, R. F., T. L. Sage, and F. Kocacinar. 2012. Photorespiration and the evolution of C<sub>4</sub>  
847 photosynthesis. *Annu Rev Plant Biol* 63: 19-47.
- 848 Sancho, R., C.P. Cantalapiedra, D. López-Alvarez, S. P. Gordon, J. P. Vogel, P. Catalán,  
849 B. Contreras-Moreira. 2018. Comparative plastome genomics and phylogenomics of  
850 *Brachypodium*: flowering time signatures, introgression and recombination in  
851 recently diverged ecotypes. *New Phytol* 218:1631-1644.
- 852 Sayers, E. W., M. Cavanaugh, K. Clark, J. Ostell, K. D. Pruitt, and I. Karsch-Mizrachi. 2019.  
853 GenBank. *Nucleic Acids Res* 47: D94–D99.
- 854 Schwarz, G. 1978. Estimating the dimension of a model. *The Annals of Statistics* 6: 461-464.
- 855 Sen L, Fares MA, Liang B, Gao L, Wang B, Wang T, Su YJ. 2011. Molecular evolution  
856 of *rbcL* in three gymnosperm families: identifying adaptive and coevolutionary  
857 patterns. *Biol Direct* 6:29.
- 858 Smith, B. N., and W. V. Brown. 1973. The Kranz syndrome in the Gramineae as indicated by  
859 carbon isotopic ratios. *American Journal of Botany* 60: 505-513.
- 860 Stegemann, S., M. Keuthe, S. Greiner, R. Bock. 2012. Horizontal transfer of chloroplast  
861 genomes between plant species. *Proc Natl Acad Sci U S A* 109:2434-8.
- 862 Stamatakis, A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of  
863 large phylogenies. *Bioinformatics* 30: 1312–1313.
- 864 Storz, J. F. 2016. Causes of molecular convergence and parallelism in protein evolution. *Nature*  
865 *Reviews Genetics* 17: 239-250.
- 866 Studer, R. A., P.-A. Christin, M. A. Williams, and C. A. Orengo. 2014. Stability-activity  
867 tradeoffs constrain the adaptive evolution of RubisCO. *Proc Natl Acad Sci U S A* 111:  
868 2223-2228.
- 869 Taylor, S. H., S. P. Hulme, M. Rees, B. S. Ripley, F. I. Woodward, and C. P. Osborne. 2010.  
870 Ecophysiological traits in C<sub>3</sub> and C<sub>4</sub> grasses: a phylogenetically controlled screening  
871 experiment. *New Phytol* 185: 780–791.
- 872 Thomas, G. W. C., and M. W. Hahn. 2015. Determining the null model for detecting adaptive  
873 convergence from genomic data: A case study using echolocating mammals. *Mol Biol*  
874 *Evol* 32: 1232–1236.
- 875 Virtanen, P., R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski,  
876 P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett M, J. Wilson, K. J.  
877 Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, I. Polat,

- 878 Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen,  
879 E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van  
880 Mulbregt; SciPy 1.0 Contributors. 2020. SciPy 1.0: fundamental algorithms for scientific  
881 computing in Python. *Nat Methods* 17: 261-272.
- 882 Walker J. F., N. Walker-Hale, O. M. Vargas, D. A. Larson, G. W. Stull. 2019. Characterizing  
883 gene tree conflict in plastome-inferred phylogenies. *PeerJ* 7:e7747.
- 884 Williams, B. P., I. G. Johnston, S. Covshoff, and J. M. Hibberd. 2013. Phenotypic landscape  
885 inference reveals multiple evolutionary paths to C4 photosynthesis. *eLife* 2: e00961.
- 886 Yang, Z. 2007. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 24:  
887 1586–1591.
- 888 Zhou, X., I. Seim, and V. N. Gladyshev. 2015. Convergent evolution of marine mammals is  
889 associated with distinct substitutions in common genes. *Sci Rep* 5: 16550.
- 890 Zou, Z., and J. Zhang. 2015. No genome-wide protein sequence convergence for echolocation.  
891 *Mol Biol Evol* 32: 1237–1241.