# Born with intronless ERF transcriptional factors: C4 photosynthesis inherits a legacy dating back 450 million years

Ming-Ju Amy Lyu[1,8], Huilong Du[2,3,8], Hongyan Yao[4,8], Zhiguo Zhang[5,8], Genyun Chen[1], Faming Chen[1], Yong-Yao Zhao[1], Qiming Tang[1], Fenfen Miao[1], Yanjie Wang[1], Yuhui Zhao[2], Hongwei Lu[2], Lu Fang[2], QiangGao[2], Yiying Qi[6], QingZhang[6], Jisen Zhang[6], Tao Yang[7], Xuean Cui[5], Chengzhi Liang[2,3$], Tiegang Lu[5$], Xin-Guang Zhu[1,9$]

1. State Key Laboratory of Plant Molecular Genetics, Center of Excellence for Molecular Plant Sciences, Chinese Academy of Sciences, Shanghai, China, 200032
2. State Key Laboratory of Plant Genomics, Institute of Genetics and Developmental Biology, Innovation Academy for Seed Design, Chinese Academy of Sciences, Beijing, China;
3. University of Chinese Academy of Sciences, Beijing, China; School of Life Sciences, Institute of Life Sciences and Green Development, Hebei University, Baoding, China.
4.State Key Laboratory of Genetic Engineering, School of Life Sciences, Fudan University, Shanghai 200438, China
5. Biotechnology Research Institute/National Key Facility for Gene Resources and Gene Improvement, Chinese Academy of Agricultural Sciences, Beijing, 100081, China
6. Center for Genomics and Biotechnology, Fujian Provincial Key Laboratory of Haixia Applied Plant Systems Biology, Key Laboratory of Sugarcane Biology and Genetic Breeding, National Engineering Research Center for Sugarcane, College of Life Sciences, Fujian Agriculture and Forestry University, Fuzhou, China
7. China National GeneBank, Shenzhen, 518120, China.
8. These authors contributed equally
9. Lead Contact
* Correspondence: zhuxg@cemps.ac.cn (X.G.Z), lutiegang@caas.cn (L.T.), cliang@genetics.ac.cn (C.L.)

## Summary

The genus *Flaveria*, containing species at different evolutionary stages of the progression from $C_3$ to $C_4$ photosynthesis, is used as a model system to study the evolution of $C_4$ photosynthesis. Here, we report chromosome-scale genome sequences for five *Flaveria* species, including $C_3$, $C_4$, and intermediate species. Our analyses revealed that both acquiring additional gene copies and recruiting ethylene responsive factor (ERF) *cis*-regulatory elements (CREs) contributed to the emergence of $C_4$ photosynthesis. ERF transcriptional factors (TFs), especially intronless ERF TFs, were co-opted in dicotyledonous $C_4$ species and monocotyledonous $C_4$ species in parallel. These $C_4$ species co-opted intronless ERF TFs originated from the Late Ordovician mass extinction that occurred ~450 million years ago in coping with environmental stress. Therefore, this study demonstrated that intronless ERF TFs were acquired during the early evolution of plants and provided the molecular toolbox facilitating multiple subsequent independent evolutions of $C_4$ photosynthesis.

**Key words**: *Flaveria* genome, $C_4$ photosynthesis, Tandem duplication, Intronless ERF TFs

## Introduction

$C_4$ photosynthesis is a complex trait that evolved from ancestral $C_3$ types in the last 35 million years (Sage, 2004; Sage et al., 2012). With high light, water, and nitrogen use efficiencies (Vogan and Sage, 2011; Zhu et al., 2008), $C_4$ photosynthesis is an ideal target to be engineered into $C_3$ crops to increase crop yield (Long et al., 2015; Maurino and Weber, 2013; Zhu et al., 2010). Compared to $C_3$ photosynthesis, $C_4$ photosynthesis dedicates more genes to carbon fixation, and these genes are compartmentalized either in mesophyll cells (MCs) or in bundle sheath cells (BSCs) (Hatch, 1987; Slack and Hatch, 1967). These MCs and BSCs form the specialized $C_4$ "Kranz anatomy" (Hatch, 1987). Therefore, the evolution of $C_4$ photosynthesis requires modifications of both metabolism and leaf anatomy in $C_3$ ancestors. Though complex, $C_4$ photosynthesis has evolved independently more than 70 times in angiosperms, making it an excellent example of convergent evolution of a complex trait (Sage, 2016). How such a complex trait emerges repeatedly remains an unresolved question in biological research.

All genes that function in $C_4$ photosynthesis have counterparts in $C_3$ species (Christin et al., 2013; Christin et al., 2009; Moreno-Villena et al., 2018; Williams et al., 2012). The same $C_4$ orthologous genes that show relatively high transcript abundances were co-opted in different $C_4$ lineages in parallel (Emms et al., 2016; Moreno-Villena et al., 2018). Moreover, the recruited $C_4$ genes often adopt pre-existing regulatory mechanisms of photosynthesis (Burgess et al., 2016), which enables the coordinated expression of $C_4$ cycle genes with other photosynthesis-related genes. This recruitment of pre-existing elements provides a mechanism for the repeated emergence of $C_4$ photosynthesis in independent lineages. In available genome sequences, conserved *cis*-regulatory elements (CREs) as well as transcription factors (TFs) have been identified that control the MCs or BSCs specificity of $C_4$ genes in different monocotyledonous $C_4$ lineages (Burgess et al., 2019; Gupta et al., 2020; John et al., 2014). Moreover, conserved TFs controlling MCs and BSCs specificity of gene expression were also identified between monocotyledonous and dicotyledonous $C_4$ species (Aubry et al.,

79　　2014). These observations nevertheless raise the questions of when and how these

80　　shared regulatory mechanisms were co-opted into the different $C_4$ species that diverged

81　　160 million years ago (mya) (Kumar et al., 2017), which is much earlier than the

82　　emergence of $C_4$ photosynthesis, *i.e.*, 35 mya (Sage et al., 2011).

83　　　　Among the dicotyledonous model systems for $C_4$ photosynthesis, the genus

84　　*Flaveria* is remarkable because it contains $C_3$, $C_4$, and many intermediate species

85　　(Powell, 1978). During the last decades, studies based on this genus have contributed

86　　to our current understanding of the evolution of $C_4$ photosynthesis (Gowik and Westhoff,

87　　2011; Powell, 1978; Sage et al., 2013). However, due to a lack of genome reference,

88　　current knowledge of the regulation of photosynthesis genes in this genus is still very

89　　limited. The first version of *Flaveria* genome references, including four species, were

90　　published recently (Taniguchi et al., 2021), and provide valuable resources for protein

91　　coding gene predictions for this genus. However, as these genomes were generated

92　　using short-read whole genome-sequencing, the assembled genomes are fragmented,

93　　which compromises their potential application (Taniguchi et al., 2021). Taking

94　　advantage of long-read genome sequencing technology, here we reported chromosome-

95　　scale genome references of five *Flaveria* species, with which we conducted a

96　　systematic study of CREs and TFs during the evolution of $C_4$ photosynthesis. We found

97　　that ethylene responsive factor (ERF) CREs were recruited by $C_4$ photosynthesis during

98　　evolution. Moreover, intronless ERF TFs that originated 450 mya to cope with

99　　environmental stress were recruited into different $C_4$ lineages; furthermore, our study

100　　highlighted the role of intronless ERF TFs in the evolution and regulation of $C_4$

101　　photosynthesis, and provided a mechanism underlying the repeated evolution of $C_4$

102　　photosynthesis.

103 **Results**

104 **Analysis of five *Flaveria* genome assemblies showed that transposable elements**

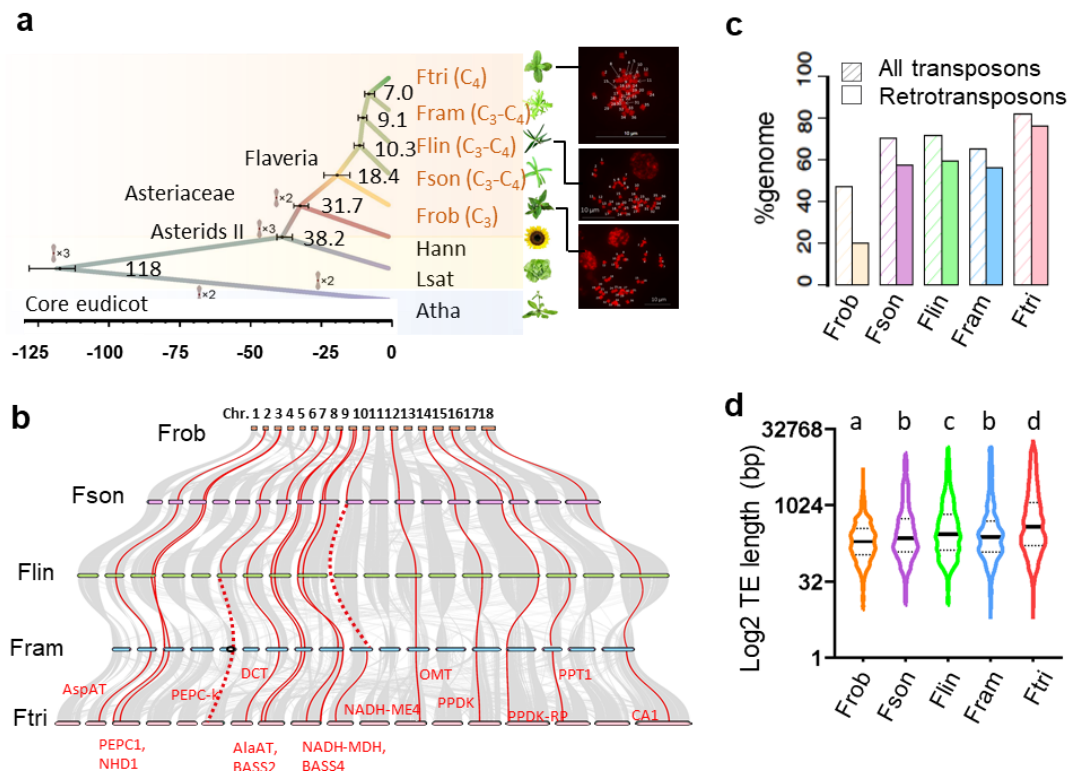105 **were more abundant in the $C_4$ species than other species**

106     The genome sequences of five *Flaveria* species, *i.e., F. robusta* (Frob, $C_3$), *F.*

107 *sonorensis* (Fson, $C_3$-$C_4$), *F. linearis* (Flin, $C_3$-$C_4$), *F. ramosissima* (Fram, $C_3$-$C_4$) and

108 *F. trinervia* (Ftri, $C_4$) were obtained with PacBio RSII single-molecule real-time

109 (SMRT) sequencing technology (Figure 1a). The assembled genome size was

110 gradually increased during the evolution of $C_4$ photosynthesis in this genus, from 0.55

111 Gb in the $C_3$ species Frob, to 1.26~1.66 in the $C_3$-$C_4$ species, and to 1.8 Gb in the $C_4$

112 species Ftri (Table S1), and these data were consistent with the analysis based on flow

113 cytometry (Supplemental Note 1). Based on chromatin conformation capture (Hi-C

114 seq), 98% to 99% of the assembled genome sequences were anchored to 18 pseudo-

115 chromosomes (Figure S1 and Supplemental Note 2), which was supported by

116 fluorescence *in situ* hybridization (FISH) results in Frob, Flin, and Ftri (Figure 1a).

117 This was consistent with the reported chromosome number of 36 (2n) in these five

118 *Flaveria* species (Powell, 1978). Genome completeness was estimated using

119 Benchmarking Universal Single-Copy Orthologues (BUSCO) genes and resulted in

120 coverage from 92.5% to 99.2% of the BUSCO genes. Additionally, an average RNA-

121 seq reads mapping rate of 94.3% (from 86.7 to 97.3%) was obtained (Supplemental

122 Note 3), suggesting completeness and high quality of genome assemblies for the five

123 *Flaveria* species.

124     Although genome size was tripled in the $C_4$ species Ftri compared to the $C_3$ species

125 Frob, the number of protein coding genes was comparable between the $C_3$ and $C_4$

126 species, with 35,875 (Frob) and 32,915 (Ftri) respectively, and 37,028 to 38,652 protein

127 coding genes were predicted in the $C_3$-$C_4$ species (Table S1). We compared the

128 predicted protein coding genes from our assembly with those from Taniguchi's

129 assembly (Taniguchi et al., 2021), we found that around 96% protein coding genes were

130    overlapped between our assembly and Taniguchi's assembly (Taniguchi et al., 2021)

131    (Supplemental Note 4). Therefore, the annotated protein-coding genes in this study can

132    be considered reliable.

133       The chromosome-scale assembly of genome sequences and reliable gene

134    annotations allowed us to study the evolution of known $C_4$ enzymes and $C_4$ transporters

135    (termed as $C_4$ genes) on location on the chromosomes. We identified eight enzymes and

136    seven transporters as $C_4$ versions by combining the gene phylogenetic tree and

137    transcript abundances (Supplemental Note 5). As $C_4$ versions of $C_4$ genes, but not their

138    orthologs, were reported to be induced by light in $C_3$ species, we thus verified the $C_4$

139    version of the $C_4$ enzymes by examining their responsiveness to light. The $C_4$ versions

140    of $C_4$ genes appeared quickly (after 2 hours) and were up-regulated after 4 hours in $C_4$

141    species after being illuminated, and such light responses were intermediate in the $C_3$-

142    $C_4$ species (Figure S2), which suggested the accuracy of identification of the $C_4$ versions

143    of $C_4$ genes, and also revealed a gradual gain of light responsiveness during $C_4$

144    evolution. The synteny of the 18 chromosomes was conserved in the five *Flaveria*

145    species; from 50% to 75% of protein coding genes were colinear between Frob and the

146    other species (Figure 1b and Supplemental Note 2). Notably, the chromosome locations

147    of all 15 $C_4$ version of $C_4$ genes were conserved during evolution (Figure 1b).

148       Transposable elements (TEs) showed the highest abundance in the $C_4$ species,

149    where they accounted for 82% of the total genome, followed by $C_3$-$C_4$ species (from

150    65.6% to 71.8%), whereas that percentage in the $C_3$ species was 47.1% (Figure1c and

151    Supplemental Note 6). In all five species, long terminal repeat retrotransposons (LTR-

152    RTs) comprised the majority of the TEs, accounting for an average 76% of the total TEs

153    (from 42% to 91%) (Figure 1c). $C_4$ genes had longer TEs on the promoter regions in

154    the $C_4$ species than their counterparts in $C_3$ and $C_3$-$C_4$ species do (Figure 1d).

155

**Figure 1. Transposon elements contributed to enlargement of genome size and promoters of C4 genes during *Flaveria* evolution.**

(a) Summary of phylogeny and timescale of the five *Flaveria* species and the three indicated outgroup species. Bars represent 95% confidence intervals of the estimated divergence time. Whole genome duplications are shown at the corresponding node/branch. Panels at the right display fluorescence *in situ* hybridization images to assess the chromosome numbers in Ftri, Flin, and Frob. (b) Collinearity of chromosomes among *Flaveria* species. C4 genes are drawn in red line. Dashed lines represent either failure in anchoring to chromosome (NADP-ME in Flin) or a deletion from the genome (PEPC-k in Fram). (c) Proportions of transposon elements, relative to the whole genome by length. (d) Assessment of 15 C4 genes (from panel c), showing that the C4 species Ftri has relatively longer TEs in the promoter region (3 kb upstream of start codon at the 5' end) of these loci. (Abbreviations: Frob: *F. robusta*, Fson: *F. sonorensis*, Flin: *F. linearis*, Fram: *F. ramosissima*, Ftri: *F. trinervia*.)

170

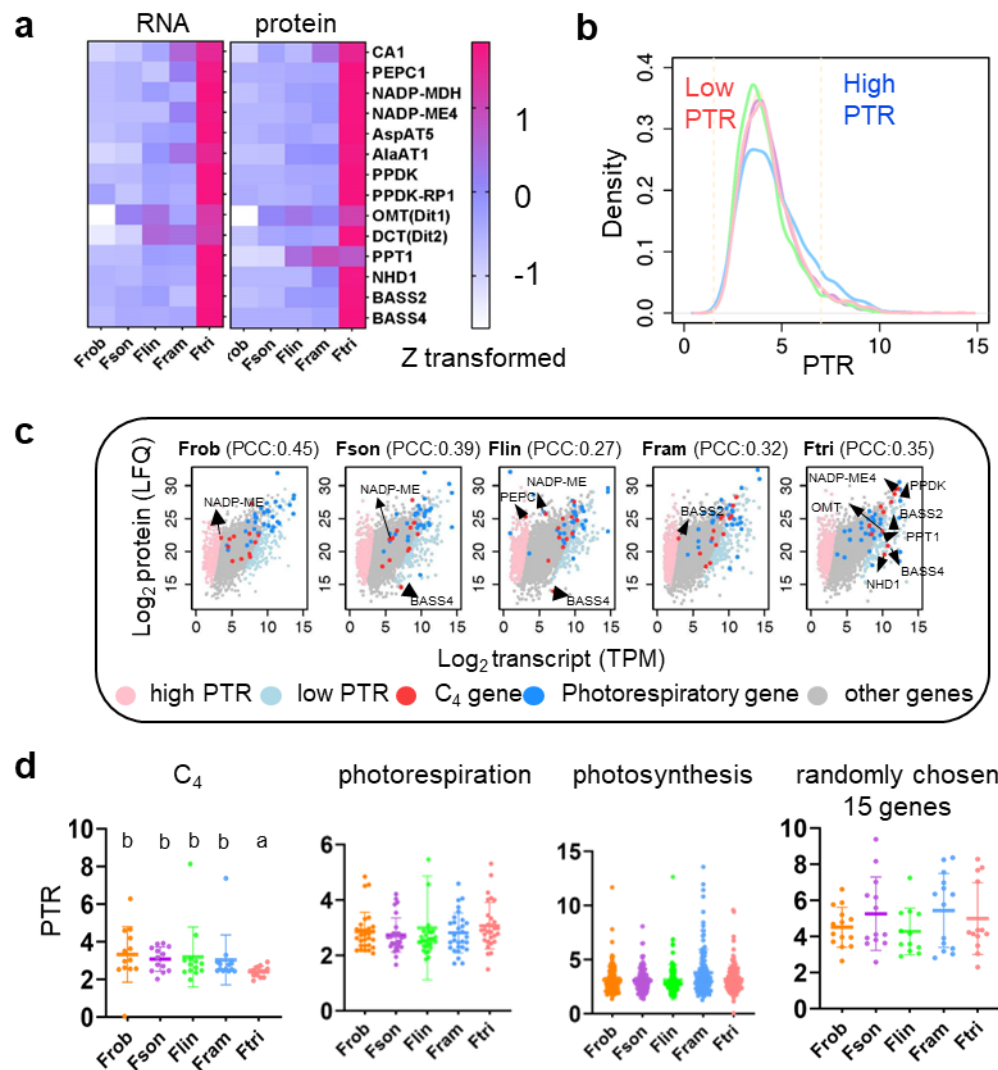**C4 genes acquired elevated protein levels during evolution, which were regulated mainly at transcriptional levels**

The transcript and protein abundances of C4 genes were generally higher in C4 species than in C3 and C3-C4 species (Figure 2a). To study whether transcriptional or

7

175    translational regulation is primarily responsible for the observed difference in protein

176    abundance between species with different photosynthetic types, we compared the

177    protein-to-mRNA ratios (PTR) between genes in five *Flaveria* species. Low PTR

178    genes and high PTR genes were defined as genes with PTR less than the mean PTR

179    minus standard deviation (SD) and higher than the mean PTR plus SD, respectively

180    (Figure 2b), and the remaining genes were defined as moderate PTR genes. In

181    general, there was a positive correlation between mRNA and protein levels, with

182    Pearson correlations ranging from 0.27 to 0.45, and most genes had moderate PTRs

183    (Figure 2c). An average of 166 low PTR genes (from 121 to 201) and 395 high PTR

184    genes (from 375 to 462) were obtained in the five species (Supplemental Notes 7~9).

185    In the $C_4$ species, seven $C_4$ genes were characterized as low PTR genes, whereas three

186    or fewer $C_4$ orthologous genes were low PTR genes in the $C_3$ and $C_3$-$C_4$ species.

187        The low PTR genes were enriched in gene ontology (GO) of photosynthesis and

188    photosynthesis related GO terms, including chloroplast, PSII, and others

189    (Supplemental Note 9), consistent with an early study in *Arabidopsis thaliana* (Atha),

190    which showed that the photosynthesis related genes had significantly lower PTRs than

191    other genes in photosynthetic functional leaf tissues (Mergner et al., 2020)

192    (Supplemental Note 9). $C_4$ genes showed significantly lower PTRs in the $C_4$ species

193    than their orthologs in the $C_3$ and $C_3$-$C_4$ species did (Figure 2d), whereas

194    photorespiratory genes or photosynthesis genes (not including $C_4$ genes) showed

195    comparable PTRs across the five *Flaveria* species. Therefore, during the evolution of

196    $C_4$ photosynthesis, $C_4$ species acquired elevated protein levels for $C_4$ genes, which

197    were regulated mainly at transcriptional levels.

198

**Figure 2. The C$_4$ species had increased transcript abundances of C$_4$ genes.**
(a) RNA-seq and proteomics data for the C$_4$ genes in the five *Flaveria* species show increased transcript and protein abundances of C$_4$ genes in the C$_4$ species Ftri. (b) The protein-to-mRNA ratio (PTR) distribution of genes from the five *Flaveria* species. High PTR and low PTR genes are defined as genes with PTR higher than the mean plus one standard deviation (SD) and with PTR values lower than the mean minus one SD respectively. (c) Scatter plot of protein versus transcript abundance of the five *Flaveria* species. low PTR and high PRT C$_4$ genes were labeled with arrows. Note the trend towards lower PTR for the C$_4$ gene set in Ftri, as compared to the C$_3$ Frob and the three intermediate *Flaveria* species. In contrast, there is no apparent shift in PTR for photorespiratory genes. (d) PTR values for the C$_4$ gene set in the five *Flaveria* species, showing that C$_4$ genes have significantly lower PTR in C$_4$ species Ftri than in the C$_3$ Frob or the three intermediate species. Note that no such decrease is showed for photorespiratory genes, photosynthesis genes and randomly chosen genes. (Abbreviations for proteins see Supplemental Note 8, Figure S14.)
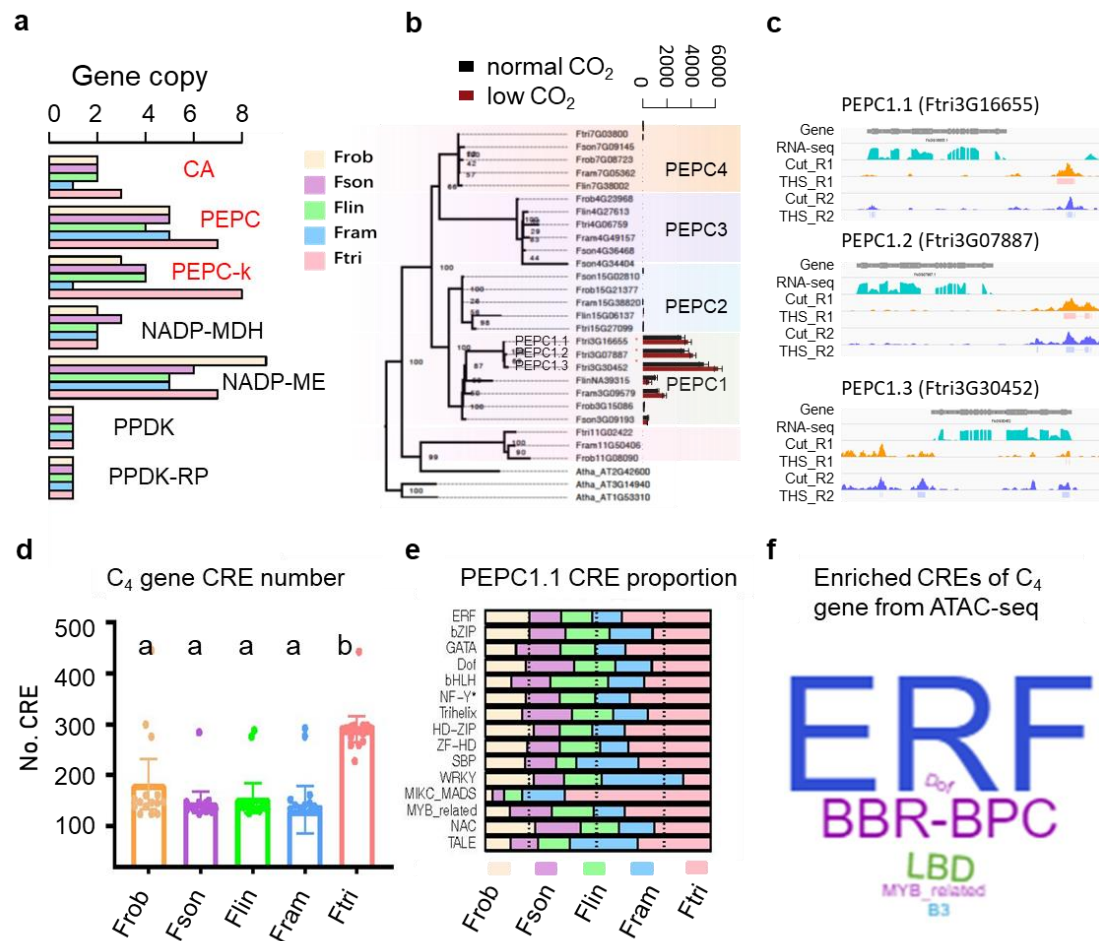
214

**Tandem duplication and recruitment of ERF *cis*-regulatory elements contributed to the increased transcript abundances of C$_4$ genes**

We then analyzed what contributed to the increased transcript abundance of C$_4$ genes in the C$_4$ species. Carbonic anhydrase (CA), phospho*enol*pyruvate carboxylase (PEPC), and PEPC kinase (PEPC-k) showed extra copies in C$_4$ species, which were derived from tandem duplications (Figure 3a and Figure S3). For example, the C$_4$ version of PEPC, termed PEPC1 because it showed the highest transcript abundance among the other paralogs in C$_4$ species, had three copies in the C$_4$ species Ftri, and only one copy in the other species. The three paralogs of PEPC1 in Ftri, termed as PEPC1.1, PEPC1.2, and PEPC1.3, were located on the same chromosome (Chr3) (Figure 3b). The existence of the three PEPC1 paralogs on the chromosome was further verified by PCR (Supplemental Note 10). In Ftri, all three PEPC1s had comparable transcript abundances, which were higher than those in the other four species (Figure 3b). Additionally, they were all upregulated under long-term low CO$_2$ treatment (100 ppm) compared to normal CO$_2$ conditions (380 ppm), suggesting that these triplets hosted shared regulatory mechanisms. Indeed, all three paralogs harbored the mesophyll expression module 1 (MEM1) CRE (Akyildiz et al., 2007) ( Supplemental Note 10); moreover, the three paralogs showed similar signatures of chromatin accessibility through transposase-accessible chromatin using sequencing (ATAC-seq) (Figure 3c and Supplemental Note 11).

We further characterized the distribution of CREs on the promoter regions of C$_4$ genes. Ftri (C$_4$) showed significantly more CREs for C$_4$ genes than the other species did (p<0.001, *t*-test) (Figure 3d). Notably, when the total numbers of CREs of each gene in the five species were compared, ERF CREs were the most abundant of all examined C$_4$ genes (Figure 3e and Supplemental Note 12). For example, for PEPC1 (PEPC1.1 from Ftri), there were 89 predicted ERF CREs from 5 species, followed by bZIP (54 CREs) and GATA (51 CREs).

242    To examine whether the ERF CREs were localized in the accessible chromatin

243    regions (ACRs) in the $C_4$ species, we analyzed the enriched CREs in the ACRs (ACR-

244    CREs) obtained from two biological ATAC-seq (Supplemental Note 11). We

245    categorized genes associated with ACRs-CREs into three types according to their

246    distance to the nearest gene, *i.e.*, genic (gACR-CREs; overlapping a gene), upstream

247    (upACR-CREs; within 3 kb upstream of the start codon of a gene) or downstream

248    (downACRs-CRES; within 3 kb downstream of the stop codon of a gene). We then

249    calculated enriched CREs in ACR-CREs. Across all three types of ACR-CREs, ERF

250    CREs had the highest abundance among enriched CREs. Moreover, ERF dominated

251    the enriched ACR-CREs of $C_4$ genes, as well as in photosynthetic and

252    photorespiratory genes (Figure 3f and Figure S4).

253    Taken together, our results suggested possible roles of both tandem duplications

254    and recruitment of ERF CREs in the elevation of transcript abundances of $C_4$ genes in

255    the $C_4$ species Ftri.

**Figure 3. Tandem duplications and recruitments of ERF *cis*-regulatory elements contributed to the increased transcript abundances of C4 genes in the C4 species Ftri**

(a) Copy number of the $C_4$ version of $C_4$ enzymes in five *Flaveria* species. Note that CA, PEPC, and PEPC-K have more copies in the $C_4$ species Ftri than other *Flaveria* species. (b) Gene tree of PEPC orthologs, PEPCs from *Arabidopsis thaliana* (Atha) are used as outgroups. PEPCs in *Flaveria* species are categorized into four groups, and PEPC1 is the $C_4$ version according to the highest expression levels among all PEPCs. PEPC1 has three copies in the $C_4$ species Ftri, showing comparable transcript abundances in leaves and uniform upregulation when plants were grown under low $CO_2$ conditions (100 ppm) compared to normal $CO_2$ conditions (380 ppm). (c) Integrated Genome Viewer (IGV) of RNA-seq reads and ATAC-seq reads of three PEPC1 in Ftri. Tn5 cuts and transposase hypersensitive sites (THS) from two biological replicates show that the three PEPC1 have shared chromatin accessibility upstream of their coding region. (d) Bar plots show the number of predicted *cis*-regulatory elements (CREs) from the promoter region (3kb upstream of start codon) of all $C_4$ genes in the five *Flaveria* species. (e) An example of the distribution of the top 15 CREs in $C_4$ genes, noting that Ftri has more CREs in PEPC1.1 than other species. CREs of the 3kb of the 5'-flank regions of $C_4$ genes were predicted applying the online tool Plantpan3.0 (score>=0.99).

12

276  The TF families from top to bottom are ordered in a decreasing rank of total number of
277  CREs from the five *Flaveria* species. (f) Word cloud represents the enriched CREs
278  associated with $C_4$ genes in the $C_4$ species Ftri based on ATAC-seq, including those
279  within 3kb upstream of start codon, within 3kb downstream of the stop codon and
280  within the gene body.

281

**282  Intronless ERF transcriptional factors were recruited in parallel in different $C_4$**

**283  species**

284      Given that many ERF CREs were recruited by $C_4$ genes in Ftri ($C_4$), we tested

285  whether cognate ERF TFs were recruited in the same manner by $C_4$ photosynthesis. We

286  constructed a genome-wide co-regulatory network (GRN) of the five *Flaveria* species

287  based on the gene expression profiles of at least 18 RNA-seq datasets either from a

288  previous work (Zhu, 2020) or generated in the current study (Supplemental Note 13).

289  We then obtained the sub GRN comprising $C_4$ genes and their co-regulated TFs

290  ($C_4$GRN). TFs that had no predicted cognate CREs within 3 kb upstream of the start

291  codon were filtered out. ERF, bHLH, MYB, NAC, and C2H2 were the top five most

292  abundant TF families in the $C_4$GRN of the five *Flaveria* species (Figure 4a and Figure

293  S5). In the $C_4$ species Ftri, 324 TFs were predicted to be co-regulated with $C_4$ genes

294  (Figure 4b), among which bHLH was the most prevalent TFs, with 29 genes, followed

295  by the MYB related and ERF TF families, with 27 and 26 genes respectively (Figure

296  4c). Notably, ERF TFs were much more abundant in the $C_4$GRN of the $C_4$ species than

297  in other species, though the number of predicted ERF TFs were comparable in all five

298  *Flaveria* species (Figure 4a and Supplemental Note 13), suggesting that ERF TFs were

299  preferentially recruited by $C_4$ genes during evolution.

300      $C_4$ photosynthesis has appeared in more than 65 evolutionary independent lineages

301  (Sage et al., 2012), and ERF CREs were previously found abundant in other $C_4$ lineages,

302  including *Zea mays* (corn; herein Zmay), *Setaria italica* (foxtail millet), and *Sorghum*

303  *bicolor* (sorghum) (Supplemental Note 12) (Burgess et al., 2019; Marand et al., 2021).

304  We investigated whether ERF TFs were also convergently recruited in other $C_4$ lineages.

13

305    We used Zmay, a model species for $C_4$ research, for this test. Specifically, we analyzed

306    a recently published leaf GRN for this species, which was constructed based on a

307    combination of Chip-seq data, gene co-expression data, and a machine-learning based

308    co-localization model(Tu et al., 2020). This genomic scale GRN included 1,475 TFs

309    from 54 TF families, in which bHLH was the most prevalent family, with 138 genes,

310    followed by ERF and MYB, with 136 and 108 genes, respectively (Tu et al., 2020)

311    (Figure S6a). The $C_4$GRN included 108 TFs from 15 TF families (Figure S6b), in which

312    ERF TFs was the most prevalent ones, with 20 genes, followed by the WRKY and

313    bHLH families, with 12 and 11 genes, respectively (Figure 4d). Therefore, ERF TFs

314    were convergently recruited in both Ftri and Zmay, whose last common ancestor

315    diverged around 160 mya (Kumar et al., 2017).

316        We further identified the shared TFs between Ftri and Zmay, *i.e.*, those recruited

317    by both species and found in the same orthologous group based on Orthofinder's

318    analysis (Methods). Shared TFs were not required to regulate the same $C_4$ genes in the

319    two species. Our analysis found shared TFs from 27 orthologous groups which included

320    63 TFs from Ftri and 47 TFs from Zmay respectively. Again, the ERF TFs were the

321    most abundant families in both species, including 14 (22.2%) and 12 (25.5%) of shared

322    TFs in Ftri and Zmay, respectively (Figure 4e), and the targeted genes of these shared

323    TFs covered 14 of the 15 $C_4$ genes (Figure 4f). Notably, among the shared ERF TFs, 12

324    out of the 14 in Ftri and 11 out of 12 in Zmay were intronless genes, which account for

325    66.7% and 73.3% of total shared intronless TFs in Ftri (18 intronless TFs) and Zmay

326    (15 intronless TFs), respectively (Figure 4f). Compared to other intronless ERF TFs in

327    Zmay, the shared intronless ERFs showed more MC preferential expression (Figure 4g).

328        We further investigated the portion of intronless genes in each of the TF families

329    in Ftri and Zmay. Intronless genes showed the most occurrences in ERF families, in

330    which, 62% and 80% of ERF TFs were intronless in Ftri and Zmay respectively,

331    accounting for 35.2% and 26.3% of the total intronless TFs in these species (Figure 4h).

332    ERF TFs were also the most abundant intronless TF family in other land plant species,

333    regardless of whether they are monocotyledonous or dicotyledonous, $C_3$ or $C_4$ species

334    (Figure S7). These intronless ERFs showed greater changes in transcript abundances in

335    response to low $CO_2$ stress compared to intron-containing ERFs in the $C_4$ species Ftri

336    but not in other non-$C_4$ *Flaveria* species. Similarly, intronless ERF TFs exhibited rapid

337    and increased changes in gene expression in response to light induction in the $C_4$ species

338    Zmay but not in the $C_3$ species *Oryza sativa, i.e.*, rice ($P<0.05$, Wilcoxon.test, Figure

339    S8). These properties of intronless ERF TFs, including MC preferential expression

340    (Figure 4g) and greater responses to low $CO_2$ and light, might have contributed to their

341    role in $C_4$ photosynthesis.

342        We further analyzed the properties of one shared intronless ERF TF, *i.e.*, EREB34

343    (Figure S9a) between the Ftri $C_4$GRN and Zmay $C_4$GRN. EREB34 showed conserved

344    expression profiles along the leaf developmental gradient or leaf age between $C_3$ and

345    $C_4$ species. However, EREB34 showed significantly higher transcript abundance in the

346    $C_4$ species than in $C_3$ species, which was shown both when we compared the

347    evolutionarily close $C_3$ and $C_4$ species pairs individually (Figure S9b ), and when we

348    compared 30 $C_4$ and 17 $C_3$ species representing 18 independent lineages of $C_4$ evolution

349    (Steven Kelly, 2018) (Figure S9c). EREB34 also showed MC preferential expression

350    in $C_4$ species (Figure S9d). All these results suggested that EREB34 may play a role in

351    the evolution of $C_4$ photosynthesis.

15

**Figure 4. Intronless ERF were recruited by C₄ genes in different C₄ lineages**

(a) The Top five most abundant TF families of all annotated TFs (top panel) and TFs that are co-regulated with C₄ genes (bottom panel) (a) The network of C₄ genes and TFs, which is termed as C₄GRN in Ftri. (b) Families of TFs from C₄GRN of Ftri. (c) Families of TFs from C₄GRN of Zmay. Zmay C₄GRN is extracted from published gene regulatory network in Tu et. al, 2019. (d) Orthologous TFs from C₄GRN of Ftri and Zmay and their distribution in TF families. Orthologous groups were predicted with Orthofinder, and orthologous TFs between Ftri and Zmay, termed as shared TFs, are those from the same orthologous groups. (e) Regulatory network of shared TFs and C₄ genes in Ftri and Zmay. Note that among the shared ERF TFs, 12 of 14 in Ftri and 11 of 12 in Zmay are intronless genes. (f) Bundle sheath cell (BSC) and mesophyll cell (MC) preferential expression of all intronless ERF TFs and shared intronless ERF TFs in Zmay. Y-axis shows log2 ratio of transcript abundance of each gene in BSC to that in MC. (g) Number of intronless genes in each TF family. (Abbreviations: GRN: gene co-regulatory network, MC: mesophyll cell, BSC: bundle sheath cell, Ftri: *Flaveria trinervia*, Zmay: *Zea mays*.)

16

**Intronless ERF TFs recruited by C$_4$ photosynthesis originated in the Late Ordovician around 450 million years ago**

ERF TFs belong to the AP2/ERF superfamily, which is one of the largest families of plant-specific TFs, and play vital roles in responses to various biotic and abiotic stresses (Feng et al., 2020; Gu et al., 2017; Xie et al., 2019). Our data showed that intronless ERF TFs were recruited as major regulators of C$_4$ photosynthesis in both monocots and dicots. Considering that monocots and dicots diverged ~160 mya (Kumar et al., 2017), while C$_4$ photosynthesis emerged ~ 35 mya (Sage et al., 2011), elements shared between the monocotyledonous and dicotyledonous C$_4$ species were likely recruited before the divergence of monocots and dicots. We hence examined the origin of the intronless ERF TFs in plants. Specifically, we first surveyed the distribution of intronless genes in all annotated TF families based on the plantTFDB online tool (Jin et al., 2017) in 23 species spanning a wide spectrum of Viridiplantae (green plants), including four species from Chlorophyta, *Marchantia polymorpha* (Mploy, liverwort), which is regarded as one of the earliest land species (Delaux et al., 2019), seven monocotyledonous species, and 11 dicotyledonous species including the five *Flaveria* species sequenced here (Figure 5a). We included five and two C$_4$ species from monocots and dicots, respectively.
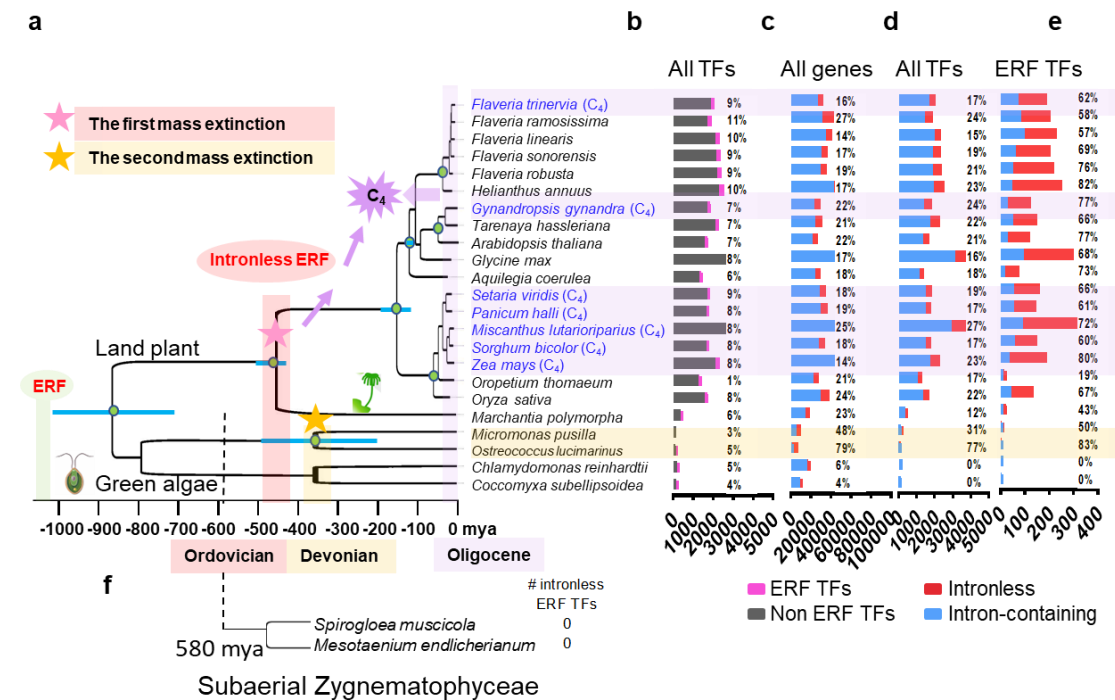
We found that ERF TFs were present in all Viridiplantae, accounting for 7% of total TFs on average (Figure 5b). Furthermore, intronless genes were also present in all Viridiplantae, accounting for 22% of total annotated genes on average (Figure 5c). Intronless ERF TFs were found in all land plants and in the clade of Chlorophyta that includes *Micromonas pusilla* (Mpus), termed Mpus clade hereafter, but not the clade that includes *Chlamydomonas reinhardtii* (Crei), termed Crei clade hereafter (Figure 5d and Figure 5e). To determine whether the intronless ERF were specifically absent from the two species in the Crei clade or from the whole clade, we examined the other two species from the Crei clade with genome sequences available in the Phytozome

17

396    database (https://phytozome-next.jgi.doe.gov), *i.e.*, *Dunaliella salina* and *Volvox*

397    *carteri*. We found that intronless ERF TFs were not present in those two species either

398    (Supplemental Note 14), implying that intronless ERF were absent from the Crei clade.

399        We then asked whether intronless ERF TFs were lost in Crei clade specifically or

400    if there were two independent gains of intronless ERF TFs in the plant kingdom. We

401    studied two species from Zygnematophyceae, which is the closest extant sister branch

402    of land plants and evolved 580 mya (Gitzendanner et al., 2018), *i.e.*, *Spirogloea*

403    *muscicola* and *Mesotaenium endlicherianum* (Cheng et al., 2019). Intronless ERF TFs

404    were also absent in these two species (Figure 5f), suggesting that intronless ERF TFs

405    in land plants and aquatic algae (Mpus clade) emerged from two independent

406    evolutionary events. As further evidence, intronless ERF TFs from the Mpus clade

407    showed nearly no orthologs in land plant species (Supplemental Note 14). Therefore,

408    there were two independent gains of intronless ERF TFs during the evolution of

409    Viridiplantae, and those evolved from the common ancestor of land plant species were

410    recruited in $C_4$ photosynthesis.

411        What might have promoted the emergence of intronless ERF TFs recruited to

412    support $C_4$ photosynthesis during the evolution of plants? To study this, we examined

413    the recorded extreme climatic events around the period of the two independent

414    occurrences of intronless ERF TFs in Viridiplantae. The first occurrence of intronless

415    ERF TFs is around 450 mya when the land plants diverged from aquatic algae

416    (Sanderson et al., 2004). This period coincided with the time of the Earth's first mass

417    extinction, around 447~444 mya during the Late Ordovician (Finnegan et al., 2012;

418    Sheehan, 2001). The second appearance of intronless ERF TFs, observed in the Mpus

419    clade, occurred around 380 mya, which coincided with the time of the second mass

420    extinction, around 372 mya during the Late Devonian (Da Silva et al., 2020; De

421    Vleeschouwer et al., 2017) (Figure 5). Dramatic climate changes such as low

422    temperature and oxygen deprivation have been proposed to underlie these two mass

423    extinctions [17,19]. Therefore, intronless ERF TFs might be the products of ancestral plants

18

424    coping with extreme climate events. Long after the first emergence of these intronless

425    ERF TFs in land species, around 35 mya (Sage et al., 2011), some of those TFs,

426    especially those with strong cell specific expression patterns (Figure 4g) were recruited

427    in C$_4$ photosynthesis.

428



429

**Figure 5. Intronless ERF TFs that were recruited in C$_4$ photosynthesis emerged much earlier than C$_4$ photosynthesis did**

(a) Phylogenetic relationships of 23 species. C$_4$ species are labeled in blue. The divergence time of each node is referenced from Timetree (http://timetree.org/). Two independent evolutionary origins of intronless ERF TFs are proposed, the occurrences of which coincide with the first mass extinction during the Late Ordovician (~482 mya, pink bar) and the second mass extinction during the Late Devonian (~358 mya, yellow bar), respectively. A star represents an independent evolutionary origin of intronless ERF TFs (b) Proportions of ERF TFs to total TFs. (c) Proportions of intronless gene to total protein coding genes. (d) Proportions of intronless TFs to total TFs. (e) Proportions of intronless ERF TFs to total ERF TFs. (f) The number of intronless ERF TFs in the two subaerial species from Zygnematophyceae, which is the sister branch of land plants and split from land plants ~580 mya. (Abbreviation: ERF: ethylene responsive factor, mya: million years ago.)

444

**Discussion**

Identifying key regulators of $C_4$ photosynthesis is a major task required for $C_4$ engineering (Cui, 2021; Hibberd and Covshoff, 2010; Schluter and Weber, 2020; Westhoff and Gowik, 2010). The high-quality genome sequences of five *Flaveria* species offered a rich resource to support evolutionary and regulatory study of $C_4$ photosynthesis. With this resource, we showed that intronless ethylene responsive factors (ERF) transcription factors (TF), a class of TFs involving in stress responses in plants (Christin et al., 2008; Ehleringer et al., 1991; Sage et al., 2011; Sage et al., 2012), played a role during the evolution of $C_4$ photosynthesis. These intronless ERF TFs, originating from ~450 mya (Sanderson et al., 2004) (Figure 5), have been repetitively recruited by different $C_4$ lineages during evolution. Therefore, our results provided a molecular mechanism underlying shared TFs and *cis*-regulatory elements (CREs) between monocotyledonous and dicotyledonous $C_4$ species that diverged 160 mya (Kumar et al., 2017), though the first $C_4$ plants emerged ~35 mya (Sage et al., 2011). The parallel recruitment of intronless ERF TFs implied that they may be used as targets during the current efforts in $C_4$ engineering.

Why intronless ERF TFs? Intronless genes, featuring short mRNA length and lower transcript abundances compared to intron-containing genes (Shabalina et al., 2010) (Supplemental Note 15), play roles in plant responses to drought and salt stress (Liu et al., 2021). Intronless genes, regardless of being TF or not, showed more changes than intron-containing genes to low $CO_2$ stress in all five *Flaveria* species, and greater and faster changes to light induction in both $C_3$ and $C_4$ species at the transcriptional level (Supplemental Note 15). Ethylene, an ancient plant hormone (Ju et al., 2015), bridges plant developmental adaptation and a changing environment (Merchante et al., 2013). ERF TFs, the last step of the ethylene signaling pathway, regulate the response of plants to environmental changes (Xie et al., 2019). Recently, one intronless ERF TF was reported to simultaneously modulate photosynthesis and nitrogen utilization in rice (Wei et al., 2022). ERF TFs showed remarkable changes to low $CO_2$ stress in all five

20

473  *Flaveria* species (Supplemental Note 15). Being evolutionary old and functioning in

474  responding to environmental changes may underlie the observation that around 70% of

475  ERF TFs evolved to be intronless genes in land plant species (Figure S7). In addition,

476  ERF TFs existed widely across the plant kingdom, with a large presence of cognate

477  CREs in plant genomes (Supplemental Note 12). The abundance of intronless ERF TFs

478  and cognate CREs provided molecular resources for the evolution of $C_4$ photosynthesis

479  in coping with environmental stressors, such as low $CO_2$, drought, and high light and

480  high temperature conditions (Christin et al., 2008; Ehleringer et al., 1991; Sage et al.,

481  2011; Sage et al., 2012) .

482      Intronless genes are also present in animals and fungi (see database:

483  http://v2.sinex.cl/) (Jorquera et al., 2021). The ancient origin of intronless genes has

484  been reported from animals. For example, intronless type I interferon (INF) in animals

485  evolved from intron-containing type I INF in fish and amphibians around 350 mya

486  during the Devonian (Gan et al., 2017), coinciding with the time when intronless ERF

487  TFs originated in the Mpus algae clade. This suggested that environmental

488  perturbations during the Devonian triggered the birth of new classes of intronless

489  genes in both animals and plants. Interestingly, in humans, the counterpart of plant

490  ERF TFs is the G-protein-coupled receptors (GPCRs). Around 50% of GPCRs are

491  intronless genes, accounting for 53% of total human intronless genes (Gentles and

492  Karlin, 1999; Grzybowska, 2012), compared to 70% ERF TFs that are intronless

493  genes, accounting for around 30% of total plant intronless genes (Figure 5). Notably,

494  ERFs in plants and GPCRs in humans have analogous functions in receiving and

495  transducing signals from the external environment (Grzybowska, 2012; Xie et al.,

496  2019). Therefore, particular types of intronless genes were retained for evolutionary

497  adaptions in both the plant and animal kingdoms (Grzybowska, 2012; Xie et al.,

498  2019).

499

**Acknowledgements**

**Authors' contributions**

XGZ, TL, CL and MJAL designed the study and wrote the paper. HD and ZG performed genome assembly and annotation, MJAL performed genome comparison analysis, qRT-PCR and RNA-seq analysis, HY conducted proteomics analysis, GC wrote the paper, FC performed gene regulatory network construction, YYZ performed PCR verification of the three paralogs of PEPC1s in Ftri, QT performed Ka/Ks analysis, FM and YW performed plasmodesmata analysis, CX performed transcriptional factor prediction, YZ and HL performed genome annotation of Fram, YT constructed *Flaveria* workspace in China National GeneBank (CNGB), LF and QG performed genome assembly of Fram, YQ perform transposon analysis, QZ and JZ performed syntenic analysis.

**Competing interests**

The authors declare no conflict of interests.

**Figure legends**

**Figure 1. Transposon elements contributed to enlargement of genome size and promoters of C$_4$ genes during *Flaveria* evolution.**

(a) Summary of phylogeny and timescale of the five *Flaveria* species and the three

22

525     indicated outgroup species. Bars represent 95% confidence intervals of the estimated

526     divergence time. Whole genome duplications are shown at the corresponding

527     node/branch. Panels at the right display fluorescence *in situ* hybridization images to

528     assess the chromosome numbers in Ftri, Flin, and Frob. (c) Collinearity of

529     chromosomes among *Flaveria* species. $C_4$ genes are drawn in red line. Dashed lines

530     represent either failure in anchoring to chromosome (NADP-ME in Flin) or a deletion

531     from the genome (PEPC-k in Fram). (b) Proportions of transposon elements, relative to

532     the whole genome by length. (d) Assessment of 15 $C_4$ genes (from panel c), showing

533     that the $C_4$ species Ftri has relatively longer TEs in the promoter region (3 kb upstream

534     of start codon at the 5' end) of these loci. (Abbreviations: Frob: *F. robusta*, Fson: *F.*

535     *sonorensis*, Flin: *F. linearis*, Fram: *F. ramosissima*, Ftri: *F. trinervia*.)

536     **Figure 2. The $C_4$ species had increased transcript abundances of $C_4$ genes.**

537     (a) RNA-seq and proteomics data for the $C_4$ genes in the five *Flaveria* species show

538     increased transcript and protein abundances of $C_4$ genes in the $C_4$ species Ftri. (b) The

539     protein-to-mRNA ratio (PTR) distribution of genes from the five *Flaveria* species. High

540     PTR and low PTR genes are defined as genes with PTR higher than the mean plus one

541     standard deviation (SD) and with PTR values lower than the mean minus one SD

542     respectively. (c) Scatter plot of protein versus transcript abundance of the five *Flaveria*

543     species. low PTR and high PRT $C_4$ genes were labeled with arrows. Note the trend

544     towards lower PTR for the $C_4$ gene set in Ftri, as compared to the $C_3$ Frob and the three

545     intermediate *Flaveria* species. In contrast, there is no apparent shift in PTR for

546     photorespiratory genes. (d) PTR values for the $C_4$ gene set in the five *Flaveria* species,

547     showing that $C_4$ genes have significantly lower PTR in $C_4$ species Ftri than in the $C_3$

548     Frob or the three intermediate species. Note that no such decrease is showed for

549     photorespiratory genes, photosynthesis genes and randomly chosen genes.

550     (Abbreviations for proteins see Supplemental Note 8, Figure S14.)

551     **Figure 3. Tandem duplications and recruitments of ERF *cis*-regulatory elements**

552     **contributed to the increased transcript abundances of $C_4$ genes in the $C_4$ species**

553 **Ftri**

554 (a) Copy number of the $C_4$ version of $C_4$ enzymes in five *Flaveria* species. Note that

555 CA, PEPC, and PEPC-K have more copies in the $C_4$ species Ftri than other *Flaveria*

556 species. (b) Gene tree of PEPC orthologs, PEPCs from *Arabidopsis thaliana* (Atha) are

557 used as outgroups. PEPCs in *Flaveria* species are categorized into four groups, and

558 PEPC1 is the $C_4$ version according to the highest expression levels among all PEPCs.

559 PEPC1 has three copies in the $C_4$ species Ftri, showing comparable transcript

560 abundances in leaves and uniform upregulation when plants were grown under low $CO_2$

561 conditions (100 ppm) compared to normal $CO_2$ conditions (380 ppm). (c) Integrated

562 Genome Viewer (IGV) of RNA-seq reads and ATAC-seq reads of three PEPC1 in Ftri.

563 Tn5 cuts and transposase hypersensitive sites (THS) from two biological replicates

564 show that the three PEPC1 have shared chromatin accessibility upstream of their coding

565 region. (d) Bar plots show the number of predicted *cis*-regulatory elements (CREs) from

566 the promoter region (3kb upstream of start codon) of all $C_4$ genes in the five *Flaveria*

567 species. (e) An example of the distribution of the top 15 CREs in $C_4$ genes, noting that

568 Ftri has more CREs in PEPC1.1 than other species. CREs of the 3kb of the 5'-flank

569 regions of $C_4$ genes were predicted applying the online tool Plantpan3.0 (score>=0.99).

570 The TF families from top to bottom are ordered in a decreasing rank of total number of

571 CREs from the five *Flaveria* species. (f) Word cloud represents the enriched CREs

572 associated with $C_4$ genes in the $C_4$ species Ftri based on ATAC-seq, including those

573 within 3kb upstream of start codon, within 3kb downstream of the stop codon and

574 within the gene body.

575 **Figure 4. Intronless ERF were recruited by $C_4$ genes in different $C_4$ lineages**

576 (a) The Top five most abundant TF families of all annotated TFs (top panel) and TFs

577 that are co-regulated with $C_4$ genes (bottom panel) (a) The network of $C_4$ genes and TFs,

578 which is termed as $C_4$GRN in Ftri. (b) Families of TFs from $C_4$GRN of Ftri. (c) Families

579 of TFs from $C_4$GRN of Zmay. Zmay $C_4$GRN is extracted from published gene

580 regulatory network in Tu et. al, 2019. (d) Orthologous TFs from $C_4$GRN of Ftri and

581 Zmay and their distribution in TF families. Orthologous groups were predicted with

582 Orthofinder, and orthologous TFs between Ftri and Zmay, termed as shared TFs, are

583 those from the same orthologous groups. (e) Regulatory network of shared TFs and $C_4$

584 genes in Ftri and Zmay. Note that among the shared ERF TFs, 12 of 14 in Ftri and 11

585 of 12 in Zmay are intronless genes. (f) Bundle sheath cell (BSC) and mesophyll cell

586 (MC) preferential expression of all intronless ERF TFs and shared intronless ERF TFs

587 in Zmay. Y-axis shows log2 ratio of transcript abundance of each gene in BSC to that

588 in MC. (g) Number of intronless genes in each TF family. (Abbreviations: GRN: gene

589 co-regulatory network, MC: mesophyll cell, BSC: bundle sheath cell, Ftri: *Flaveria*

590 *trinervia*, Zmay: *Zea mays*.)

591 **Figure 5. Intronless ERF TFs that were recruited in $C_4$ photosynthesis emerged**

592 **much earlier than $C_4$ photosynthesis did**

593 (a) Phylogenetic relationships of 23 species. $C_4$ species are labeled in blue. The

594 divergence time of each node is referenced from Timetree (http://timetree.org/). Two

595 independent evolutionary origins of intronless ERF TFs are proposed, the occurrences

596 of which coincide with the first mass extinction during the Late Ordovician (~482 mya,

597 pink bar) and the second mass extinction during the Late Devonian (~358 mya, yellow

598 bar), respectively. A star represents an independent evolutionary origin of intronless

599 ERF TFs (b) Proportions of ERF TFs to total TFs. (c) Proportions of intronless gene to

600 total protein coding genes. (d) Proportions of intronless TFs to total TFs. (e) Proportions

601 of intronless ERF TFs to total ERF TFs. (f) The number of intronless ERF TFs in the

602 two subaerial species from Zygnematophyceae, which is the sister branch of land plants

603 and split from land plants ~580 mya. (Abbreviation: ERF: ethylene responsive factor,

604 mya: million years ago.)

605

## Methods

### Plant materials and fluorescence in situ hybridization assay

*F. robusta* (Frob, C$_3$) and *F. ramosissima* (Fram, C$_3$-C$_4$) were provided by Prof. Peter Westhoff (Heinrich Heine University, Germany). Seeds of *F. sonorensis* (Fson, C$_3$-C$_4$), *F. linearis* (Flin, C$_3$-C$_4$) and *F. trinervia* (Ftri, C$_4$) were obtained from Prof. Rowan F. Sage (University of Toronto, Canada). Plants were grown in soil in green house as depicted in (Lyu et al., 2020).

The chromosome numbers of Frob, Flin and Ftri were investigated applying fluorescence in situ hybridization assay (FISH). Mitotic metaphase spreads of meristem root tip cells were prepared following (Deng et al., 2012). FISH was performed following (Li et al., 2019) with slight modifications, which is briefly depicted in Supplemental Note 2.

### Genome sequencing

Total DNA was extracted from young leaves. PacBio sequencing libraries were constructed following the tips of Pacific Biosciences (USA). DNA fragments of 0.5-18kb were chosen using BluePippin electrophoresis (Sage Science, USA). Libraries were then sequenced on the PacBio Sequel platform (PacBio, USA). The N50 of PacBio reads were from 16.4 to 21.9 kbp. Around 120 GB data were produced for each species on average. Genome coverage is from 66.9-fold (Ftri) to 232.2-fold (Frob). Besides, short reads were sequenced in Illumima X Ten platform in paired-end 150 bp mode. Around 200 million short reads were obtained for each species, which were used for genome assembly polishing as well as genome assembly completeness estimation. Hi-C libraries were constructed following (Mascher et al., 2017). Two Hi-C libraries were constructed for each species, with an inserted size of ~350 bp, libraries were sequenced in Illumima X Ten platform. From 291Gb to 325Gb 150-bp paired-ended cleans data were generated for each species.

26

**De novo assembly**

*Flaveria* nuclear genome sequences were assembled into 18 pseudochromosomes in a step-wise way. Sequencing adaptors were removed, and reads with low quality and short length were filtered applying PacBio SMRT Analysis package with following parameters: readScore, 0.75; minSubReadLength 50. The remained high-quality PacBio subreads were then corrected and contigs were assembled using Canu (v1.8) (Koren et al., 2017) with following parameters: useGrid = true, minThreads=4, genomeSize=1200m, minOverlapLength = 500, minReadLength = 1000. For contig polishing, the Illumina paired-end reads were mapped to assembled contigs applying bwa mem (bwa v0.7.17) (Li and Durbin, 2009), low qualified mapped reads were filtered off applying samtools (v1.11) (Li et al., 2009) with q30 setting. Pilon (v1.22) (Walker et al., 2014) were applied to polish with the following parameters: --mindepth 10 --changes --fix bases.

For Fram specifically, the BioNano next-generation mapping system was used to help high-quality genome assembly. DNA was labelled at Nt.BspQI sites applying the IrysPrep kit (BioNano Genomics, USA). Molecules collected from BioNano chips (BioNano Genomics, USA) were de novo assembled applying RefAligne and Assembler offered on the BioNano (Pendleton et al., 2015) using following parameters: -U -d -T 20 -j 4 -N 10 -i 5, which resulted in the optical genome maps. Next, genome assembly resulting from Pilon (v1.22) (Walker et al., 2014) mentioned above were then evaluated and corrected by aligning with the optical genome maps. Corrected contigs and optical genome maps were aligned and merged applying hybridScaffold.pl (Pendleton et al., 2015) which resulted in hybrid scaffolds. Next, HERA(Du and Liang, 2019) was used to fill gaps of obtained hybrid scaffold in following parameters: InterIncluded_Side=30000, InterIncluded_Identity=99, InterIncluded_Coverage=99, MinIdentity=97, MinCoverage=90, MinLength=5000, MinIdentity_Overlap=97, MinOverlap_Overlap=1000, MaxOverhang_Overlap=100, MinExtend_Overlap=500. Obtained hybrid scaffolds were then used for following

27

660    assembly.

661    Followed, assembled genome sequences were improved using Hi-C data in two

662    steps. First, contigs were corrected using Hi-C data. Briefly, low-quality Hi-C data

663    (over 10% N base pairs or Q10 < 50%) were removed, and remained reads were

664    mapped to assembled contigs applying bwa (v0.7.17) (Li and Durbin, 2009) with 'aln'

665    settings and other parameters were in default. Only uniquely mapped reads were used

666    to perform re-assembly. Invalid mapping was filtered off applying HiC-Pro (v2.11.1)

667    (Servant et al., 2015) with following settings: mapped_2hic_fragments.py -v -S -s 100

668    -l 1000 -a -f -r -o. Next, corrected contigs were re-assembled into scaffold applying

669    LACHESIS(Burton et al., 2013) with following parameters: CLUSTER MIN RE

670    SITES = 770, CLUSTER MAX LINK DENSITY=2, CLUSTER

671    NONINFORMATIVE RATIO = 2, ORDER MIN N RES IN TRUNK=578, ORDER

672    MIN N RES IN SHREDS=593.

673    **Annotation of transposable elements**

674    To predict transposable elements (TEs), whole genome sequences of the five

675    *Flaveria* species were searched for repetitive sequences individually. A de novo repeat

676    sequence library was constructed by RepeatModeler (RepeatModeler-Open-1.0.5) with

677    the following parameters: RepeatModeler -database database_name -engine ncbi -pa

678    [int]. Then, we used RepeatMasker (RepeatMasker-Open-4.1.0) to search for similar

679    TEs against the de novo library with the following parameters: RepeatMasker

680    genome.fa -lib de_novo_library -nolow -no_is -q -engine rmblast -pa [int] –norna.

681    Intact long terminal repeat retrotransposons (LTR-RTs) were identified using

682    LTR_FINDER (v1.07) (Xu and Wang, 2007) and LTRharvest (v1.5.10) (Ellinghaus et

683    al., 2008) with the default parameters. And Then LTR_Retriever (v2.9.0) (Ou and Jiang,

684    2018) was used to merge the above results with the parameters: LTR_retriever -genome

685    genome.fa -inharvest species.harvest.scn -infinder species.finder.scn –nonTGCA

686    species.harvest.nonTGCA.scn. The insertion time of intact LTR-RT was extracted from

687    LTR-Retriever analysis.

688 **Annotation of protein coding genes**

689      Gene models were predicted by combining de novo prediction, homology-based

690 and transcriptome-based strategies. Briefly, Augustus (v2.4) (Stanke and

691 Morgenstern, 2005), GlimmerHMM (v3.0.4) (Majoros et al., 2004), GeneID (v1.4)

692 (Parra et al., 2000) and Genscan (http://genes.mit.edu/GENSCAN.html) were used in

693 combination for de novo prediction. GeMoMa (v1.3.1) [Keilwagen et al., 2019] was used for

694 homology-based prediction. To facilitate gene annotation, from 18 to 32 Illumina

695 RNA-seq datasets were generated either in this study (for Flin, as depicted below) or

696 generated in our previous work (Zhu, 2020). Clean RNA-seq reads were mapped to

697 genome applying Hisat2 (v2.0.4) (Kim et al., 2019) and genome-based transcript

698 assembly was performed applying StringTie (v1.2.3) (Pertea et al., 2015) in default

699 parameters. Besides, de novo transcript assembly was conducted based on RNA-seq

700 data applying PASA (v2.0.2) (Haas et al., 2003) in default parameters. All predicted

701 gene structures were integrated into consensus gene models using EVidenceModeler

702 (v1.1.1) (Haas et al., 2008), and pseudo genes were predicted applying GeneWise

703 (v2.4.1) [Birney et al., 2004]，Coding sequence (CDS) failed to be translated either lacking

704 an open reading frame (ORF) or having premature stop codons were removed.

705      The completeness of protein repertoire was estimated in different aspects: 1) using

706 BUSCO (v3.0.2) (Seppey et al., 2019) against to viridiplantae reference, 2) RNA-seq

707 reads mapping to genome applying STAR (v2.7.3a) (Dobin et al., 2013), and 3) 150-bp

708 paired-ended DNA sequencing reads mapping to genome apply bowtie2 (v2.3.4.3)

709 (Langmead and Salzberg, 2012) (Supplemental Note 3).

710      Putative gene functions were assigned using the best match to GO, KEGG,

711 Swiss-Prot, TrEMBL and a non-redundant protein database (NR) using BLASTP

712 (v2.2.31+) (Camacho et al., 2009) with the E value threshold of 1e-5.

713      Transcriptional Factors were predicted using online website PlantTFDB (v5.0) (Jin

714 et al., 2017; Tian et al., 2020) (http://planttfdb.gao-lab.org/prediction.php). *Cis-*

715 regulatory elements (CREs) of promoter regions (3kb upstream of the start codon) were

716    predicted using Plantpan (v3.0) (Chow et al., 2019) with a score threshold of 0.99.

717    **Orthologous genes prediction and gene evolution**

718    To predict orthologous groups, protein coding genes from the five *Flaveria* species,

719    *Arabidopsis thaliana* (Atha), *Helianthus annuus* (Hann, sun flower), and *Lactuca sativa*

720    (Lsat, lettuce) were predicted applying Orthofinder (v2.3.11) (Emms and Kelly, 2019)

721    using default parameters. The protein sequences of Atha (TAIR10), Hann (v1.0) and

722    Lsat    (v7)    were    downloaded    from    Phytozome    (v13)

723    (https://phytozome.jgi.doe.gov/pz/portal.html). In case where there were multiple

724    alternative transcripts, the longest one was kept to represent the protein-coding gene.

725    **Phylogeny and divergence time analysis**

726    To construct the phylogenetic tree, CDS sequences of 1:1 orthologous genes were

727    aligned applying MUSCLE (v3.8.31) (Edgar, 2004) in default parameters. Alignments

728    of all the CDS were linked to make a super matrix, and RAxML (v7.9.3) (Stamatakis,

729    2006) was then applied for inferring phylogenetic tree using the following model: GTR

730    (General Time Reversible nucleotide substitution model) + GAMMA (variations in

731    sites follow GAMMA distribution) + I (a portion of Invariant sites in a sequence). To

732    calibrated the evolutionary time, CDS were aligned codon-wisely guided by protein

733    alignment using pal2nal (v14) (Suyama et al., 2006).The evolutionary time was

734    calibrated applying mcmctree in PAML package (v4.9) (Yang, 2007) using the

735    following parameters: seqtype=0 (nucleotides), clock=2 (independent), model = 0

736    (JC69). The reported fossil time between Hann and Lsat, *i.e.*, 34~40 million years as

737    inferred from timetree (http://timetree.org/) was used for calibration. The phylogenetic

738    tree    and    calibrated    evolutionary    time    were    displayed    using    FigTree

739    (http://tree.bio.ed.ac.uk/software/Figuretree/).

740    **Synteny analysis between *Flaveria* species**

741    To identify syntenic gene blocks in each species and between Frob with other

742     four species, all-against-all BLASTP (E value < 1e−10, top five matches) (v2.2.31+)

743     (Camacho et al., 2009) was performed for protein coding genes for each genome

744     pairs. Syntenic blocks were determined according to the presence of at least five

745     synteny gene pairs applying MCScanX (v0.8) (Wang et al., 2012) with default

746     parameters. Colinearity of the five species were drawn with JCVI

747     (https://github.com/tanghaibao/jcvi).Circular graphic was plotted using Circos (v0.69-

748     5).

749     **Investigation of light responsiveness of $C_4$ genes using qRT-PCR**

750     Consider that $C_4$ genes showed fast light responsiveness in $C_4$ species but not in

751     $C_3$ species (Burgess et al., 2016; Lyu et al., 2020), to verify the identified $C_4$ version

752     of $C_4$ genes, we investigated the changes of gene expression in response to light

753     induction using quantitative real time PCR (qRT-PCR). *Flaveria* species were put to

754     dark room at 6:00 pm. The dark-adapted plants were illuminated at 9:00 am the next

755     day. Fully expanded leaves, usually the $2^{nd}$ or $3^{rd}$ leaf pair counted from the top, were

756     cut after the leaves were illuminated for different time periods, *i.e.*, 0, 2, and 4 h, and

757     then flashed into liquid nitrogen. Samples were stored at −80°C before processing.

758     RNA isolation and qRT-PCR were performed as described earlier in(Lyu et al., 2020).

759     Relative transcript abundances were calculated by comparing to ACTIN7, the primers

760     used here were as depicted in our previous work (Zhu, 2020) .

761     **RNA-seq and transcriptional quantification for *Flaveria* species**

762     RNA-seq data of Flin were obtained from plant grown under low $CO_2$ (100 ppm)

763     *vs* normal $CO_2$ (380 ppm) for two weeks and four weeks respectively, and plant grown

764     under high light (with PPFD of 1400 μmol m$^{-2}$ s$^{-1}$) *vs* control light condition (500 μmol

765     m$^{-2}$ s$^{-1}$) were sequenced independently. Growth conditions were as depicted in (Zhu,

766     2020). For RNA extraction, the young fully expanded leaf usually situated on the $2^{nd}$ or

767     $3^{rd}$ pair of leaves counting started from the top was used. The chosen leaves were cut

768     and immediately frozen into liquid nitrogen and stored thereafter at -80 °C until further

769     processing. Total RNA was then isolated following the protocol of the PureLInk[TM] RNA

770     kit (ThermoFisher Scientific, USA). The RNA-sequencing was performed in the

771     Illumina platform in the paired-end mode with a read length of 150 bp. RNA-seq data

772     of other four species were from our previous work (Zhu, 2020).

773     To quantify the expression level of *Flaveria* genes, raw reads were trimmed

774     applying fastp (v0.20.0) (Chen et al., 2018) using default parameters. Transcript

775     abundance of gene were calculated by mapping RNA-seq reads to assembly genome

776     sequence of corresponding species using RSEM (v1.3.3) (Li and Dewey, 2011) in

777     default parameters, where STAR (v2.7.3a) (Dobin et al., 2013) was selected as the

778     mapping tool.

779     **Proteomics**

780     Mature leaves were cut from one-month old plant as depicted above, and leaves

781     were put into liquid nitrogen quickly. Frozen leaf samples were grinded thoroughly and

782     then incubated in lysis buffer (50 mM ammonium bicarbonate, 8M urea, 1mM DTT,

783     complete EDTA-free protease inhibitor cocktail (PIC) (Roche)). Samples were

784     centrifuged at 14,000g for 10 min at 4 °C. The supernatant was kept for total protein

785     samples. Total protein concentration was measured with a Bradford assay (Bradford,

786     1976).

787     The process of protein digestion, HPLC Fractionation, LC-MS/MS analysis and

788     data processing were detailed in Supplemental Note 8. Briefly, to generate data

789     dependent acquisition (DDA) library, peptides were prefractionated. Fractionated

790     peptides were mixed from all the 30 samples (a total of 200 μg). The mixture was

791     separated by a linear gradient, and finally, 30 fractions were mixed into 15

792     components. Raw data from each species were used to construct library based on

793     protein sequence from such species. As a result, five peptide libraries were obtained

794     with one for each species. Finally, data independent acquisition (DIA) was performed

795     using Spectronaut (version 14.7, Biognosys, Zurich, Switzerland). Default settings for

796  quantification at MS1 level were employed for quantification. The mass spectrometry

797  proteomics data have been deposited to the PRteomics IDEntifications Database

798  (PRIDE).

799  **ATAC-seq for the C$_4$ species Ftri**

800  To isolate nuclei from C$_4$ species Ftri, fully expanded mature leaves were harvest

801  at 1:00 pm. Around 3g fresh leaves from 5 plants were used for each of the two

802  biological replicates. Leaf materials were grinded in ice in 10 ml 4xNE buffer (40 mM

803  MES -KOH, PH5.4, 40 mM NaCl, 40 mM KCl, 10mM EDTA, 1M Sucrose, 0.1 mM

804  spermidine, 0.5mM spermine and 1mM DTT). Next, the debris was removed by

805  sieving through two layers of 70 μm nylon cell strainer into precooled flasks and then

806  the fluid were centrifuged at 200g at 4 °C for 3min to further remove debris. The

807  supernatant was centrifuged at 2000g at 4 °C for 5min to spin down Nuclei. Nuclei

808  were lysed by adding 1X NE buffer, 0.1% (v/v) NP40 and 0.1 (v/v) Tween-20 and

809  incubated on ice for 3 min. Nuclei were pelleted by centrifugation at 2000g at 4°C for

810  5 min. Pellets were then incubated in RS buffer (Tn5 mix, 10 mM Tris-HCL, PH 7.4,

811  10 mM NaCl, 3mM MgCl$_2$, 0.01% digitonin, 0.1% OM and 0.1% Tween-20) at 37 °C

812  for 30 min. The Tn5 tagmentation was then terminated under 95 °C for 2 min. DNA

813  was purified using a spin column (Qiagen, Germany) and then amplified using index

814  primers matching the Illumina Nextra adapter.

815  ATAC-seq libraries containing DNA insert between 50 and 150 bp were gel

816  purified and sequenced in Illumima X Ten platform in paired-end 150 bp mode. Raw

817  reads were trimmed using fastp (v0.20.0) (Chen et al., 2018) in default parameters.

818  Sequencing reads were mapped to genome sequence of Ftri (C$_4$) using bowtie2

819  (v2.3.4.3) (Langmead and Salzberg, 2012) in default parameters. Mapping result were

820  sorted using "sort" function in samtools (v1.11) (Li et al., 2009), and mapping with low

821  quality was filtered off using "view" function in samtools with -q=10. Duplicated

822  mapped reads were removed using "rmdup" function in samtools. Mapping peaks were

823    then called using macs2 (v2.2.7.1) (Zhang et al., 2008) using the following parameters:

824    -f BAMPE, -g 1.7e9 -q 0.05, --broad --nomodel --min-length 50. The parameter "broad"

825    was used to allow closed peaks merging into a broad peak. We referred peaks predicted

826    in this study as Tn5 hyper sensitive site (THS).

827    Peaks associated genes were assessed using "closest" function in bedtools (v2.29.2)

828    (Quinlan and Hall, 2010) with -k 2, considering the closest two genes (both upstream

829    and downstream). The distribution of THS relative to genome feature were assessed

830    using "computeMatrix" function in deeptools (v3.5.0) (Ramirez et al., 2014) with the

831    following parameters: --skipZeros –reference Point TSS -a 3000 -b 3000, the result was

832    then plotted using "plotHeatmap" in the same tool. To predict known motif of the THS,

833    the function fimo within meme package (v5.0.2) (Grant et al., 2011) was applied to scan

834    known motifs annotated in Plantpan 3.0 (Chow et al., 2019) through the sequences of

835    THS with default parameters.

836    **Comparison of intron-containing and intronless genes in *Flaveria* and other**

837    **species**

838    To calculate proportions of intronless genes in different species, we classified

839    intronless genes in different species based on gene annotations. A gene was classified

840    as intronless gene if all of its transcriptional isoforms contains exact one exon. For non-

841    *Flaveria* species studied here, genome sequences, gene annotation files and protein

842    sequences were downloaded from accessible databases. Briefly, those of Zea mays (v5)

843    were downloaded from Maize GDB (https://maizegdb.org/), those of *Spirogloea*

844    *muscicola* and *Mesotaenium endlicherianum* were downloaded from figshare

845    (figshare.com) referencing from (Cheng et al., 2019), and those of *Miscanthus*

846    *lutarioriparius* (Mlut) was downloaded from figshare referencing from (Miao et al.,

847    2021), and those of the rest species were downloaded from Phytozome (v13)

848    (https://phytozome-next.jgi.doe.gov/), with genome versions as following: Atha

849    (TAIR10), *Hann* (v1.2), *Glycine max* (v2.0), *Aquilegia coerulea* (v3.1), *Oropetium*

850 *thomaeum* (v1), *Setaris viridis* (Svir, v2.1), *Panicum Hallii* (Phal, v2.1), *Sorghum*

851 *bicolor* (Sbic, v3.1.1), Osat (v7), *Marchantia* polymorpha (v3.1), *Panicum halli* (v2.0),

852 *Chlamydomonas reinharditii* (Crei, v5.5), *Micromonas* pusilla (v3.0), *Coccomyxa*

853 subellipsoidea (v2.0), *Dunaliella* salina (v1.0), *Ostreococcus* lucimarinus (v2.0) and

854 *Volvox carteri* (v2.1).

855      To compare the transcript abundance of intron-containing and intronless genes for

856 non-*Flaveria* species, and to compare the expressional pre ferences of intron-containing

857 and intronless genes in mesophyll cells, we either inferred transcript abundances of

858 genes from published references or calculated gene expression levels based on

859 published RNA-seq datasets as detailed in Supplemental note 14. Specifically, we thank

860 Eric Schranz (Wageningen University), Andreas Weber (Heinrich Heine University)

861 and Julian Hibberd (Cambridge University) for access to the *Ggyn* genome sequence

862 and the updated Thas genome assembly (Cheng et al., 2013) for classifying intronless

863 genes and performing the RNA-seq quantification.

864 **Data availability.**

865      The genome assemblies, gene annotations, transcriptome data, proteomics data

866 and raw reads are available at China National GeneBank (CNGB)

867 (https://db.cngb.org/codeplot/datasets/public_dataset?id=flaveria) with project ID of

868 CPN0003058. The genome assemblies, gene annotations, transcriptome data,

869 proteomics data are also available at figshare

870 (https://figshare.com/account/home#/projects/114567). The genome assemblies are

871 also available at National Center for Biotechnology Information (NCBI) with accession

872 number SAMN14943594 for *F. robusta*, SAMN14943595 for *F. sonorensis*,

873 SAMN14943597 for *F. linearis*, SAMN14943596 for *F. ramosissima* and

874 SAMN14943598 for *F. trinervia*. The mass spectrometry proteomics data were

875 submitted to PRoteomics IDEntifications Database (PRIDE) with accession number

876 PXD024720 (username: reviewer_pxd024720@ebi.ac.uk, password: M6E7WzlM).

877    RNA-seq data of Flin were submitted to Gene Expression Omnibus (GEO) in the NCBI

878    database available with accession number: PRJNA827625. RNA-seq data of Frob, Fson,

879    Fram and Ftri are from published data with project accession PRJNA600545.

880

881    **Supplemental information**

882    1. **Supplemental Notes:** including supplemental note 1 ~ supplemental note 20, which

883       contain methods and results that support the main conclusion of the work.

884    2. **Supplemental Table and Figures**: including Table S1 and Figure S1 ~ Figure S9.

885

## References

886

887 Akyildiz, M., Gowik, U., Engelmann, S., Koczor, M., Streubel, M., and Westhoff, P.
888 (2007). Evolution and function of a cis-regulatory module for mesophyll-specific gene
889 expression in the $C_4$ dicot Flaveria trinervia. Plant Cell *19*, 3391-3402.

890 Aubry, S., Kelly, S., Kumpers, B.M., Smith-Unna, R.D., and Hibberd, J.M. (2014).
891 Deep evolutionary comparison of gene expression identifies parallel recruitment of
892 trans-factors in two independent origins of C4 photosynthesis. PLoS genetics *10*,
893 e1004365.

894 Billakurthi., K., Wrobel., T.J., Bräutigam., A., Weber., A.P.M., Westhoff., P., and Gowik.,
895 U. (2020). Transcriptome dynamics in developing leaves from C3 and C4 Flaveria
896 species reveal determinants of Kranz anatomy. BioRxiv.

897 Birney, E., Clamp, M., and Durbin, R. (2004). GeneWise and genomewise. Genome
898 research *14*, 988-995.

899 Bradford, M.M. (1976). Rapid and Sensitive Method for Quantitation of Microgram
900 Quantities of Protein Utilizing Principle of Protein-Dye Binding. Analytical
901 Biochemistry *72*, 248-254.

902 Burgess, S.J., Granero-Moya, I., Grange-Guermente, M.J., Boursnell, C., Terry, M.J.,
903 and Hibberd, J.M. (2016). Ancestral light and chloroplast regulation form the
904 foundations for $C_4$ gene expression. Nat Plants *2*, 16161.

905 Burgess, S.J., Reyna-Llorens, I., Stevenson, S.R., Singh, P., Jaeger, K., and Hibberd,
906 J.M. (2019). Genome-Wide Transcription Factor Binding in Leaves from $C_3$ and $C_4$
907 Grasses. Plant Cell *31*, 2297-2314.

908 Burton, J.N., Adey, A., Patwardhan, R.P., Qiu, R., Kitzman, J.O., and Shendure, J.
909 (2013). Chromosome-scale scaffolding of de novo genome assemblies based on
910 chromatin interactions. Nature biotechnology *31*, 1119-1125.

911 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and
912 Madden, T.L. (2009). BLAST plus : architecture and applications. BMC bioinformatics
913 *10*.

914 Chang, Y.M., Liu, W.Y., Shih, A.C., Shen, M.N., Lu, C.H., Lu, M.Y., Yang, H.W., Wang,
915 T.Y., Chen, S.C., Chen, S.M*., et al.* (2012). Characterizing regulatory and functional
916 differentiation between maize mesophyll and bundle sheath cells by transcriptomic
917 analysis. Plant physiology *160*, 165-177.

918 Chen, S., Zhou, Y., Chen, Y., and Gu, J. (2018). fastp: an ultra-fast all-in-one FASTQ
919 preprocessor. Bioinformatics *34*, i884-i890.

920 Cheng, S., van den Bergh, E., Zeng, P., Zhong, X., Xu, J., Liu, X., Hofberger, J., de
921 Bruijn, S., Bhide, A.S., Kuelahoglu, C*., et al.* (2013). The Tarenaya hassleriana genome
922 provides insight into reproductive trait and genome evolution of crucifers. Plant Cell
923 *25*, 2813-2830.

924 Cheng, S.F., Xian, W.F., Fu, Y., Marin, B., Keller, J., Wu, T., Sun, W.J., Li, X.L., Xu,
925 Y., Zhang, Y*., et al.* (2019). Genomes of Subaerial Zygnematophyceae Provide Insights
926 into Land Plant Evolution. Cell *179*, 1057-+.

927    Chow, C.N., Lee, T.Y., Hung, Y.C., Li, G.Z., Tseng, K.C., Liu, Y.H., Kuo, P.L., Zheng,
928    H.Q., and Chang, W.C. (2019). PlantPAN3.0: a new and updated resource for
929    reconstructing transcriptional regulatory networks from ChIP-seq experiments in plants.
930    Nucleic acids research *47*, D1155-D1163.

931    Christin, P.A., Besnard, G., Samaritani, E., Duvall, M.R., Hodkinson, T.R., Savolainen,
932    V., and Salamin, N. (2008). Oligocene CO2 decline promoted C4 photosynthesis in
933    grasses. Current biology : CB *18*, 37-43.

934    Christin, P.A., Boxall, S.F., Gregory, R., Edwards, E.J., Hartwell, J., and Osborne, C.P.
935    (2013). Parallel recruitment of multiple genes into C4 photosynthesis. Genome Biol
936    Evol *5*, 2174-2187.

937    Christin, P.A., Petitpierre, B., Salamin, N., Buchi, L., and Besnard, G. (2009). Evolution
938    of C4 phosphoenolpyruvate carboxykinase in Grasses, from genotype to phenotype.
939    Mol Biol Evol *26*, 357-365.

940    Cui, H. (2021). Challenges and Approaches to Crop Improvement Through C3-to-C4
941    Engineering. Frontiers in plant science *12*, 715391.

942    Da Silva, A.C., Sinnesael, M., Claeys, P., Davies, J., de Winter, N.J., Percival, L.M.E.,
943    Schaltegger, U., and De Vleeschouwer, D. (2020). Anchoring the Late Devonian mass
944    extinction in absolute time by integrating climatic controls and radio-isotopic dating.
945    Sci Rep *10*, 12940.

946    De Vleeschouwer, D., Da Silva, A.C., Sinnesael, M., Chen, D., Day, J.E., Whalen, M.T.,
947    Guo, Z., and Claeys, P. (2017). Timing and pacing of the Late Devonian mass extinction
948    event regulated by eccentricity and obliquity. Nature communications *8*, 2268.

949    Delaux, P.M., Hetherington, A.J., Coudert, Y., Delwiche, C., Dunand, C., Gould, S.,
950    Kenrick, P., Li, F.W., Philippe, H., Rensing, S.A.*, et al.* (2019). Reconstructing trait
951    evolution in plant evo-devo studies. Current biology : CB *29*, R1110-R1118.

952    Deng, C.L., Qin, R.Y., Gao, J., Cao, Y., Li, S.F., Gao, W.J., and Lu, L.D. (2012).
953    Identification of sex chromosome of spinach by physical mapping of 45s rDNAs by
954    FISH. Caryologia *65*, 322-327.

955    Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P.,
956    Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner.
957    Bioinformatics *29*, 15-21.

958    Du, H., and Liang, C. (2019). Assembly of chromosome-scale contigs by efficiently
959    resolving repetitive sequences with long reads. Nature communications *10*, 5360.

960    Edgar, R.C. (2004). MUSCLE: multiple sequence alignment with high accuracy and
961    high throughput. Nucleic acids research *32*, 1792-1797.

962    Ehleringer, J.R., Sage, R.F., Flanagan, L.B., and Pearcy, R.W. (1991). Climate Change
963    and the Evolution of C4 Photosynthesis. Trends Ecol Evol *6*, 95-99.

964    Ellinghaus, D., Kurtz, S., and Willhoeft, U. (2008). LTRharvest, an efficient and flexible
965    software for de novo detection of LTR retrotransposons. BMC bioinformatics *9*, 18.

966    Emms, D.M., Covshoff, S., Hibberd, J.M., and Kelly, S. (2016). Independent and
967    Parallel Evolution of New Genes by Gene Duplication in Two Origins of C4
968    Photosynthesis Provides New Insight into the Mechanism of Phloem Loading in C4

bioRxiv preprint doi: https://doi.org/10.1101/2022.10.14.512192; this version posted October 18, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY-NC-ND 4.0 International license.

969     Species. Mol Biol Evol *33*, 1796-1806.

970     Emms, D.M., and Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for

971     comparative genomics. Genome biology *20*, 238.

972     Feng, K., Hou, X.L., Xing, G.M., Liu, J.X., Duan, A.Q., Xu, Z.S., Li, M.Y., Zhuang, J.,

973     and Xiong, A.S. (2020). Advances in AP2/ERF super-family transcription factors in

974     plant. Crit Rev Biotechnol *40*, 750-776.

975     Finnegan, S., Heim, N.A., Peters, S.E., and Fischer, W.W. (2012). Climate change and

976     the selective signature of the Late Ordovician mass extinction. Proceedings of the

977     National Academy of Sciences of the United States of America *109*, 6829-6834.

978     Gan, Z., Chen, S.N., Huang, B., Hou, J., and Nie, P. (2017). Intronless and intron-

979     containing type I IFN genes coexist in amphibian Xenopus tropicalis: Insights into the

980     origin and evolution of type I IFNs in vertebrates. Dev Comp Immunol *67*, 166-176.

981     Gentles, A.J., and Karlin, S. (1999). Why are human G-protein-coupled receptors

982     predominantly intronless? Trends in genetics : TIG *15*, 47-49.

983     Gitzendanner, M.A., Soltis, P.S., Wong, G.K., Ruhfel, B.R., and Soltis, D.E. (2018).

984     Plastid phylogenomic analysis of green plants: A billion years of evolutionary history.

985     Am J Bot *105*, 291-301.

986     Gowik, U., and Westhoff, P. (2011). The path from $C_3$ to $C_4$ photosynthesis. Plant

987     physiology *155*, 56-63.

988     Grant, C.E., Bailey, T.L., and Noble, W.S. (2011). FIMO: scanning for occurrences of

989     a given motif. Bioinformatics *27*, 1017-1018.

990     Grzybowska, E.A. (2012). Human intronless genes: functional groups, associated

991     diseases, evolution, and mRNA processing in absence of splicing. Biochem Biophys

992     Res Commun *424*, 1-6.

993     Gu, C., Guo, Z.H., Hao, P.P., Wang, G.M., Jin, Z.M., and Zhang, S.L. (2017). Multiple

994     regulatory roles of AP2/ERF transcription factor in angiosperm. Bot Stud *58*, 6.

995     Gupta, S.D., Levey, M., Schulze, S., Karki, S., Emmerling, J., Streubel, M., Gowik, U.,

996     Paul Quick, W., and Westhoff, P. (2020). The $C_4$ Ppc promoters of many $C_4$ grass species

997     share a common regulatory mechanism for gene expression in the mesophyll cell. The

998     Plant journal : for cell and molecular biology *101*, 204-216.

999     Haas, B.J., Delcher, A.L., Mount, S.M., Wortman, J.R., Smith, R.K., Jr., Hannick, L.I.,

1000    Maiti, R., Ronning, C.M., Rusch, D.B., Town, C.D.*, et al.* (2003). Improving the

1001    Arabidopsis genome annotation using maximal transcript alignment assemblies.

1002    Nucleic acids research *31*, 5654-5666.

1003    Haas, B.J., Salzberg, S.L., Zhu, W., Pertea, M., Allen, J.E., Orvis, J., White, O., Buell,

1004    C.R., and Wortman, J.R. (2008). Automated eukaryotic gene structure annotation using

1005    EVidenceModeler and the program to assemble spliced alignments. Genome biology *9*.

1006    Hatch, M.D. (1987). $C_4$ photosynthesis - a unique blend of modified biochemistry,

1007    anatomy and ultrastructure. Biochimica Et Biophysica Acta *895*, 81-106.

1008    Hibberd, J.M., and Covshoff, S. (2010). The regulation of gene expression required for

1009    $C_4$ photosynthesis. Annual review of plant biology *61*, 181-207.

1010    Jin, J., Tian, F., Yang, D.C., Meng, Y.Q., Kong, L., Luo, J., and Gao, G. (2017).

39

PlantTFDB 4.0: toward a central hub for transcription factors and regulatory interactions in plants. Nucleic acids research *45*, D1040-D1045.

John, C.R., Smith-Unna, R.D., Woodfield, H., Covshoff, S., and Hibberd, J.M. (2014). Evolutionary convergence of cell-specific gene expression in independent lineages of C$_4$ grasses. Plant physiology *165*, 62-75.

Jorquera, R., Gonzalez, C., Clausen, P., Petersen, B., and Holmes, D.S. (2021). SinEx DB 2.0 update 2020: database for eukaryotic single-exon coding sequences. Database (Oxford) *2021*.

Ju, C.L., Van de Poel, B., Cooper, E.D., Thierer, J.H., Gibbons, T.R., Delwiche, C.F., and Chang, C.R. (2015). Conservation of ethylene as a plant hormone over 450 million years of evolution. Nature Plants *1*.

Keilwagen, J., Hartung, F., and Grau, J. (2019). GeMoMa: Homology-Based Gene Prediction Utilizing Intron Position Conservation and RNA-seq Data. Methods Mol Biol *1962*, 161-177.

Kim, D., Paggi, J.M., Park, C., Bennett, C., and Salzberg, S.L. (2019). Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. Nature biotechnology *37*, 907-915.

Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H., and Phillippy, A.M. (2017). Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. Genome research *27*, 722-736.

Kulahoglu, C., Denton, A.K., Sommer, M., Mass, J., Schliesky, S., Wrobel, T.J., Berckmans, B., Gongora-Castillo, E., Buell, C.R., Simon, R.*, et al.* (2014). Comparative transcriptome atlases reveal altered gene expression modules between two Cleomaceae C3 and C4 plant species. Plant Cell *26*, 3243-3260.

Kumar, S., Stecher, G., Suleski, M., and Hedges, S.B. (2017). TimeTree: A Resource for Timelines, Timetrees, and Divergence Times. Mol Biol Evol *34*, 1812-1819.

Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. Nature Methods *9*, 357-U354.

Li, B., and Dewey, C.N. (2011). RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. BMC bioinformatics *12*, 323.

Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics *25*, 1754-1760.

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., and Genome Project Data Processing, S. (2009). The Sequence Alignment/Map format and SAMtools. Bioinformatics *25*, 2078-2079.

Li, S.F., Guo, Y.J., Li, J.R., Zhang, D.X., Wang, B.X., Li, N., Deng, C.L., and Gao, W.J. (2019). The landscape of transposable elements and satellite DNAs in the genome of a dioecious plant spinach (Spinacia oleracea L.). Mob DNA *10*, 3.

Liu, H., Lyu, H.M., Zhu, K., Van de Peer, Y., and Max Cheng, Z.M. (2021). The emergence and evolution of intron-poor and intronless genes in intron-rich plant gene families. The Plant journal : for cell and molecular biology *105*, 1072-1082.

Long, S.P., Marshall-Colon, A., and Zhu, X.G. (2015). Meeting the global food demand

1053    of the future by engineering crop photosynthesis and yield potential. Cell *161*, 56-66.

1054    Lyu, M.J., Wang, Y., Jiang, J., Liu, X., Chen, G., and Zhu, X.G. (2020). What Matters

1055    for C4 Transporters: Evolutionary Changes of Phosphoenolpyruvate Transporter for C4

1056    Photosynthesis. Frontiers in plant science *11*, 935.

1057    Majoros, W.H., Pertea, M., and Salzberg, S.L. (2004). TigrScan and GlimmerHMM:

1058    two open source ab initio eukaryotic gene-finders. Bioinformatics *20*, 2878-2879.

1059    Marand, A.P., Chen, Z.L., Gallavotti, A., and Schmitz, R.J. (2021). A cis-regulatory

1060    atlas in maize at single-cell resolution. Cell *184*, 3041-+.

1061    Mascher, M., Gundlach, H., Himmelbach, A., Beier, S., Twardziok, S.O., Wicker, T.,

1062    Radchuk, V., Dockter, C., Hedley, P.E., Russell, J.*, et al.* (2017). A chromosome

1063    conformation capture ordered sequence of the barley genome. Nature *544*, 427-433.

1064    Maurino, V.G., and Weber, A.P. (2013). Engineering photosynthesis in plants and

1065    synthetic microorganisms. J Exp Bot *64*, 743-751.

1066    Merchante, C., Alonso, J.M., and Stepanova, A.N. (2013). Ethylene signaling: simple

1067    ligand, complex regulation. Current opinion in plant biology *16*, 554-560.

1068    Mergner, J., Frejno, M., List, M., Papacek, M., Chen, X., Chaudhary, A., Samaras, P.,

1069    Richter, S., Shikata, H., Messerer, M.*, et al.* (2020). Mass-spectrometry-based draft of

1070    the Arabidopsis proteome. Nature *579*, 409-414.

1071    Miao, J., Feng, Q., Li, Y., Zhao, Q., Zhou, C., Lu, H., Fan, D., Yan, J., Lu, Y., Tian, Q.*,

1072    et al.* (2021). Chromosome-scale assembly and analysis of biomass crop Miscanthus

1073    lutarioriparius genome. Nature communications *12*, 2458.

1074    Moreno-Villena, J.J., Dunning, L.T., Osborne, C.P., and Christin, P.A. (2018). Highly

1075    Expressed Genes Are Preferentially Co-Opted for C4 Photosynthesis. Mol Biol Evol

1076    *35*, 94-106.

1077    Ou, S., and Jiang, N. (2018). LTR_retriever: A Highly Accurate and Sensitive Program

1078    for Identification of Long Terminal Repeat Retrotransposons. Plant physiology *176*,

1079    1410-1422.

1080    Parra, G., Blanco, E., and Guigo, R. (2000). GeneID in Drosophila. Genome research

1081    *10*, 511-515.

1082    Pendleton, M., Sebra, R., Pang, A.W., Ummat, A., Franzen, O., Rausch, T., Stutz, A.M.,

1083    Stedman, W., Anantharaman, T., Hastie, A.*, et al.* (2015). Assembly and diploid

1084    architecture of an individual human genome via single-molecule technologies. Nat

1085    Methods *12*, 780-786.

1086    Pertea, M., Pertea, G.M., Antonescu, C.M., Chang, T.C., Mendell, J.T., and Salzberg,

1087    S.L. (2015). StringTie enables improved reconstruction of a transcriptome from RNA-

1088    seq reads. Nature biotechnology *33*, 290-295.

1089    Powell, A.M. (1978). Systematics of Flaveria (Flaveriinae Asteraceae). Ann Mo Bot

1090    Gard *65*, 590-636.

1091    Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for

1092    comparing genomic features. Bioinformatics *26*, 841-842.

1093    Ramirez, F., Dundar, F., Diehl, S., Gruning, B.A., and Manke, T. (2014). deepTools: a

1094    flexible platform for exploring deep-sequencing data. Nucleic acids research *42*, W187-

41

1095    191.

1096    Sage, R.F. (2004). The evolution of C4 photosynthesis. New Phytologist, 341-370.

1097    Sage, R.F. (2016). A portrait of the C4 photosynthetic family on the 50th anniversary
1098    of its discovery: species number, evolutionary lineages, and Hall of Fame. J Exp Bot
1099    *67*, 4039-4056.

1100    Sage, R.F., Christin, P.A., and Edwards, E.J. (2011). The C4 plant lineages of planet
1101    Earth. J Exp Bot *62*, 3155-3169.

1102    Sage, R.F., Sage, T.L., and Kocacinar, F. (2012). Photorespiration and evolution of $C_4$
1103    photosynthesis. Annual Review of Plant Biologist *63*, 19-47.

1104    Sage, T.L., Busch, F.A., Johnson, D.C., Friesen, P.C., Stinson, C.R., Stata, M.,
1105    Sultmanis, S., Rahman, B.A., Rawsthorne, S., and Sage, R.F. (2013). Initial events
1106    during the evolution of C4 photosynthesis in C3 species of Flaveria. Plant physiology
1107    *163*, 1266-1276.

1108    Sanderson, M.J., Thorne, J.L., Wikstrom, N., and Bremer, K. (2004). Molecular
1109    evidence on plant divergence times. American Journal of Botany *91*, 1656-1665.

1110    Schluter, U., and Weber, A.P.M. (2020). Regulation and Evolution of $C_4$ Photosynthesis.
1111    Annual review of plant biology *71*, 183-215.

1112    Seppey, M., Manni, M., and Zdobnov, E.M. (2019). BUSCO: Assessing Genome
1113    Assembly and Annotation Completeness. Methods Mol Biol *1962*, 227-245.

1114    Servant, N., Varoquaux, N., Lajoie, B.R., Viara, E., Chen, C.J., Vert, J.P., Heard, E.,
1115    Dekker, J., and Barillot, E. (2015). HiC-Pro: an optimized and flexible pipeline for Hi-
1116    C data processing. Genome biology *16*, 259.

1117    Shabalina, S.A., Ogurtsov, A.Y., Spiridonov, A.N., Novichkov, P.S., Spiridonov, N.A.,
1118    and Koonin, E.V. (2010). Distinct patterns of expression and evolution of intronless and
1119    intron-containing mammalian genes. Mol Biol Evol *27*, 1745-1749.

1120    Sheehan, P.M. (2001). The Late Ordovician mass extinction. Annual Review of Earth
1121    and Planetary Sciences *29*, 331-364.

1122    Slack, C.R., and Hatch, M.D. (1967). Comparative studies on the activity of
1123    carboxylases and other enzymes in relation to the new pathway of photosynthetic
1124    carbon dioxide fixation in tropical grasses. The Biochemical journal *103*, 660-665.

1125    Stamatakis, A. (2006). RAxML-VI-HPC: Maximum likelihood-based phylogenetic
1126    analyses with thousands of taxa and mixed models. Bioinformatics *22*, 2688-2690.

1127    Stanke, M., and Morgenstern, B. (2005). AUGUSTUS: a web server for gene prediction
1128    in eukaryotes that allows user-defined constraints. Nucleic acids research *33*, W465-
1129    467.

1130    Steven Kelly, S.C., Samart Wanchana, Vivek Thakur, W. Paul Quick,YuWang, Martha
1131    Ludwig, Richard Bruskiewich, Alisdair R. Fernie, Rowan F. Sage8,Zhijian
1132    Tian,Zixiang Yan, Jun Wang, Yong Zhang, Xin-Guang Zhu, Gane Ka-Shu Wong, Julian
1133    M. Hibberd (2018). Wide sampling of natural diversity identifies novel molecular
1134    signatures of $C_4$ photosynthesis **BioRxiv**.

1135    Suyama, M., Torrents, D., and Bork, P. (2006). PAL2NAL: robust conversion of protein
1136    sequence alignments into the corresponding codon alignments. Nucleic acids research

1137    *34*, W609-612.

1138    Taniguchi, Y.Y., Gowik, U., Kinoshita, Y., Kishizaki, R., Ono, N., Yokota, A., Westhoff,

1139    P., and Munekage, Y.N. (2021). Dynamic changes of genome sizes and gradual gain of

1140    cell-specific distribution of C4 enzymes during C4 evolution in genus Flaveria. Plant

1141    Genome, e20095.

1142    Tian, F., Yang, D.C., Meng, Y.Q., Jin, J., and Gao, G. (2020). PlantRegMap: charting

1143    functional regulatory maps in plants. Nucleic acids research *48*, D1104-D1113.

1144    Tu, X., Mejia-Guerra, M.K., Valdes Franco, J.A., Tzeng, D., Chu, P.Y., Shen, W., Wei,

1145    Y., Dai, X., Li, P., Buckler, E.S.*, et al.* (2020). Reconstructing the maize leaf regulatory

1146    network using ChIP-seq data of 104 transcription factors. Nature communications *11*,

1147    5089.

1148    Vogan, P.J., and Sage, R.F. (2011). Water-use efficiency and nitrogen-use efficiency of

1149    $C_3$ -$C_4$ intermediate species of Flaveria Juss. (Asteraceae). Plant, cell & environment

1150    *34*, 1415-1430.

1151    Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo,

1152    C.A., Zeng, Q., Wortman, J., Young, S.K.*, et al.* (2014). Pilon: an integrated tool for

1153    comprehensive microbial variant detection and genome assembly improvement. PLoS

1154    One *9*, e112963.

1155    Wang, Y., Tang, H., Debarry, J.D., Tan, X., Li, J., Wang, X., Lee, T.H., Jin, H., Marler,

1156    B., Guo, H.*, et al.* (2012). MCScanX: a toolkit for detection and evolutionary analysis

1157    of gene synteny and collinearity. Nucleic acids research *40*, e49.

1158    Wei, S., Li, X., Lu, Z., Zhang, H., Ye, X., Zhou, Y., Li, J., Yan, Y., Pei, H., Duan, F.*, et*

1159    *al.* (2022). A transcriptional regulator that boosts grain yields and shortens the growth

1160    duration of rice. Science *377*, eabi8455.

1161    Westhoff, P., and Gowik, U. (2010). Evolution of C4 photosynthesis--looking for the

1162    master switch. Plant physiology *154*, 598-601.

1163    Williams, B.P., Aubry, S., and Hibberd, J.M. (2012). Molecular evolution of genes

1164    recruited into C4 photosynthesis. Trends in plant science *17*, 213-220.

1165    Xie, Z., Nolan, T.M., Jiang, H., and Yin, Y. (2019). AP2/ERF Transcription Factor

1166    Regulatory Networks in Hormone and Abiotic Stress Responses in Arabidopsis.

1167    Frontiers in plant science *10*, 228.

1168    Xu, J., Brautigam, A., Weber, A.P., and Zhu, X.G. (2016). Systems analysis of cis-

1169    regulatory motifs in $C_4$ photosynthesis genes using maize and rice leaf transcriptomic

1170    data during a process of de-etiolation. J Exp Bot *67*, 5105-5117.

1171    Xu, Z., and Wang, H. (2007). LTR_FINDER: an efficient tool for the prediction of full-

1172    length LTR retrotransposons. Nucleic acids research *35*, W265-W268.

1173    Yang, Z. (2007). PAML 4: phylogenetic analysis by maximum likelihood. Mol Biol

1174    Evol *24*, 1586-1591.

1175    Zhang, Y., Liu, T., Meyer, C.A., Eeckhoute, J., Johnson, D.S., Bernstein, B.E., Nusbaum,

1176    C., Myers, R.M., Brown, M., Li, W.*, et al.* (2008). Model-based analysis of ChIP-Seq

1177    (MACS). Genome biology *9*, R137.

1178    Zhu, M.-J.A.L.J.E.F.C.G.C.X.-G. (2020). Evolution of co-regulatory network of C4

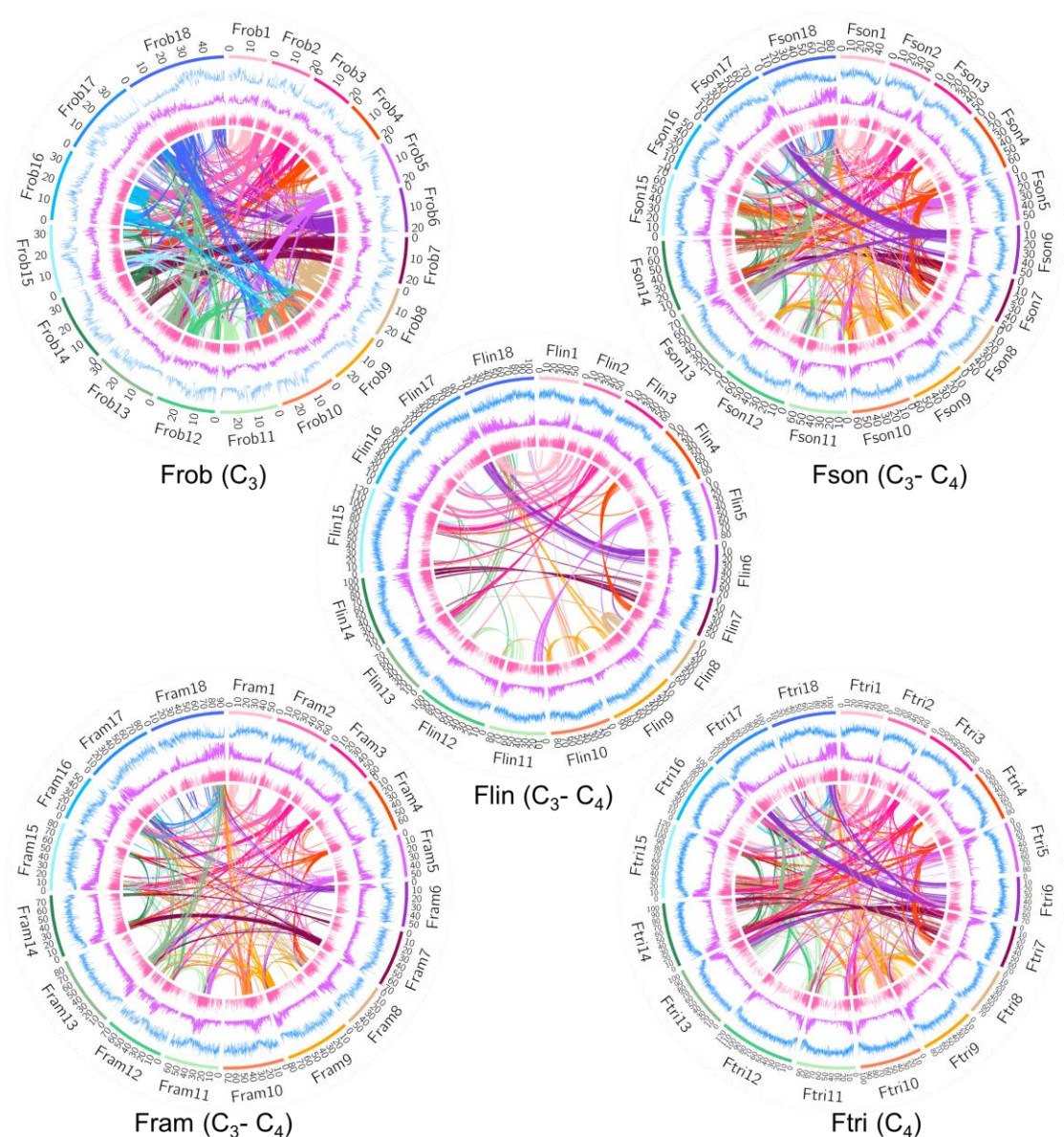1179 metabolic genes and TFs in the genus   Flaveria: go anear or away in the intermediate
1180 species? BioRxiv.

1181 Zhu, X.-G., Shan, L., Wang, Y., and Quick, W.P. (2010). $C_4$ Rice - an ideal arena for
1182 systems biology research. Journal of Integrative Plant Biology *52*, 762-770.

1183 Zhu, X.G., Long, S.P., and Ort, D.R. (2008). What is the maximum efficiency with
1184 which photosynthesis can convert solar energy into biomass? Current opinion in
1185 biotechnology *19*, 153-159.

1186
1187

## 1188 **Supplemental Table and Figures**

## 1189 **Table S1. Statistics of genome assemblies and annotations**

| Species | *F. robusta* | *F. sonorensis* | *F. linearis* | *F. ramosissima* | *F. trinervia* |
|---|---|---|---|---|---|
| Photosynthetic type | $C_3$ | $C_3$-$C_4$ | $C_3$-$C_4$ | $C_3$-$C_4$ | $C_4$ |
| Genome size (GB) | 0.55 | 1.26 | 1.66 | 1.42 | 1.8 |
| Genome size estimated by flow cytometry (GB) | 0.45 | 1.2 | 1.86 | 1.62 | 1.65 |
| anchored to chromosome (%) | 92.2 | 92 | 91.2 | 94.3 | 93.4 |
| Contig N50 (MB) | 7.9 | 1.8 | 1.2 | 0.76 | 1.6 |
| GC content (%) | 32.87 | 36.06 | 37.45 | 37.3 | 37.85 |
| BUSCO% | 99.2 | 98.1 | 92.5 | 97 | 95.1 |
| Gene number | 35,875 | 37,028 | 38,652 | 34,029 | 32,915 |
| Average gene length (bp, intron+exon) | 3564.67 | 3670.95 | 3973.67 | 3971.15 | 3555.47 |
| Average Exons number per gene | 5.53 | 5.54 | 5.44 | 5.93 | 5.66 |

1190

45

**Figure S1. Genome features of five *Flaveria* species**

The circular representation of pseudochromosomes. From outer to inner side: blue: LTR density per million base pair (Mb), purple: exon density per Mb, pink: transcript abundance per gene in log10 TPM (transcript per million mapped reads). Lines in the inner circle represent links between synteny-selected paralogs. (Abbreviations: Frob: *F. robusta*, Fson: *F. sonorensis*, Flin: *F. linearis*, Fram: *F. ramosissima*, Ftri: *F. trinervia*)

**Figure S2. $C_4$ gene gradually gained light responses during evolution**

Real-time quantitative (qRT)-PCR was used to quantify the transcript abundance of $C_4$ enzymes in mature leaves after 0, 2 and 4h upon illumination. significance levels were calculated using $t$-test. (*: 0.05–0.01, **: 0.01–0.001, ***: < 0.001) (Abbreviations: CA1, carbonic anhydrase 1; PEPC1, phospho*enol*pyruvate carboxylase 1; PEPC-k: PEPC kinase; NADP-MDH, NADP-dependent malate dehydrogenase; NADP-ME4, NADP-dependent malic enzyme 4; PPDK, pyruvate/orthophosphate dikinase; PPDK-RP, PPDK regulatory protein)

**Figure S3. $C_4$ version of CA and PEPC-k show more copies in the $C_4$ species Ftri resulting from tandem duplication**

(a) and (b) illustrate the gene tree of CA and PEPC-k respectively. Gene tree were constructed based alignment of protein sequences. Bootstrap scores were from 100 bootstrap samplings. Bars show gene expressions of leaves from two-month-old plants grown in low $CO_2$ condition (100 ppm) *vs* normal $CO_2$ condition (380 ppm) for four weeks. Three biological replicates were performed for each condition. (c) and (d) Integrated Genome Viewer (IGV) of RNA-seq reads and ATAC-seq reads of two copies of CA1 and four copies of PEPC-k1 anchored to chromosomes in Ftri respectively. Tn5 cuts and transposase hypersensitive sites (THS) from two biological replicates are showed. (Abbreviations: CA1: carbonic anhydrase1; PEPC-k1: phospho*enol*pyruvate carboxylase kinase1.)

**Figure S4. Predicted *cis*-regulatory elements in the C₄ species Ftri by applying ATAC-seq**

(a) Integrated Genome Viewer (IGV) of RNA-seq reads and two-biological replicates of ATAC-seq reads in Ftri are showed in three spatial resolutions, *i.e.*, genome scale with chromosome number and location indicated (top), chromosomal scale (middle), and million-base genomic region including PEPC1. (b) Enriched *cis*-regulatory elements (CREs) in three types of accessible chromatin regions (ACR-CREs), *i.e.*, genic (gACR-CREs: overlapping a gene), upstream (upACR-CREs: within 3kb upstream of the start codon of a gene) and downstream (down ACRs-CRES: within 3kb downstream of the stop codon of a gene). (c) Enriched ACR-CREs associated with photosynthetic genes and photorespiratory genes. (Abbreviations: ATAC-seq: transposase-accessible chromatin using sequencing; ACR: accessible chromatin regions; CREs: *cis*-regulatory elements)
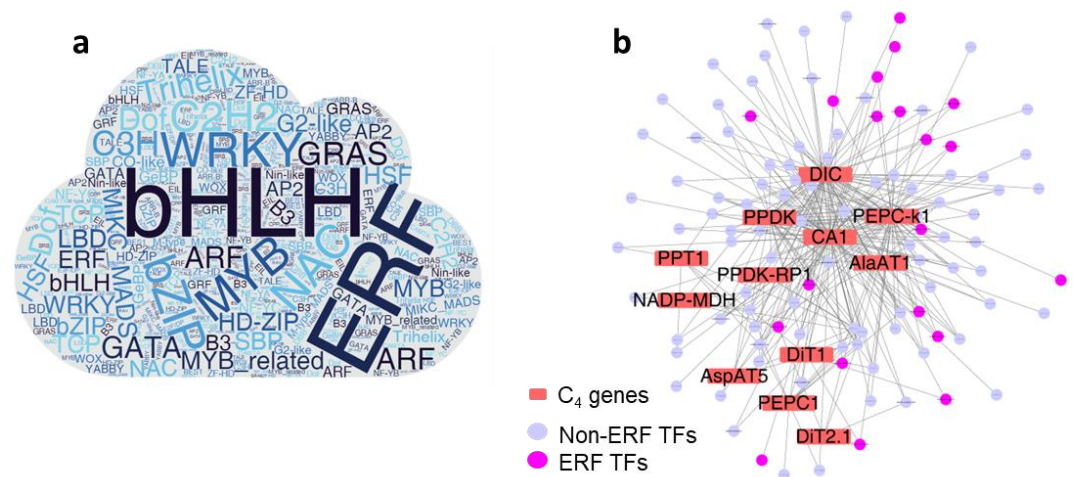
49

**Figure S5. The number and transcript abundances of intron-containing and intro, genes in each TF family**
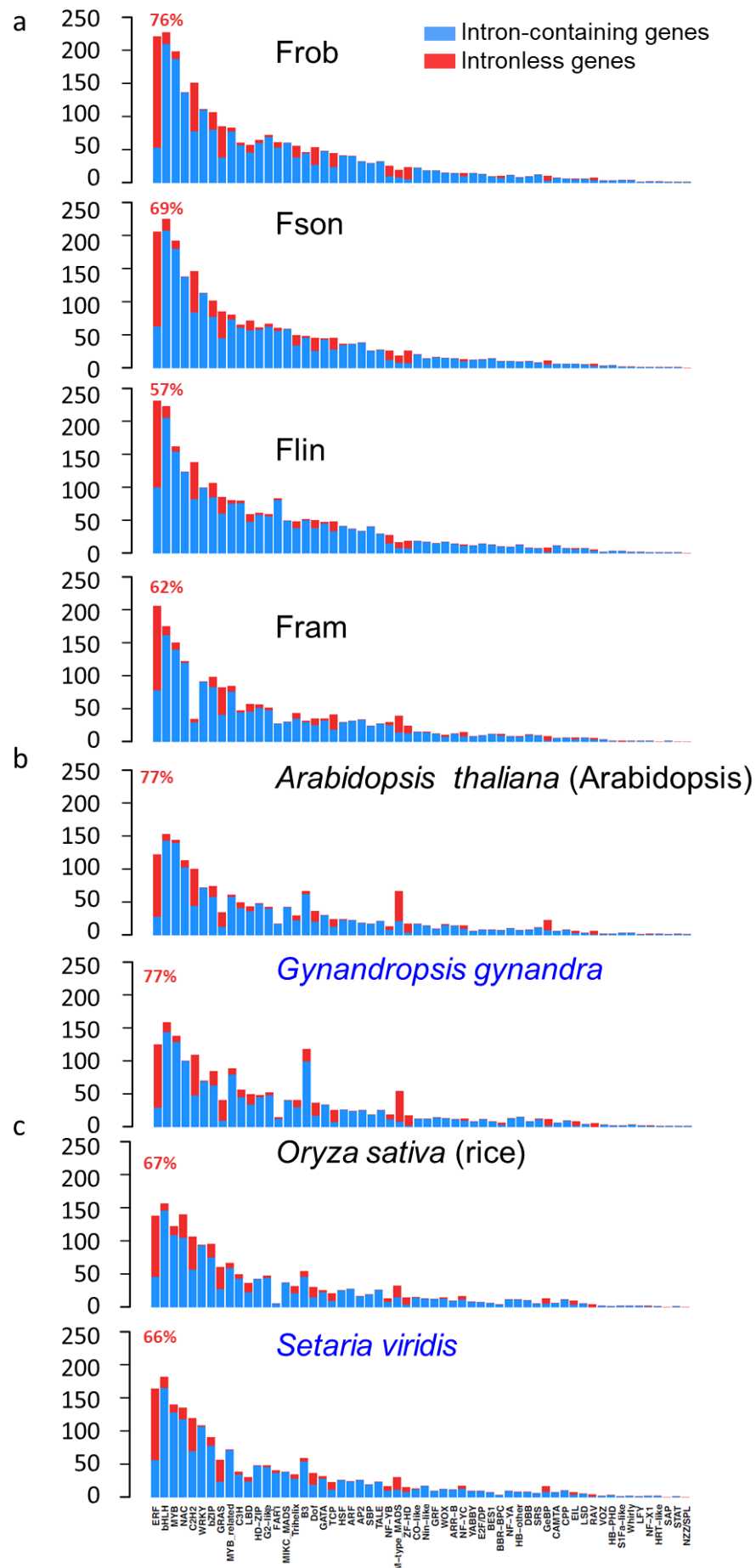
Heat maps show the gene number and transcript abundance of genes in each TF family from all the annotated TFs (left panel) and from $C_4$GRN (right panel). The size of circle represents the number of genes, and the color represents the log2 transformed transcript abundances in transcript per million mapped reads (TPM).

1246

**Figure S6. ERF TF were recruited by C$_4$ genes in *Zea mays***

(a) Word cloud shows the frequencies of TFs in each TF families in the leaf gene regulatory network (GRN) of *Zea mays* (Zmay, corn) from (Tu et al., 2020). In the whole leaf GRN, bHLH is the most prevalent TFs with 138 genes. (b) The C$_4$GRN in Zmay that includes C$_4$ genes and their regulatory TFs. ERF is the most abundant TFs in the C$_4$GRN as showed in purple circle. 12 C$_4$ gene are included within the whole leaf GRN. (Abbreviation: Zmay: *Zea mays*)

1254

1255

1256    **Figure S7. The number of intronless genes in each TF family in different species**
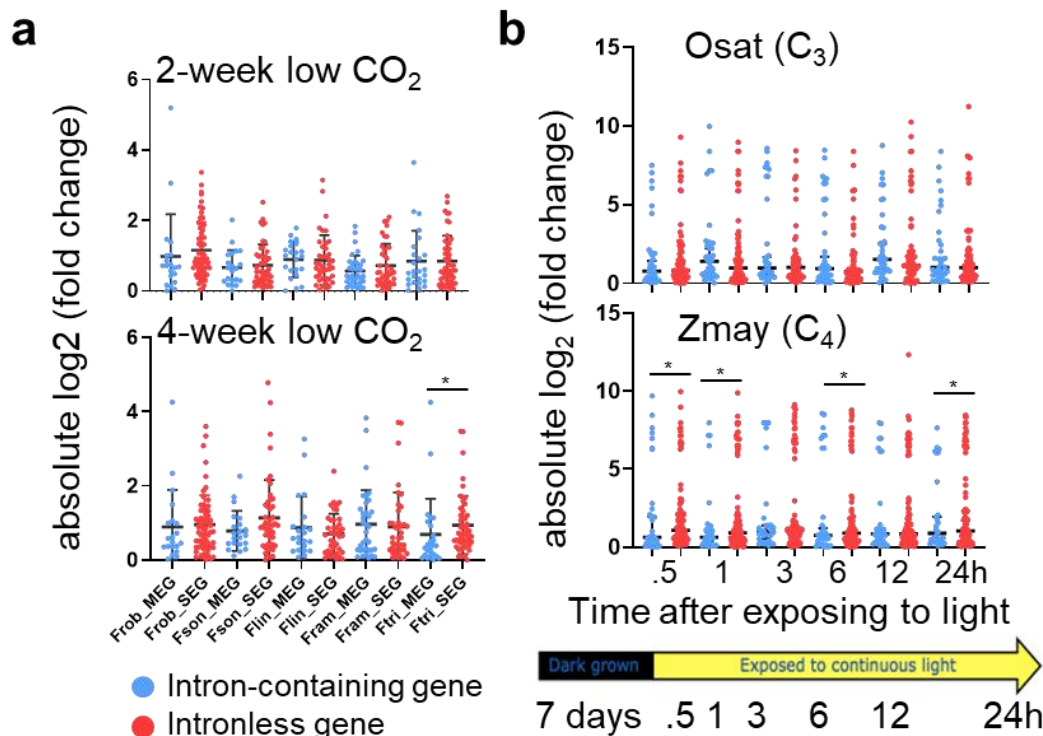
1257    The number of intronless gene and intron-contain genes are showed in each TF family

1258    from (a) four *Flaveria* species, (b) two dicotyledonous species and (c) two

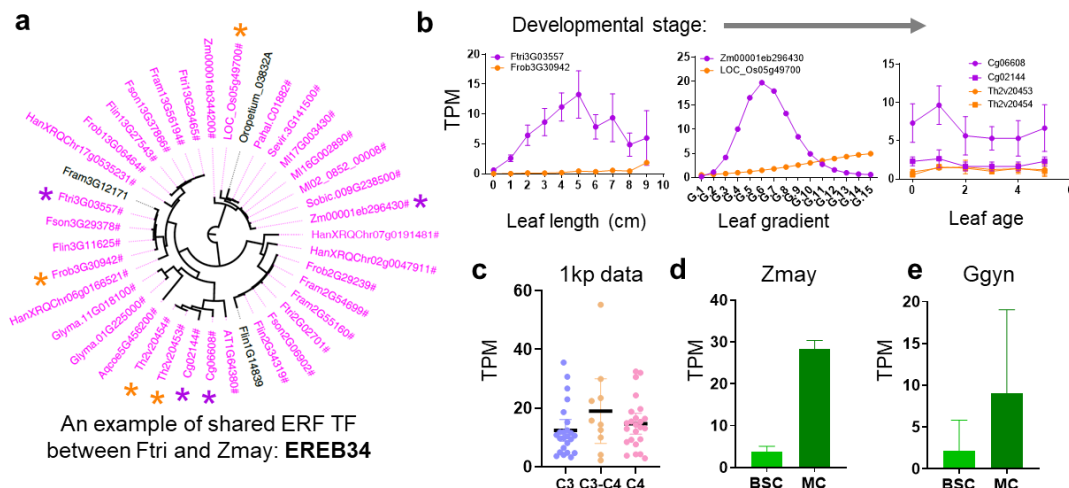1259    monocotyledonous species. $C_4$ species are labeled in blue font. Proportions of intronless

1260    genes in ERF TF family are showed in red fond for each species.

1261

**Figure S8. The response of intronless and intron-containing ERF TFs to environmental changes**

(a) The changes on transcript abundances of intronless and intron-containing ERF TFs in five *Flaveria* species in response to low $CO_2$ (100 ppm) compared to normal $CO_2$ (380 ppm). RNA-seq data of both low $CO_2$ and normal $CO_2$ grown plants were taken from leaves after plants being grown under each condition for two weeks and four weeks respectively. (b) The response of transcript abundances of intronless and intron-containing ERF TFs in *Oryza sativa* (Osat, $C_3$) and *Zea mays* (Zmay, $C_4$) under light induction. The gene expression data of Osat and Zmay are from (Xu et al., 2016). Seeds of both species were germinated and grown under dark for 7 days. RNA-seq data of leaves were taken before light, 0.5h, 1, 3h, 6h, 12h and 24h after light respectively. The fold change of each time point was calculated as the ratio of gene expression level of this time point to that of the prior time point. Gene expression levels for all analysis was showed in transcript per million mapped reads (TPM). (Abbreviations: MEG: multi-exon genes, *i.e.*, intron-containing genes; SEG: single exon genes, *i.e.*, intronless genes; Osat: *Oryza sativa*; Zmay: *Zea mays*.)

54

**Figure S9. An example of intronless ERF TF that was recruited to regulate C4 genes in both Ftri and Zmay**

**(a)** Gene tree of EREB34. Intronless genes are labeled in purple. EREB34 orthologous genes among 23 species (see Figure5) were predicted applying Orthofiner. EREB34 present in higher plants but not in algae or liverwort. Genes marked with orange and purple stars are from $C_3$ and $C_4$ plants that compared in transcript abundances in **(b)**. **(b)** Comparisons of EREB34 in transcript abundances between Frob ($C_3$) *vs* Ftri ($C_4$) (lef), *Oryza sative* (Osat, $C_3$) *vs Zea mays* (Zmay, $C_4$) (middle), and *Tarenaya hassleriana* (Thas, $C_3$) vs *Gynandropsis gynandra* (Ggyn, $C_4$) (right) along leaf developmental gradient. Genes from $C_3$ species are labeled in orange, and those from $C_4$ species are in purple. RNA-seq data of these species are from published sources, *i.e.*, data of Frob and Ftri are from (Billakurthi. et al., 2020), data of Osat and Zmay are from (Xu et al., 2016), data of Thas and Ggyn are from (Kulahoglu et al., 2014) . The first points in Frob and Ftri (0) represent meristem. Leaf ages of Thas and Ggyn are as following: **0**: 0-2 days (d); **1**: 2-4 d; **2**: 4-6 d; **3**: 6-8 d; **4**: 8-10 d and **5**: 10-12 d. **(c)** Transcript abundances of EREB34 in $C_3$, $C_3$-$C_4$ and $C_4$ species from one thousand plants (1kp) project, covering 18 independent $C_4$ lineages. RNA-seq data are from (Steven Kelly, 2018). **(d)** Transcript abundance of EREB34 in Zmay bundle sheath cell (BSC) and mesophyll cell (MC). Expressional data are from (Chang et al., 2012). **(e)** Transcript abundances of EREB34 in Ggyn BSC and MC. Expressional data are from (Aubry et al., 2014). (Abbreviations: MC: mesophyll cell; BSC: bundle sheath cell.)