



# Research Repository

## **Ancient duons may underpin spatial patterning of gene expression in C<sub>4</sub> leaves**

Accepted for publication in Proceedings of the National Academy of Sciences.

Research Repository link: <https://repository.essex.ac.uk/40006/>

### **Please note:**

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the published version if you wish to cite this paper.

<https://doi.org/10.1073/pnas.1720576115>

1    **Ancient duons may underpin spatial patterning of gene expression in C<sub>4</sub> leaves**  
2  
3    Ivan Reyna-Llorens, Steven J. Burgess, Gregory Reeves, Pallavi Singh, Sean R. Stevenson, Ben  
4    P. Williams, Susan Stanley and Julian M. Hibberd  
5  
6    Department of Plant Sciences, Downing Street, University of Cambridge, Cambridge CB2 3EA, UK.  
7  
8    I-RL - suallorems@gmail.com  
9    SJB - sburgess011@gmail.com  
10   GR - gr360@cam.ac.uk  
11   PS - ps753@cam.ac.uk  
12   SRS - srs62@cam.ac.uk  
13   BPM - bp.williams0@gmail.com  
14   SS - ss383@cam.ac.uk  
15   JMH (corresponding) - jmh65@cam.ac.uk. Tel - 44(0) 1223 766547  
16

## 17    **Abstract**

18    If the highly efficient C<sub>4</sub> photosynthesis pathway could be transferred to crops with the C<sub>3</sub> pathway  
19    there could be yield gains of up to 50%. It has been proposed that the multiple metabolic and  
20    developmental modifications associated with C<sub>4</sub> photosynthesis are underpinned by relatively few  
21    master regulators that have allowed the evolution of C<sub>4</sub> photosynthesis more than sixty times in  
22    flowering plants. Here we identify a component of one such regulator that consists of a pair of *cis*-  
23    elements located in coding sequence of multiple genes that are preferentially expressed in bundle  
24    sheath cells of C<sub>4</sub> leaves. These motifs represent duons as they play a dual role in coding for amino  
25    acids as well as controlling the spatial patterning of gene expression associated with the C<sub>4</sub> leaf.  
26    They act to repress transcription of C<sub>4</sub> photosynthesis genes in mesophyll cells. These duons are  
27    also present in the C<sub>3</sub> model *Arabidopsis thaliana*, and in fact, are conserved in all land plants and  
28    even some algae that use C<sub>3</sub> photosynthesis. C<sub>4</sub> photosynthesis therefore appears to have co-opted  
29    an ancient regulatory code to generate the spatial patterning of gene expression that is a hallmark  
30    of C<sub>4</sub> photosynthesis. This novel intragenic transcriptional regulatory sequence could be exploited  
31    in the engineering of efficient photosynthesis of crops.

32

## 33    **Significance statement**

34    We describe an ancient regulatory code that patterns the gene expression required for the efficient  
35    C<sub>4</sub> pathway. This code is based on two novel regulatory elements located in exonic sequence that  
36    act co-operatively to repress transcription in specific cells. As these regulators are located in gene  
37    bodies they determine amino acid sequence as well as gene expression and so are known as duons.  
38    They are found in multiple genes, are not restricted to C<sub>4</sub> species, but rather found in all land plants  
39    and even some algae. The prevalence of these motifs has likely facilitated the repeated evolution of  
40    the complex C<sub>4</sub> system, and they also provide a mechanism that could be used to engineer  
41    photosynthesis.

42

43

44 /body

45

## 46 **Introduction**

47 Photosynthesis forms the basis of life on earth. When plants moved onto land they inherited a  
48 photosynthetic system developed by bacteria in which Ribulose Bisphosphate Carboxylase  
49 Oxygenase (RuBisCO) generates phosphoglyceric acid (PGA)<sup>1</sup>. As PGA contains three carbon  
50 atoms, this is referred to as C<sub>3</sub> photosynthesis. A side-reaction of RuBisCO fixes O<sub>2</sub> rather than CO<sub>2</sub>  
51 and generates the toxic compound phosphoglycolate. Although plants use photorespiration to  
52 remove phosphoglycolate, both carbon and energy are lost in the process<sup>2</sup>. Around 30 million years  
53 ago, some species evolved a system in which CO<sub>2</sub> is concentrated around RuBisCO such that  
54 oxygenation is minimised, and photosynthetic efficiency increases by around 50%<sup>3,4</sup>. These species  
55 now represent the most productive vegetation on the planet<sup>5,6</sup>, and because they initially generate a  
56 C<sub>4</sub> acid in the photosynthetic process, are known as C<sub>4</sub> plants.

57 The C<sub>4</sub> mechanism depends on spatial separation of photosynthetic reactions, most commonly  
58 between mesophyll (M) and bundle sheath (BS) cells. How this complex process has evolved in over  
59 66 independent lineages remains a mystery. The finding that expression of multiple *NAD-dependent*  
60 *MALIC ENZYME (NAD-ME)* genes in the BS of C<sub>4</sub> *Gynandropsis gynandra* is dependent on coding  
61 sequence<sup>7</sup> hinted at the importance of exons in this process, which are more highly conserved than  
62 intragenic sequences and therefore a potential hotspot underpinning repeated evolution. Although  
63 genome-wide studies commonly report transcription factor binding within genes<sup>8,9</sup>, there is little  
64 functional validation to support these findings and their general importance remains unclear.

65

## 66 **Results**

67 To better understand mechanisms responsible for generating preferential expression in BS cells of  
68 C<sub>4</sub> plants analysis was focussed on the coding region of *GgNAD-ME1*. Although a 240 nucleotide  
69 fragment under control of the constitutive CaMV35S promoter confers BS accumulation in *G.*  
70 *gynandra* (Figure 1A&B; Figure S1,2) it was unclear if this was due to transcriptional and post-

transcriptional mechanisms. These can be distinguished using an antisense construct that maintains DNA sequence, but when transcribed generates a complementary mRNA. An antisense construct under the constitutive CaMV35S promoter maintained preferential accumulation in the BS (Figure 1A&B) indicating that DNA sequence is recognised by *trans*-acting factors regardless of the orientation of the sequence. The preferential rather than exclusive accumulation of GUS in BS cells reported here is consistent with quantification of transcripts encoding components of the C<sub>4</sub> cycle in BS and M cells of C<sub>4</sub> leaves<sup>7,10</sup>, and also previous analysis of elements controlling gene expression<sup>11–13</sup>. Thus, the coding sequence of *GgNAD-ME1* controls transcription, and as it suppresses activity of the constitutive 35SCaMV promoter in M cells, the simplest explanation is that this region interacts with a repressive transcription factor.

Without definition of the motifs within *NAD-ME* genes specifying expression in the BS, it was not possible to determine if they control expression of additional genes, nor to understand if the same elements are used in other species. To identify specific nucleotides responsible for BS expression a deletion series was generated (Figure 1A). Deletion of 24 or 78 nucleotides from the 3' end did not affect preferential accumulation in BS cells (Figure 1A-B) but removal of 90 nucleotides did (Figure 1A-C). Similarly, deletion of the first 63 nucleotides from the 5' end did not abolish preferential accumulation of GUS in BS cells, but removing 78 nucleotides did (Figure 1A-C). A fragment incorporating bases 64 to 162 was sufficient for preferential accumulation in the BS after microprojectile bombardment (Figure 1B) and stable transformation (Figure 1C; Figure S2). We conclude that one region composed of nucleotides TTGGGTGAA (64 to 79 downstream of the translational start codon) and another of GATCCTTG (141 to 162 nucleotides downstream of the translational start codon) are necessary for preferential accumulation of *GgNAD-ME1* in BS cells of C<sub>4</sub> *G. gynandra*. These two regions are separated by 75 nucleotides and will hereafter be referred to as Bundle Sheath Motif 1a (BSM1a) and Bundle Sheath Motif 1b (BSM1b).

In order to refine efforts to identify the presence of this *cis*-regulatory logic in other genes, each motif was subjected to site-directed mutagenesis and the spacer separating BSM1a and BSM1b was replaced with exogenous sequence (Figure S3). This showed that both regions are necessary but

98 that the intervening sequence is not required for preferential accumulation of GUS in BS cells (Figure  
99 S3). The fact that BSM1a and BSM1b alone are sufficient to repress transcription in M cells indicates  
100 that they do not function as promoters for long antisense mRNAs. Although the exact sequence  
101 separating BSM1a and BSM1b does not impact on their function, the distance separating them does.  
102 BSM1a and BSM1b do not generate preferential accumulation in BS cells when fused together  
103 directly, or when separated by 999 nucleotides (Figure S4). However, when the intervening  
104 sequence was between 240 and 550 nucleotides preferential accumulation in BS cells occurred  
105 (Figure S4). We conclude that in the coding region of *NAD-ME1*, two sequences separated by a  
106 spacer are necessary and sufficient to generate strong expression in BS cells.

107 Although thousands of genes are differentially expressed between M and BS cells of *C<sub>4</sub>*  
108 plants<sup>14,13,12,11</sup>, to our knowledge no DNA motifs that determine the patterning of more than one gene  
109 in BS cells have been identified. To test whether BSM1a and BSM1b operate more widely to  
110 generate preferential expression in BS cells, coding sequence of other genes relevant to *C<sub>4</sub>*  
111 photosynthesis was scanned. Sequences similar to BSM1a and BSM1b were identified in two such  
112 genes encoding mitochondrial MALATE DEHYDROGENASE (*mMDH*) and GLYCOLATE OXIDASE  
113 1 (*GOX1*). Fragments from *mMDH* and *GOX1* containing each motif were sufficient to drive BS  
114 accumulation of GUS (Figure 2A), and when either was deleted preferential accumulation in BS cells  
115 was lost (Figure 2A). The identification of BSM1a and BSM1b in these additional genes allowed  
116 consensus sequences to be defined (Figure 2B). These data suggest that multiple gene families  
117 involved in *C<sub>4</sub>* photosynthesis and photorespiration have been recruited into BS expression using a  
118 regulatory network based on these two motifs.

119 Sequences similar to BSM1a and BSM1b were identified near the predicted translational start  
120 sites of *GgNAD-ME2*, but also in the orthologues *AtNAD-ME1* and *AtNAD-ME2* from *C<sub>3</sub> Arabidopsis*  
121 *thaliana* (Figure 3A). In these additional genes BSM1a is located in the mitochondrial transit peptide  
122 and its position varies relative to the translational start site. BSM1b is located in the mature  
123 processed protein and its position appears invariant (Figure 3A). When either motif was removed  
124 preferential accumulation in BS cells was lost (Figure 3B). Thus, sequences defined by BSM1a and

BSM1b from eight genes (Figure 3C, Figure S6A) are necessary and sufficient to generate BS expression in the C<sub>4</sub> leaf. Although synonymous substitutions to the BSM1a and BSM1b sequences from *GgNAD-ME1* failed to abolish preferential expression in BS cells (Figure S5) amino acid sequence encoded by BSM1a and BSM1b vary considerably for two reasons. First, because BSM1a is found on either DNA strand (Figure S6). Second, in the eight genes studied amino acids encoded by each motif differ because codons are not in identical frames (Figure S6). Combined with the fact that both BSM1a and BSM1b are functional when present in antisense orientation, this variation in amino acid sequence supports the notion that these motifs do not act post-translationally through amino acid sequence, but rather function transcriptionally via transcription factor binding.

To provide orthogonal evidence that BSM1a and BSM1b are indeed the targets for transcription factors binding *in vivo* two additional approaches were pursued. First, the consensus motifs for BSM1a and BSM1b were used to search databases that document transcription factor binding specificities<sup>15,16</sup>. Six members of the MYOBLASTOMA (MYB) family have been reported to bind sequences that are similar but not identical to the BSM1a motif ( $p < e^{-3}$ , Supplementary Table IV) but there were no clear matches for BSM1b. It is therefore possible that a member of the MYB family binds BSM1a, but the cognate factor binding BSM1b remains unclear. Second, DNaseI sequencing of *G. gynandra* was used to define regions of chromatin that are accessible for transcription factor binding. Genome-wide DNaseI sequencing has previously been used to define the landscape of DNase I hypersensitive sites (DHS) as well as digital genomic footprints (DGF) that mark sequence motifs subject to transcription factor binding<sup>9,16,18,19</sup>. Consistent with the reporter analysis described above, DNaseI analysis showed that BSM1a and BSM1b from *GgNAD-ME1*, *GgNAD-ME2*, *GgmMDH* and *GgGOX1* were associated with DHS and so accessible to *trans*-acting factors (Figure S7). For *GgNAD-ME1*, *GgmMDH* and *GgGOX1* these DHS were present in germinating seedlings grown in the dark, whilst for *GgNAD-ME2* the DHS appeared within 4 hours of transfer from dark to light. The four DHS remained in adult leaf tissue (Figure S7). Moreover, BSM1a and BSM1b either overlapped with, or were adjacent to DGF present in these DHS (Figure S7). Whilst DNA sequence bound by a transcription factor can be found in the centre of a DGF, variation in how individual

transcription factors bind and contort DNA and thus make sequence available for DNaseI digestion result in different cleavage profiles<sup>20–23</sup>. The similarity of BSM1a to a MYB binding site, the location of both BSM1a and BSM1b in DNA accessible for transcription factor binding, and the presence of DGF in close proximity to each sequence supports the contention that these motifs function as duons and recruit transcription factors that lead to preferential expression in the BS of C<sub>4</sub> *G. gynandra*.

In contrast, analysis of publicly available DNaseI data for Arabidopsis<sup>9,20</sup> indicated that for *AtNAD-ME1*, *AtNAD-ME2* or *AtGOX1* BSM1a and BSM1b are not located within DHS (Figure S8). These data are consistent with a model in which BSM1a and BSM1b are present in the ancestral C<sub>3</sub> state, but their relative inaccessibility inhibits binding by cognate transcription factors, and so cell specific gene expression is not achieved. In contrast, in the C<sub>4</sub> leaf BSM1a and BSM1b are accessible to transcription factor binding and this leads to preferential gene expression in BS compared with M cells.

To determine the extent to which BSM1a and BSM1b could control gene expression more widely, their distribution across the genome was quantified. More than 4000 genes contained both motifs in *G. gynandra*, but fewer than half of these genes were separated by 35–550 nucleotides required to repress expression in M cells (S Table III). Compared with randomly generated hits of the same size both BSM1a and BSM1b were over-represented in coding sequence DHS (Figure S9), supporting the notion that they carry out a specialised role in transcription factor binding within gene bodies. Of the genes containing BSM1a and BSM1b separated by 35–550 nucleotides, 250 are found in genes that are preferentially expressed in BS cells of *G. gynandra*. We therefore estimate that these duons control BS expression of between four and 250 genes in the C<sub>4</sub> leaf. It was noticeable that BSM1a and BSM1b were also widespread in C<sub>3</sub> *A. thaliana* (S Table III). The simplest explanation for the presence of BSM1a and BSM1b in genes that are not preferentially expressed in BS cells may be that they contain additional *cis*-elements that over-ride their function. The highly combinatorial nature of transcription factor binding and impact that this has on cell-specific gene expression in the root has previously been documented<sup>25</sup>. It is also possible that chromatin environment also influences the function of these motifs and that this can repress their role in restricting expression to BS cells.



179 We next investigated the extent to which BSM1a and BSM1b are conserved across 1135 wild  
180 inbred *A. thaliana* accessions<sup>26</sup>. No Single Nucleotide Polymorphisms (SNP) were detected within  
181 either BSM1a or BSM1b (Figure 4A&B). Despite the lack of chromatin accessibility observed in the  
182 orthologues from *A. thaliana* (Figure S7) the BSM1a and BSM1b motifs and sequence around them  
183 are highly conserved suggesting a role for these duons in C<sub>3</sub> *A. thaliana* that is independent cell  
184 preferential expression in the leaf. To determine whether these motifs are also found more widely in  
185 *NAD-ME* genes, homologues were retrieved from all sequenced land plants (Figure S10). All  
186 dicotyledons contained at least one *NAD-ME* gene carrying sequences that define BSM1a and  
187 BSM1b (Figure 4C). In monocotyledons, BSM1a was completely conserved in rice, *Brachypodium*  
188 and *Panicum*. Although BSM1b showed one nucleotide substitution in all monocotyledonous  
189 genomes available, its conservation in spikemoss and moss argues for it being ancient (Figure 4C,  
190 S10). Both BSM1a and BSM1b are highly conserved in *GOX1* and *MDH* genes of land plants, but  
191 BSM1b appears more ancient as it is found in *GOX1* genes from all land plants and even chlorophyte  
192 algae.

193 In view of the importance of exonic sequences in regulating BS accumulation of *NAD-ME*, *MDH*  
194 and *GOX1*, the wider importance of exonic regulation in controlling C<sub>4</sub> gene expression was  
195 investigated. Seven of twelve core C<sub>4</sub> cycle genes contained regulatory elements located in  
196 transcribed sequence that were sufficient to generate preferential accumulation in M or BS cells  
197 (Figure 5). Except for *PPCk1*, which encodes the kinase that post-translationally modifies PEPC in  
198 M cells of C<sub>4</sub> plants, coding sequences from orthologues of these C<sub>4</sub> genes in C<sub>3</sub> *A. thaliana* led to  
199 preferential expression in either M or BS cells of the C<sub>4</sub> leaf (Figure 5). Overall, these data indicate  
200 that spatial patterning of gene expression in the C<sub>4</sub> leaf is largely derived from *cis*-regulatory  
201 elements present in genes found in the ancestral C<sub>3</sub> state.

202

## 203 Discussion

204 The data presented here, combined with previous reports<sup>7,27,28</sup> portray an overview of the contribution  
205 that untranslated regions (UTRs) and coding sequences make to the generation of cell-specific gene

206 expression in leaves of C<sub>4</sub> *G. gynandra*. Seven of the twelve core C<sub>4</sub> cycle genes possess regulatory  
207 elements in their transcript sequences that are sufficient for preferential accumulation in either M or  
208 BS cells of the C<sub>4</sub> leaf. These data strongly imply that, in addition to promoters being involved in  
209 generating cell-specificity in C<sub>4</sub> leaves<sup>29,30</sup>, coding sequences and UTRs play a widespread role in  
210 the preferential accumulation of C<sub>4</sub> transcripts to either M or BS cells. It remains to be seen whether  
211 this high degree of regulation from genic sequence is critical for the spatial control of gene expression  
212 in other tissues and other species. However, as genome-wide studies of transcription factor  
213 recognition sites in organisms as diverse as *A. thaliana* and human cells<sup>18,31</sup> have reported that  
214 significant binding occurs in genic sequence, we anticipate many more examples of spatial regulation  
215 of gene expression being associated with *cis*-elements outside of promoter sequences.

216 The accumulation of *GgNADME1*, *GgNADME2*, *mMDH* and *GOX1* transcripts in BS cells is  
217 dependent on the co-operative function of two *cis*-elements that are separated by a spacer  
218 sequence, all of which are located in the first exon of these genes. There are a number of options  
219 relating to how such regulatory sequence found within genes could operate. For example, they could  
220 impact on transcription by either interacting with a transcription factor or by acting as promoters to  
221 drive expression of an antisense RNA. In addition, they could act post-transcriptionally via the RNA  
222 molecule once the gene is transcribed. As the BSM1a and BSM1b motifs without additional  
223 sequence are sufficient to restrict gene expression to BS, this excludes them acting as promoters to  
224 antisense RNA molecules. Furthermore, two pieces of evidence indicate that they do not function  
225 post-transcriptionally via the transcripts generated from the endogenous gene. First, in different  
226 genes BSM1a is found on each strand of DNA so that it would not be present in all the mRNAs that  
227 accumulate preferentially in BS cells. Second, an antisense construct inserted as a transgene into  
228 *G. gynandra* still generates BS expression again indicating that the transcribed RNA is not involved.  
229 The simplest hypothesis is that BSM1a and BSM1b therefore act transcriptionally to suppress activity  
230 of the constitutive 35SCaMV promoter in M cells. The proposal that BSM1a and BSM1b act as duons  
231 does not exclude additional levels of regulation contributing to the cell specific accumulation of  
232 proteins in the C<sub>4</sub> leaf. Indeed, it is notable that these *cis*-elements alone do not completely restrict

233 accumulation of GUS to BS cells, and so post-transcriptional and/or post-translational regulation  
234 may well act to reinforce their function.

235 BSM1a and BSM1b are conserved both in their sequences, but also in the number of nucleotides  
236 that separates them. In orthologous *NAD-ME* genes from C<sub>3</sub> *A. thaliana*, which diverged from the  
237 Cleomaceae ~38 million years ago<sup>32,33</sup>, although these motifs are present, they are not sufficient to  
238 generate cell preferential expression in the C<sub>3</sub> leaf<sup>7</sup>. This finding indicates that for *NAD-ME* genes to  
239 be preferentially expressed in BS cells of C<sub>4</sub> plants, a change in the behaviour of one or more *trans*-  
240 factors was a fundamental event. At least in *G. gynandra*, evolution appears to have repeatedly  
241 made use of *cis*-elements that exist in genes of C<sub>3</sub> species that are orthologous to those recruited  
242 into C<sub>4</sub> photosynthesis<sup>7,27,28,34</sup>. The alteration in *trans*-factors such that they recognise ancestral  
243 elements in *cis* in the M or BS therefore appears to be an important and common mechanism  
244 associated with evolution of the highly complex C<sub>4</sub> system.

245 The dual role of exons in protein coding as well as the regulation of gene expression has received  
246 significant attention in vertebrates<sup>18,35–38</sup>. Although, 11% of transcription factor binding sites are  
247 located in exonic sequence in *A. thaliana*<sup>31</sup>, to our knowledge, the identification of BSM1a and  
248 BSM1b represents the first functional evidence for *cis*-elements in plant exons. The fact that these  
249 motifs are present in C<sub>3</sub> *A. thaliana*, and in fact, also found in the genomes of many land plants and  
250 some chlorophyte algae, indicates that these duons play ancient and conserved roles in  
251 photosynthetic organisms. The role of such regulatory elements within coding sequences has  
252 previously been proposed to be associated with constraints on both protein coding function and  
253 codon bias. For example, mutation to these *cis*- elements could be deleterious to both the correct  
254 function of the protein, but also to codon usage and so translational efficiency<sup>39–41</sup>. If this is the case,  
255 BSM1a and BSM1b could be highly conserved across deep phylogeny because of strong selection  
256 pressure on these elements that impact on translation, and this conservation is then co-opted to also  
257 regulate transcription during the evolution of C<sub>4</sub> photosynthesis to generate cell-specific gene  
258 expression. Establishing the role of BSM1a and BSM1b in C<sub>3</sub> plants would provide insight into the  
259 extent to which their role has altered during the transition from C<sub>3</sub> to C<sub>4</sub> photosynthesis.

260 Duons under strong selection pressure may represent a rich resource of *cis*-elements upon which  
261 the C<sub>4</sub> pathway has evolved. Although C<sub>4</sub> photosynthesis is a complex trait that requires multiple  
262 changes to gene expression, the repeated evolution of C<sub>4</sub> species across multiple plant lineages  
263 suggests that a relatively low number of changes may be required to acquire the C<sub>4</sub> syndrome<sup>42–44</sup>.  
264 A single C<sub>4</sub> master switch has been proposed<sup>44</sup> but despite multiple comparative transcriptomic  
265 studies<sup>13,45–47</sup>, there is as yet no evidence for it. The only transcription factor associated with  
266 photosynthesis in C<sub>4</sub> species is Golden-like<sup>148</sup>, which controls expression of genes associated with  
267 chlorophyll biosynthesis and the light harvesting complexes in both C<sub>3</sub> and C<sub>4</sub> species<sup>49</sup>. Given the  
268 repeated and highly convergent evolution of the C<sub>4</sub> pathway, as well as evidence that separate  
269 lineages can arrive at the C<sub>4</sub> state via different routes<sup>50</sup>, it appears more plausible that C<sub>4</sub>  
270 photosynthesis made use of a number of gene sub-networks. This is now supported by a number of  
271 findings. First, just as core photosynthesis genes encoding the light harvesting complexes and  
272 Calvin-Benson-Bassham cycle are regulated by light, the vast majority of genes that encode proteins  
273 of the C<sub>4</sub> cycle in C<sub>3</sub> *A. thaliana* are also regulated by light signalling, yet, during the evolution of C<sub>4</sub>  
274 photosynthesis there was a significant gain of responsiveness to chloroplast signalling<sup>51</sup>. Second, it  
275 has been suggested that evolution of the C<sub>4</sub> pathway is associated with the recruitment of  
276 developmental motifs into leaves that in C<sub>3</sub> species operate in roots<sup>47</sup>. Lastly, the identification of the  
277 *cis*-element MEM2<sup>28</sup>, which controls preferential expression of multiple genes in C<sub>4</sub> M cells, and now  
278 BSM1a and BSM1b in four different genes that are strongly expressed in BS cells, indicates that that  
279 C<sub>4</sub> evolution has made use of small-scale recruitment of gene sub-networks in both cell-types.

280 In summary, the data indicate regulatory elements located in transcript sequences are of central  
281 importance in patterning gene expression in the C<sub>4</sub> leaf. Expression of multiple genes in BS cells is  
282 regulated by highly conserved duons that appear to play ancient roles in photosynthetic organisms.  
283 Our data also indicate that evolution of the highly C<sub>4</sub> complex trait is built upon ancient *cis*-regulatory  
284 architecture that is common to all land plants and some green algae.

285    **Materials and Methods**

286    *G. gynandra* seed were germinated in the dark at 30°C for 24 h. For microprojectile bombardment  
287    seedlings were then transferred to Murashige and Skoog (MS) medium with 1% (w/v) sucrose and  
288    0.8% (w/v) agar (pH 5.8) and grown for a further 13 days in a growth room at 22°C and 200  $\mu\text{mol m}^{-2}$   
289     $\text{s}^{-1}$  photon flux density (PFD) with a photoperiod of 16 h light. For DNaseI-seq, germinating seeds  
290    were placed in the dark to promote hypocotyl extension for 3 days or transferred a growth cabinet  
291    maintained at 25°C $\pm$ 0.5°C, 60% relative humidity, ambient CO<sub>2</sub> and 300  $\mu\text{mol m}^{-2} \text{s}^{-1}$  PFD with  
292    photoperiod of 16 h light. For DNaseI analysis tissue was flash frozen in liquid nitrogen and stored  
293    at -80°C before processing. After bombardment or stable transformation, plant tissue was GUS  
294    stained and then chlorophyll removed in 70% (v/v) ethanol.

295 **Figure legends**

296 **Figure 1: Two regions within the coding sequence of *GgNAD-ME1* are necessary for**  
297 **preferential gene expression in the bundle sheath.** An antisense construct, as well as a deletion  
298 series from the 5' and 3' ends of *GgNAD-ME1* coding sequence were translationally fused to the  
299 *uidA* reporter under the control of the CaMV35S promoter **(A)**. Percentage of cells containing GUS  
300 in the bundle sheath (BS) after microprojectile bombardment of *G. gynandra* leaves. Bars represent  
301 the percentage of stained cells in BS cells, error bars denote the standard error. \* represents  
302 statistically significant differences with P-values <0.05 and CI = 95% determined by a one-tailed t  
303 test **(B)**. GUS in *G. gynandra* transformants containing *uidA* fused to 1-240, 1-141, 79-240 and 64-  
304 162 base pairs from the translational start site of *GgNAD-ME1* **(C)**. Scale bars, 100  $\mu$ m.

305

306 **Figure 2: BSM1a and BSM1b drive the expression of additional genes in  $C_4$  photosynthesis**  
307 **and photorespiration.** Sequences similar to BSM1a and BSM1b were identified in coding  
308 sequences of *mMDH* and *GOX1* genes of *G. gynandra*. Deleting the motifs resulted in the loss of  
309 preferential accumulation of GUS in the bundle sheath (BS) **(A)**. Consensus sequences were defined  
310 from motifs operational in *NAD-ME1*, *mMDH* and *GOX1* **(B)**. Error bars denote the standard error. \*  
311 represents statistically significant differences with P-values <0.05 and CI = 95% determined by a  
312 one-tailed t test.

313

314 **Figure 3: Functional versions of BSM1a and BSM1b are present in additional *NAD-MEs*.**  
315 BSM1a and BSM1b are found in *GgNAD-ME2* and in orthologs of *GgNAD-ME1&2* from the  $C_3$   
316 species *A. thaliana* **(A)**. Translational fusions carrying these fragments confer bundle sheath (BS)  
317 preferential expression in *G. gynandra* leaves. When BSM1a or BSM1b were removed this pattern  
318 of GUS was lost **(B)**. Consensus sequences derived from all versions of BSM1a and BSM1b tested  
319 experimentally **(C)**. Error bars denote the standard error. \* represents statistically significant  
320 differences with P-values <0.05 and CI = 95% determined by a one-tailed t test.

321

322 **Figure 4: BSM1a and BSM1b are highly conserved in land plants.** Single nucleotide  
 323 polymorphisms (SNP) in *AtNAD-ME1* (A) and *AtNAD-ME2* (B) genes from 1135 wild inbred *A.*  
 324 *thaliana* accessions. On the left, the position of BSM1a and BSM1b are highlighted by dashed blue  
 325 lines, UTRs, exons and introns are denoted by black, grey and white bars respectively on the X-axis.  
 326 To the right an expanded area representing exon 1, intron 1 and exon 2 is shown, with BSM1a and  
 327 BSM1b marked within the blue dashed lines. For both genes, no SNP were detected in either motif.  
 328 The presence of each motifs was investigated in gene sequences of *NAD-ME1*, *mMDH* and *GOX1*  
 329 retrieved from 44 species in Phytozome (v10.1)<sup>52</sup>. Each Pie-chart shows the percentage of motif  
 330 instances that were identical (green), or had 1 base pair (yellow), 2 base pair (orange) substitutions  
 331 or no similarity (white) detected.

332

333 **Figure 5: Intragenic regulatory sequences play a major role controlling C<sub>4</sub> photosynthesis**  
 334 **genes.** Coding sequences encoding core proteins of the C<sub>4</sub> pathway from *G. gynandra* together with  
 335 orthologs from *A. thaliana* were translationally fused to *uidA* and placed under control of the  
 336 CaMV35S promoter. After introduction into *G. gynandra* leaves by microprojectile bombardment  
 337 mesophyll preferential expression of *CA2*, *CA4*, *PPDK* and *PPCk1*, together with bundle sheath (BS)  
 338 preferential expression of *mMDH*, *NAD-ME1* and *NAD-ME2* were observed (A). With the exception  
 339 of *PPCk* these regulatory elements are conserved in orthogues from *A. thaliana* (B). The contribution  
 340 of intragenic sequences controlling gene regulation of the C<sub>4</sub> pathway is summarized in (C), *CA2*,  
 341 *CA4*, *PPDK* and *PPCk1* (blue) and *mMDH*, *NAD-ME1&2* (red) denote genes where intragenic  
 342 sequences control cell preferential gene expression. Error bars denote the standard error. \*  
 343 represents statistically significant differences with P-values <0.05 and CI = 95% determined by a  
 344 one-tailed t test.

345

## 346 Acknowledgements

347 I-RL was supported by CONACyT and BBSRC grant BB/L014130; SJB by the *3to4* grant from the  
 348 EU and BB/I002243 from the BBSRC; GR by a Gates Cambridge Trust PhD Fellowship; PS and

349 SRS by Advanced ERC grant 694733 Revolution to JMH; and BPW by a BBSRC PhD Studentship.  
350 We thank Chris Bournnell for help with data retrieval. I-RL designed and undertook experiments to  
351 identify BSM1a and BSM1b; I-RL and SJB identified these motifs in MDH and GOX1; SJB undertook  
352 analysis of the additional coding regions; GR, PS and SRS performed DNaseI analysis and BPW  
353 designed and undertook the antisense experiment. SS made and maintained transgenic plants. I-  
354 RL, SJB and JMH wrote the manuscript and prepared the figures.

355

#### 356 **Competing Interests**

357 The authors have no competing interests.

358



1. Anbar, A. D. *et al.* A Whiff of Oxygen before the Great Oxidation Event? *Source Sci. New Ser. Nat. Sci. Philos. Trans. R. Soc London Ser. B Sci. Rapid Commun. Mass Spectrom S. Ono ai Earth Planet Sei Lett. N. J. Beukes, Sediment Geol Econ. Geol* **317**, 1903–1906 (2007).
2. Bauwe, H., Hagemann, M. & Fernie, A. R. Photorespiration: players, partners and origin. *Trends Plant Sci.* **15**, 330–6 (2010).
3. Hatch, M. D. & Slack, C. R. Photosynthesis by sugar-cane leaves. A new carboxylation reaction and the pathway of sugar formation. *Biochem. J.* **101**, 103–11 (1966).
4. Sage, R. F., Wedin, D. a & Li, M. The Biogeography of C<sub>4</sub> Photosynthesis: Patterns and Controlling Factors. *C<sub>4</sub> Plant Biol.* 313–373 (1999). doi:10.1016/B978-012614440-6/50011-2
5. Sage, R. F. Tansley review: The evolution of C<sub>4</sub> photosynthesis. *New Phytol* **161**, 30 (2004).
6. Ray, D. K., Ramankutty, N., Mueller, N. D., West, P. C. & Foley, J. A. Recent patterns of crop yield growth and stagnation. *Nat. Commun.* **3**, 1293 (2012).
7. Brown, N. J. *et al.* Independent and parallel recruitment of preexisting mechanisms underlying C<sub>4</sub> photosynthesis. *Science* **331**, 1436–1439 (2011).
8. Stergachis, A. B. *et al.* Exonic transcription factor binding directs codon choice and affects protein evolution. *Science* **342**, 1367–72 (2013).
9. Sullivan, A. M. *et al.* Mapping and dynamics of regulatory DNA and transcription factor networks in *A. thaliana*. *Cell Rep.* **8**, 2015–2030 (2014).
10. Xu, T., Purcell, M., Zucchi, P., Helentjaris, T. & Bogorad, L. TRM1, a YY1-like suppressor of *RbcS-m3* expression in maize mesophyll cells. *Proc. Natl. Acad. Sci. U. S. A.* **98**, 2295–300 (2001).
11. John, C. R., Smith-Unna, R. D., Woodfield, H., Covshoff, S. & Hibberd, J. M. Evolutionary convergence of cell-specific gene expression in independent lineages of C<sub>4</sub> grasses. *Plant Physiol.* **165**, 62–75 (2014).
12. Chang, Y.-M. *et al.* Characterizing regulatory and functional differentiation between maize mesophyll and bundle sheath cells by transcriptomic analysis. *Plant Physiol.* **160**, 165–77 (2012).
13. Li, P. *et al.* The developmental dynamics of the maize leaf transcriptome. *Nat. Genet.* **42**, 1060–1067 (2010).
14. Aubry, S., Kelly, S., Kumpers, B. & Smith-Unna, R. Deep evolutionary comparison of gene expression identifies parallel recruitment of trans-factors in two independent origins of C<sub>4</sub> photosynthesis. *PLoS* (2014).
15. Franco-Zorrilla, J. M. *et al.* DNA-binding specificities of plant transcription factors and their potential to define target genes. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 2367–72 (2014).
16. O'Malley, R. C. *et al.* Cistrome and Epicistrome Features Shape the Regulatory DNA Landscape. *Cell* **165**, 1280–1292 (2016).
17. Song, L. & Crawford, G. E. DNase-seq: A High-Resolution Technique for Mapping Active Gene Regulatory Elements across the Genome from Mammalian Cells. doi:10.1101/pdb.prot5384
18. Stergachis, A. B. *et al.* Conservation of *trans*-acting circuitry during mammalian regulatory evolution. *Nature* **515**, 365–370 (2014).
19. Neph, S. *et al.* An expansive human regulatory lexicon encoded in transcription factor footprints. *Nature* **489**, 83 (2012).
20. Sung, M.-H., Guertin, M. J., Baek, S. & Hager, G. L. DNase footprint signatures are dictated by factor dynamics and DNA sequence. *Mol. Cell* **56**, 275–85 (2014).
21. Piper, J. *et al.* Wellington: a novel method for the accurate identification of digital genomic footprints from DNase-seq data. *Nucleic Acids Res.* **41**, e201 (2013).
22. Hesselberth, J. R. *et al.* Global mapping of protein-DNA interactions in vivo by digital genomic footprinting. *Nat. Methods* **6**, 283–289 (2009).
23. Kulakovskiy, I. V., Favorov, A. V. & Makeev, V. J. Motif discovery and motif finding from

- genome-mapped DNase footprint data. *Bioinformatics* **25**, 2318–2325 (2009).
24. Zhang, T., Marand, A. P. & Jiang, J. PlantDHS: a database for DNase I hypersensitive sites in plants. *Nucleic Acids Res.* **44**, D1148–D1153 (2016).
  25. Sparks, E. E. *et al.* Establishment of Expression in the SHORTROOT-SCARECROW Transcriptional Cascade through Opposing Activities of Both Activators and Repressors. *Dev. Cell* **39**, 585–596 (2016).
  26. Consortium, 1001 Genomes. 1,135 Genomes Reveal the Global Pattern of Polymorphism in *Arabidopsis thaliana*. *Cell* (2016).
  27. Kajala, K. *et al.* Multiple *Arabidopsis* genes primed for recruitment into C<sub>4</sub> photosynthesis. *Plant J* **69**, 47–56 (2012).
  28. Williams, B. P. *et al.* An untranslated *cis*-element regulates the accumulation of multiple C<sub>4</sub> enzymes in *Gynandropsis gynandra* mesophyll cells. *Plant Cell* (2016).
  29. Gowik, U. *et al.* *cis*-Regulatory elements for mesophyll-specific gene expression in the C<sub>4</sub> plant *Flaveria trinervia*, the promoter of the C<sub>4</sub> phosphoenolpyruvate carboxylase gene. *Plant Cell* **16**, 1077–1090 (2004).
  30. Sheen, J. C<sub>4</sub> gene expression. *Ann. Rev Plant Physiol. Plant Mol Biol* **50**, 187–217 (1999).
  31. Sullivan, A. M. *et al.* Mapping and Dynamics of Regulatory DNA and Transcription Factor Networks in *A. thaliana*. *Cell Rep.* 1–16 (2014). doi:10.1016/j.celrep.2014.08.019
  32. Schranz, M. E. & Mitchell-Olds, T. Independent ancient polyploidy events in the sister families Brassicaceae and Cleomaceae. *Plant Cell* **18**, 1152–1165 (2006).
  33. Couvreur, T. L. P. *et al.* Molecular phylogenetics, temporal diversification, and principles of evolution in the mustard family (Brassicaceae). *Mol. Biol. Evol.* **27**, 55–71 (2010).
  34. Reyna-Llorens, I. & Hibberd, J. M. Recruitment of pre-existing networks during the evolution of C<sub>4</sub> photosynthesis. *Philos. Trans. R. Soc. London B Biol. Sci.* **372**, (2017).
  35. Lang, G., Gombert, W. & Gould, H. A transcriptional regulatory element in the coding 529 sequence of the human *Bcl-2* gene. *Immunology* **114**, 25–36 (2005).
  36. Goren, A. *et al.* Comparative Analysis Identifies Exonic Splicing Regulatory Sequences—The Complex Definition of Enhancers and Silencers. *Mol. Cell* **22**, 769–781 (2006).
  37. Tümpel, S., Cambroner, F., Sims, C., Krumlauf, R. & Wiedemann, L. M. A regulatory module embedded in the coding region of *Hoxa2* controls expression in rhombomere 2. *Proc. Natl. Acad. Sci.* **105**, 20077–20082 (2008).
  38. Dong, X. *et al.* Exonic remnants of whole-genome duplication reveal *cis*-regulatory function of coding exons. *Nucleic Acids Res.* **38**, 1071–1085 (2009).
  39. Robinson, M. *et al.* Codon usage can affect efficiency of translation of genes in *Escherichia coli*. *Nucleic Acids Res.* **12**, 6663–6671 (1984).
  40. Tuller, T., Waldman, Y. Y., Kupiec, M. & Ruppin, E. Translation efficiency is determined by both codon bias and folding energy. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 3645–50 (2010).
  41. Nakahigashi, K. *et al.* Effect of codon adaptation on codon-level and gene-level translation efficiency in vivo. *BMC Genomics* **15**, 1115 (2014).
  42. Sinha, N. R. & Kellogg, E. A. Parallelism and diversity in multiple origins of C<sub>4</sub> photosynthesis in grasses. *Am. J. Bot.* **83**, 1458–1470 (1996).
  43. Hibberd, J. M., Sheehy, J. E. & Langdale, J. A. Using C<sub>4</sub> photosynthesis to increase the yield of rice - rationale and feasibility. *Curr. Opin. Plant Biol.* **11**, 228–231 (2008).
  44. Westhoff, P. & Gowik, U. Evolution of C<sub>4</sub> Photosynthesis—Looking for the Master Switch. *Plant Physiol.* **154**, 598–601 (2010).
  45. Brautigam, A. *et al.* An mRNA blueprint for C<sub>4</sub> photosynthesis derived from comparative transcriptomics of closely related C<sub>3</sub> and C<sub>4</sub> species. *Plant Physiol* (2010).
  46. Aubry, S., Kelly, S., Kümpers, B. M. C., Smith-Unna, R. D. & Hibberd, J. M. Deep Evolutionary Comparison of Gene Expression Identifies Parallel Recruitment of Trans-Factors in Two Independent Origins of C<sub>4</sub> Photosynthesis. *PLoS Genet* **10**, e1004365 (2014).
  47. Kulahoglu, C. *et al.* Comparative Transcriptome Atlases Reveal Altered Gene Expression Modules between Two Cleomaceae C<sub>3</sub> and C<sub>4</sub> Plant Species. *Plant Cell* **26**, 3243–3260

464 (2014).

465 48. Rossini, L., Cribb, L., Martin, D. J. & Langdale, J. A. The maize *GOLDEN2* gene defines a

466 novel class of transcriptional regulators in plants. *Plant Cell* **13**, 1231–1244 (2001).

467 49. Waters, M. T. *et al.* GLK transcription factors coordinate expression of the photosynthetic

468 apparatus in *Arabidopsis*. *Plant Cell* **21**, 1109–1128 (2009).

469 50. Williams, B. P., Johnston, I. G., Covshoff, S. & Hibberd, J. M. Phenotypic landscape

470 inference reveals multiple evolutionary paths to C<sub>4</sub> photosynthesis. *Elife* **2**, e00961 (2013).

471 51. Burgess, S. J. *et al.* Ancestral light and chloroplast regulation form the foundations for C<sub>4</sub>

472 gene expression. *Nat. Plants* **2**, (2016).

473 52. Goodstein, D. M. *et al.* Phytozome: a comparative platform for green plant genomics.

474 *Nucleic Acids Res.* **40**, D1178-86 (2012).

475 53. Gibson, D. G. *et al.* Enzymatic assembly of DNA molecules up to several hundred kilobases.

476 *Nat. Methods* **6**, 343–5 (2009).

477 54. Newell, C. A. *et al.* *Agrobacterium tumefaciens* mediated transformation of *Cleome*

478 *gynandra* L., a C<sub>4</sub> dicotyledon that is closely related to *Arabidopsis thaliana*. *J. Exp. Bot.* **61**,

479 1311–1319 (2010).

480 55. Schindelin, J. *et al.* Fiji: an open-source platform for biological-image analysis. *Nat. Methods*

481 **9**, 676–682 (2012).

482 56. Bailey, T. L. *et al.* MEME S UITE : tools for motif discovery and searching. *Conflict* **37**, 202–

483 208 (2009).

484 57. Gupta, S., Stamatoyannopoulos, J. A., Bailey, T. L. & Noble, W. Quantifying similarity

485 between motifs. *Genome Biol.* **8**, R24 (2007).