

1 **Genome-wide transcription factor binding in leaves from C₃ and C₄ grasses**

2

3 Steven J. Burgess^{1†} (0000-0003-2353-7794), Ivan Reyna-Llorens^{1†} (0000-0001-7964-
4 7306), Sean R. Stevenson¹ (0000-0001-5635-4340), Pallavi Singh¹ (0000-0003-3694-
5 6378), Katja Jaeger² (0000-0002-4153-7328) and Julian M. Hibberd (0000-0003-0662-
6 7958)^{1*}

7

8

9 ¹Department of Plant Sciences, Downing Street, University of Cambridge, Cambridge CB2
10 3EA, UK.

11 ²Sainsbury Laboratory, Bateman Street, University of Cambridge, Cambridge, CB2 1LR,
12 UK.

13

14 I-RL - suallorems@gmail.com

15 SJB - sjb287@illinois.edu

16 SRS - srs62@cam.ac.uk

17 PS - ps753@cam.ac.uk

18 KJ - katja.jaeger@slcu.cam.ac.uk

19 JMH (corresponding) - jmh65@cam.ac.uk

20

21 †These authors contributed equally to this work

22

23

24

25

26

27

28

29

30 **Short title:** Transcription factor binding atlas of grass leaves

31

32 **Keywords:** C₄ photosynthesis, gene regulation, evolution, cereals

33

34 **One sentence summary:** Genome-wide patterns of transcription factor binding *in vivo*
35 defined by DNaseI for leaves of C₃ and C₄ grasses

36 **Abstract**

37 The majority of plants use C₃ photosynthesis, but over sixty independent lineages of
38 angiosperms have evolved the C₄ pathway. In most C₄ species, photosynthesis gene
39 expression is compartmented between mesophyll and bundle sheath cells. We performed
40 DNaseI-SEQ to identify genome-wide profiles of transcription factor binding in leaves of
41 the C₄ grasses *Zea mays*, *Sorghum bicolor* and *Setaria italica* as well as C₃ *Brachypodium*
42 *distachyon*. In C₄ species, while bundle sheath strands and whole leaves shared similarity
43 in the broad regions of DNA accessible to transcription factors, the short sequences bound
44 varied. Transcription factor binding was prevalent in gene bodies as well as promoters,
45 and many of these sites could represent duons that impact gene regulation in addition to
46 amino acid sequence. Although globally there was little correlation between any individual
47 DNaseI footprint and cell-specific gene expression, within individual species transcription
48 factor binding to the same motifs in multiple genes provided evidence for shared
49 mechanisms governing C₄ photosynthesis gene expression. Furthermore, interspecific
50 comparisons identified a small number of highly conserved transcription factor binding
51 sites associated with leaves from species that diverged around 60 million years ago.
52 These data therefore provide insight into the architecture associated with C₄
53 photosynthesis gene expression in particular and characteristics of transcription factor
54 binding in cereal crops in general.

55 **Introduction**

56 Most photosynthetic organisms, including crops of global importance such as wheat,
57 rice and potato use the C₃ photosynthesis pathway in which Ribulose-Bisphosphate
58 Carboxylase Oxygenase (RuBisCO) catalyses the primary fixation of CO₂. However,
59 carboxylation by RuBisCO is competitively inhibited by oxygen binding the active site
60 (Bowes et al., 1971). This oxygenation reaction generates toxic waste-products that are
61 recycled by an energy-demanding series of metabolic reactions known as photorespiration
62 (Bauwe et al., 2010; Tolbert, 1971). The ratio of oxygenation to carboxylation increases
63 with temperature (Jordan and Ogren, 1984; Sharwood et al., 2016) and so losses from
64 photorespiration are particularly high in the tropics.

65 Multiple plant lineages have evolved mechanisms that suppress oxygenation by
66 concentrating CO₂ around RuBisCO. One such strategy is known as C₄ photosynthesis.
67 Species that use the C₄ pathway include maize, sorghum and sugarcane, and they
68 represent the most productive crops on the planet (Sage and Zhu, 2011). In C₄ leaves,
69 additional expenditure of ATP, alterations to leaf anatomy and cellular ultrastructure, as
70 well as spatial separation of photosynthesis between compartments (Hatch, 1987) allows
71 CO₂ concentration to be increased around tenfold compared with that in the atmosphere
72 (Furbank, 2011). Despite the complexity of C₄ photosynthesis, it is found in over 60
73 independent plant lineages (Sage et al., 2011). In most C₄ plants the initial RuBisCO-
74 independent fixation of CO₂ and the subsequent RuBisCO-dependent reactions take place
75 in distinct cell-types known as mesophyll and bundle sheath cells. Although the spatial
76 patterning of gene expression that generates these metabolic specialisations is
77 fundamental to C₄ photosynthesis very few examples of *cis*-elements or *trans*-factors that
78 restrict gene expression to mesophyll or bundle sheath cells of C₄ plants have been
79 identified (Brown et al., 2011; Gowik et al., 2004; Williams et al., 2016; Reyna-Llorens et
80 al., 2018). Moreover, in grasses more generally the DNA-binding properties of relatively
81 few transcription factors have been validated (Bolduc and Hake, 2009; Yu et al., 2015;
82 Eveland et al., 2014; Pautler et al., 2015). In summary, in both C₃ and C₄ species, work
83 has focussed on analysis of mechanisms controlling the expression of individual genes,
84 and so our understanding of the overall landscape associated with photosynthesis gene
85 expression is poor.

86 In yeast and animal systems, the high sensitivity of open chromatin to DNaseI (Zentner
87 and Henikoff, 2014) has allowed comprehensive, genome-wide characterization of
88 transcription factor binding sites at single nucleotide resolution (Hesselberth et al., 2009;
89 Neph et al., 2012; Thurman et al., 2012). In plants, DNaseI-SEQ and more recently Assay

90 for Transposase-Accessible Chromatin (ATAC-SEQ) have been employed in C₃ species
91 and provided insight into the patterns of transcription factor binding associated with
92 development (Zhang et al., 2012a; Pajoro et al., 2014; Zhang et al., 2012b, 2016), heat
93 stress (Sullivan et al., 2014) and root cell differentiation (Maher et al., 2017). By carrying
94 out DNaseI-SEQ on grass leaves that use either C₃ or C₄ photosynthesis, we aimed to
95 provide insight into the transcription factor binding repertoire associated with each form of
96 photosynthesis. Our data indicate more transcription factor binding sites are found in gene
97 bodies than promoters, and up to 25% of the footprints represent 'duons' – sequences
98 located in exons that have an impact on both gene regulation as well as the amino acid
99 sequence of the protein they encode. It is also clear that specific cell types from leaf tissue
100 make use of a markedly distinct *cis*-regulatory code and that despite significant turnover in
101 the cistrome of grasses, a small number of transcription factor motifs are conserved across
102 60 million years of evolution. Comparison of sites bound by transcription factors in both C₃
103 and C₄ leaves demonstrates that the repeated evolution of C₄ photosynthesis is built on
104 both the *de novo* gain of *cis*-elements and the exaptation of highly conserved regulatory
105 elements found in the ancestral C₃ system.

106 **Results**

107 **A *cis*-regulatory atlas for grasses**

108 To provide insight into the regulatory architecture associated with C₃ and C₄
109 photosynthesis in cereal crops, four grass were selected. *Brachypodium distachyon* uses
110 the ancestral C₃ pathway (Figure 1A). *Sorghum bicolor*, *Zea mays* and *Setaria italica* all
111 use C₄ photosynthesis, they were chosen as phylogenetic reconstructions indicate that *S.*
112 *italica* represents an independent evolutionary origin of the C₄ pathway (Figure 1A) and
113 comparison of these species can provide insight into parallel and convergent evolution of
114 C₄ gene expression. Nuclei from a minimum of duplicate samples of *S. italica* (C₄), *S.*
115 *bicolor* (C₄), *Z. mays* (C₄) and *B. distachyon* (C₃) leaves were treated with DNaseI
116 (Supplemental Figure 1) and subjected to deep sequencing. A total of 806,663,951 reads
117 could be uniquely mapped to the respective genome sequences of these species
118 (Supplemental Table 1). From all four genomes, 159,396 DNaseI-hypersensitive sites
119 (DHS) of between 150-15,060 base pairs representing broad regulatory regions accessible
120 to transcription factor binding were identified (Figure 1B). Between 20,817 and 27,746
121 genes were annotated as containing at least one DHS (Supplemental Table 2). For
122 subsequent analysis, only DHS that were consistent between replicates as determined by
123 the Irreproducible Discovery Rate framework (Li and Dewey, 2011) were used.

124 DNaseI footprinting is a well-established technique for detecting DNA-protein
125 interactions at base pair resolution and as such has been used to generate Digital
126 Genomic Footprints (DGF) to predict transcription factor binding sites. DGF are obtained
127 by pooling all replicates to maximise the number of reads that map within each DHS, and
128 then modelling differential accumulation of reads mapping to positive or negative strands
129 around transcription factor binding sites within the DHS (Piper et al., 2013). However, the
130 DNaseI enzyme possesses some sequence bias that can affect prediction of transcription
131 factor binding sites (He et al., 2014; Yardimci et al., 2014). After performing DNaseI-SEQ
132 on “naked DNA” that is devoid of nucleosomes from each species, we identified hundreds
133 of DGF that likely represent false positives (Supplemental Figure 2A). For all species,
134 analysis of the DGF derived from naked DNA showed that treatment with DNaseI led to
135 similar sequences being preferentially digested (Supplemental Figure 2B). However,
136 because false positive DGF predicted from this approach will be influenced by the number
137 of reads that map to each genome, and in the case of maize fewer reads mapped in total,
138 the number of false positive DGF varied between species (Supplemental Figure 2A). To
139 overcome this issue, we implemented a more conservative pipeline that rather than
140 defining false positives at specific locations within the genome, calculates DNaseI cutting

141 bias for all hexamers across each genome. By employing a mixture model framework,
142 these data are then used to generate a background signal to estimate footprint likelihood
143 scores for each putative DGF (Yardimici et al., 2014; Supplemental Figure 2B). This
144 approach removed between 15% and 30% of DGF from each sample (Supplemental
145 Figure 2C) and left a total of 430,205 DGF corresponding to individual transcription factor
146 binding sites between 11 and 25 base pairs being identified (Figure 1B&C; Supplemental
147 Table 3). At least one transcription factor footprint was identified in >75% of the broader
148 regions defined by DHS (Supplemental Table 2).

149 We attempted to saturate the number of predicted DGF by sequencing each species at
150 high depth (Supplemental Table 1). *In silico* sub-sampling of these data indicated that for
151 *S. bicolor*, *S. italica* and *B. distachyon*, the total number of DGF was close to saturation,
152 but for maize despite obtaining 251,955,063 reads from whole leaves this was not the
153 case (Supplementary Figure 3). Consequently, fewer DGF were predicted in maize
154 compared with the other species (Figure 1C). Since maize has a similar gene number to
155 the other species analysed, it is possible that the reduced ability to map reads to unique
156 loci was associated with the high amount of repetitive DNA in the maize genome. Another
157 contributing factor to the poor mapping rate in maize may be the low complexity found in
158 one of the libraries as reflected by the PCR bottleneck coefficient (Supplemental Table 1).
159 According to the Encyclopaedia of DNA elements (ENCODE), large number of reads from
160 low complexity libraries decreases the chances of identifying the majority of transcription
161 factor binding sites. However, despite these differences in coverage and in certain quality
162 metrics, for all four species DHS and DGF were primarily located in gene-rich regions and
163 depleted around centromeres (Figure 1D). Individual transcription factor binding
164 sequences were resolved in all chromosomes from each species (Figure 1D). On a
165 genome-wide basis, the distribution of DHS was similar between species, with the highest
166 proportion of such sites located in promoter, coding sequence and intergenic regions
167 (Figure 1E). Notably, in all four grasses, genic sequences contained more DHS than
168 promoters (Figure 1F).

169 To provide additional evidence confirming that DGF identified in our analysis derive
170 from protein-DNA interactions, they were compared with previously identified motifs from
171 maize. Maize is the most appropriate choice for this analysis as there are more data on
172 transcription factor binding sites than in *S. bicolor* and *S. italica*. Moreover, support from
173 previous work goes some way to supporting the smaller number of DGF that we identified
174 in this species. Therefore, the literature was assessed for validated transcription factor
175 binding sites in maize. These have previously been associated with flowering (Kozaki et

176 al., 2004; Vollbrecht et al., 2005; Eveland et al., 2014), meristem development (Bolduc et
177 al., 2012), Gibberellin catabolism (Bolduc and Hake, 2009), sugar signalling (Niu et al.,
178 2002) and leaf development of maize (Yu et al., 2015), but in all cases, DGF matching
179 these motifs were found in our dataset (Figure 2A, FDR < 0.001). In addition, a larger
180 ChIP-SEQ dataset of 117 transcription factors from maize leaves obtained from the pre-
181 release maize cistrome (<http://www.epigenome.cuhk.edu.hk/C3C4.html>, Supplemental
182 Figure 4, Supplemental Table 4) was compared with our data. Differences between
183 specific binding sites are likely because in all cases growth conditions will have varied from
184 ours, and in some cases different tissues were sampled. Despite this, 66% and 29% of the
185 ChIP-SEQ peaks overlapped with our DHS and DGF respectively. Although only 29% of
186 DGF overlapped with motifs defined by ChIP-SEQ, permutation tests performed using the
187 “regioner” package (Gel et al., 2015) indicated a statistically greater overlap than would
188 be expected by chance (pvalue = 0.0099, 100 permutations). Moreover, when both
189 features were systematically shifted from their original position, the local z-score, which
190 represents the strength of the association at any particular position showed a sharper
191 decrease for DGF than DHS suggesting the association between ChIP-SEQ peaks and
192 DGFs is more strongly linked to the exact position of the DGF (Supplemental Figure 4B).
193 In summary, despite detecting fewer DGF in maize than in the other species, the DGF we
194 found are supported by publicly available ChIP-SEQ, EMSA and Selex datasets.

195 Consistent with the distribution of DHS (Figure 1E), annotated DGF were most common
196 in promoter, coding sequence and intergenic regions (Figure 2B) and genic sequences
197 contained more DGF than promoters (Figure 2C). Distribution plots showed that the
198 highest density of DGF was close to the annotated transcription start sites but indicated a
199 slightly skewed distribution favouring genic sequence including exons (Supplemental
200 Figure 5). A similar pattern was observed for the ChIP-SEQ signal peaks (Supplemental
201 Figure 4C). Transcription factor binding sites located in exons have been termed duons
202 because they could impact both on the regulation of transcription and amino acid
203 sequence. Whilst in general synonymous mutations not affecting amino acid sequence
204 should be under relaxed purifying selection, because of transcription factor recognition all
205 nucleotides in duons should be under purifying selection, and thus show lower mutation
206 rates. We therefore investigated the nucleotide substitution rate at Four-Fold Degenerate
207 Sites (FFDS) using variation data from 1218 maize lines (Bukowski et al., 2018) and found
208 that it was statistically significantly lower in duons than in surrounding coding sequence
209 (Figure 2E, $p = 7.04e-9$). This contrasts with the density of polymorphisms in non-
210 synonymous sites (Figure 2D). Although it has been proposed that GC bias of duons

211 constrains FFDS (Xing and He, 2014) we found no such bias between duon and exon
212 sequences used in this analysis (Supplemental Figure 6). Taken together, we conclude
213 that in these cereals a significant proportion of transcription factor binding likely takes
214 place within genes.

215

216 **A distinct *cis*-regulatory lexicon for specific cells within the leaf**

217 The above analysis provides a genome-wide overview of the *cis*-regulatory architecture
218 associated with leaves of grasses. However, as with other complex multicellular systems,
219 leaves are composed of many specialised cell types. Because DGF are defined by the
220 differential DNA cleavage between protected and unprotected regions of DNA within a
221 DHS, a negative distribution compared with the larger DHS is produced (Figure 3A). Thus,
222 transcription factor binding signal from a low abundance cell-type is likely to be obscured
223 by overall signal from a tissue-level analysis (Figure 3A). Since bundle sheath strands can
224 be separated (Covshoff et al., 2013; Leegood, 1985; Furbank et al., 1985) C₄ species
225 provide a simple system to study transcription factor binding in specific cells of leaves
226 (Figure 3B). After bundle sheath isolation from *S. bicolor*, *S. italica* and *Z. mays*, and
227 naked DNA correction for inherent bias in DNaseI cutting, a total of 129,137 DHS were
228 identified (Figure 3B; Supplemental Table 5) containing 244,554 DGF (Figure 3B;
229 Supplemental Table 5; FDR<0.01). Of these, 138,075 were statistically enriched in the
230 bundle sheath samples compared with whole leaves (Figure 3B; Supplemental Table 5).
231 The number of these statistically enriched DGF in bundle sheath strands of C₄ species
232 was large and ranged from 14,250 to 73,057 in maize and *S. italica* respectively
233 (Supplemental Table 5). The lower number in maize is likely due to the reduced
234 sequencing depth achieved. Genome-wide, the number of broad regulatory regions
235 defined by DHS in the bundle sheath that overlapped with those present in whole leaves
236 ranged from 71 to 84% in *S. italica* and *S. bicolor* respectively (Supplemental Table 6).
237 However, only 6-20% of the narrower DGF found in the bundle sheath were also identified
238 in whole leaves (Supplemental Table 7). Taken together, these findings indicate that
239 specific cell types of cereal leaves share similarity in the broad regions of DNA that are
240 accessible to transcription factors (DHS), but that the short sequences actually bound by
241 transcription factors (DGF) vary dramatically.

242 To provide evidence that DGF predicted after analysis of separated bundle sheath
243 strands are of functional importance, they were compared with previously validated
244 sequences. In C₄ grasses, to our knowledge there are no such examples in *S. bicolor* or *S.*
245 *italica*, but in the *RbcS* gene from maize, which is preferentially expressed in bundle

246 sheath cells, an I-box (GATAAG) is essential for light-mediated activation (Giuliano et al.,
247 1988) and a HOMO motif (CCTTTTTTCCTT) is important in driving bundle sheath
248 expression (Xu et al., 2001) (Figure 3C). Despite not reaching saturation in DGF prediction
249 in maize (Supplemental Figure 3) both elements were detected in our pipeline.
250 Interestingly, a signal suggesting TF binding to the HOMO motif was enriched in the
251 bundle sheath strands (Figure 3C), and whilst the I-box was detected in both bundle
252 sheath strands and whole leaves its position was slightly different in each cell type (Figure
253 3C). These findings are therefore consistent with the biochemical data implicating the I-box
254 in control of abundance and the HOMO box in control of cell-specific accumulation of
255 *RbcS* transcripts.

256 The *ZmPEPC* gene (GRMZM2G083841) encodes the phosphoenolpyruvate
257 carboxylase responsible for producing C₄ acids used in the C₄ pathway and is
258 preferentially expressed in mesophyll cells. Previous reports showed that a region of 600
259 nucleotides upstream of the transcription start site carrying repeated C-rich sequences
260 was sufficient to drive expression in mesophyll cells of maize (Shaffner & Sheen, 1992;
261 Matzuoka et al., 1994). Although no DGF were detected with these C-rich sequences, they
262 are located within a DHS indicating that they are available for transcription factor binding
263 (Figure 3D). Thus, despite the fact that we had not reached saturation of DGF in maize, for
264 both *RbcS* and *PEPC* the regions of DNA accessible to transcription factor binding are
265 consistent with previous reports, and in the case of *RbcS* DGF were detected that coincide
266 with known *cis*-elements.

267 To investigate the relationship between cell specific gene expression and the position of
268 DHS and DGF, the DNaseI data were interrogated using RNA-SEQ datasets from
269 mesophyll and bundle sheath cells of C₄ leaves (Chang et al., 2012; John et al., 2014;
270 Emms et al., 2016). At least three mechanisms associated with cell specific gene
271 expression operating around individual genes were identified and can be exemplified using
272 three co-linear genes found on chromosome seven of *S. bicolor*. First, in the *NADP-malate*
273 *dehydrogenase (MDH)* gene, which is highly expressed in mesophyll cells and encodes a
274 protein of the core C₄ cycle (Figure 3E) a broad DHS site and two DGF were present in
275 whole leaves, but not in bundle sheath strands (Figure 3F). Whilst presence of this site
276 indicates accessibility of DNA to transcription factors that could activate expression in
277 mesophyll cells, global analysis of all genes strongly and preferentially expressed in
278 bundle sheath strands versus whole leaves indicates that presence/absence of a DHS in
279 one cell type is not sufficient to generate cell specificity (Supplemental Figure 7 & 8).
280 Second, in the next contiguous gene that encodes an additional isoform of MDH also

281 preferentially expressed in mesophyll cells (Figure 3E), a DHS was found in both whole
282 leaf and bundle sheath strands but DGF within this region differed between cell types
283 (Figure 3F). Thus, despite similarity in DNA accessibility, the binding of particular
284 transcription factors varied between cell types. However, once again, genome-wide
285 analysis indicated that alterations to individual DGF were not sufficient to explain cell
286 specific gene expression. For example, only 30 to 40% of all enriched DGF in the bundle
287 sheath were associated with differentially expressed genes (Supplemental Table 8).
288 Lastly, in the third gene in this region, which encodes a NAC domain transcription factor
289 preferentially expressed in bundle sheath strands (Figure 3E), differentially enriched DGF
290 were associated both with regions of the gene that have similar DHS in each cell type, but
291 also a region lacking a DHS in whole leaves compared with bundle sheath strands (Figure
292 3E). These three classes of alteration to transcription factor accessibility and binding were
293 detectable in genes encoding core components of the C₄ cycle in all three species
294 (Supplemental Figure 9-11). Overall, we conclude that differences in transcription factor
295 binding between cells of C₄ leaves is associated with both DNA accessibility defined by
296 broad DHS, as well as fine-scale alterations to transcription factor binding defined by DGF.
297 Moreover, bundle sheath strands possessed a distinct regulatory landscape compared
298 with the whole leaf, and in genes encoding enzymes of the C₄ pathway multiple
299 transcription factor binding sites differed between bundle sheath and whole leaf samples.
300 This finding implies that cell specific gene expression in C₄ leaves is mediated by
301 combinatorial effects derived from alterations to gene accessibility as defined by DHS as
302 well as changes to binding of multiple transcription within these regions.

303

304 **DNA motifs associated with cell specific expression**

305 To provide an overview of the transcription factors most likely associated with DGF
306 ChIP-SEQ data from maize (Figure 2B) together with motifs from JASPAR plants (Khan et
307 al., 2018) and an additional 529 transcription factor motifs validated in Arabidopsis
308 (O'Malley et al., 2016) were used to annotate the DGF from *Z. mays*, *S. bicolor*, *S. italica*
309 and *B. distachyon* (Figure 4A). To increase the number of annotated DGF *de novo*
310 prediction was used to identify sequences over-represented in DGF compared with those
311 across the whole genome. This resulted in an additional 524 motifs being annotated
312 (Figure 4A), but in fact all of these were previously detected after *de novo* prediction from
313 DNaseI-SEQ of rice (Zhang et al., 2012b). As would be expected from *bona fide*
314 transcription factor binding, inspection of these motifs predicted *de novo* demonstrated
315 clear strand bias in DNaseI cuts (Figure 4B). By combining previously known motifs and

316 those predicted *de novo*, the percentage of DGF that could be annotated in each species
317 increased from about 60% to more than 75% (Figure 4C, Supplemental Table 9).

318 To define the most common sequences bound by transcription factors in mature leaves
319 undertaking C₃ and C₄ photosynthesis and to investigate whether C₄ photosynthesis is
320 controlled by an increase in binding of sets of transcription factors, individual motifs were
321 ranked by frequency and the Kendall rank correlation coefficient used to compare species
322 (Figure 4D). In both C₃ and C₄ species, the most prevalent transcription factor binding
323 motifs were associated with AP2-EREBP and MYB transcription factor families (p-value <
324 2.2⁻¹⁶; Figure 4D). Next, to identify regulatory factors associated with gene expression in
325 the C₄ bundle sheath, transcription factor motifs located in DGF enriched in either the
326 bundle sheath or in whole leaf samples of *S. bicolor* were identified (Figure 4E). There was
327 little difference in the ranking of the most commonly used motifs between these cell types
328 (Kendall's tau=0.815; p-value < 2.2⁻¹⁶), indicating cell-specificity is not associated with
329 large-scale changes in the abundance of many transcription factor families (Figure 4E).
330 After performing hypergeometric tests for enrichment of individual motifs in differentially
331 occupied DGF we found 133 and 106 motifs enriched in whole leaves and bundle sheath
332 strands respectively (p < 0.001). Of these 239 motifs, 37 were enriched in all C₄ species
333 with 10 and 27 enriched in the bundle sheath and whole leaf respectively (Figure 4F,
334 Supplemental Table 10) 66 were only enriched in bundle sheath strands and 91 to whole
335 leaf tissue (Supplemental Table 11). Some of these conserved and cell specific motifs
336 have been previously described to have a relevant role in photosynthesis. For instance, in
337 whole leaves of maize and *Setaria*, we found significant enrichment of the bHLH129 motif
338 (Supplemental Table 11) that has been proposed to act as a negative regulator of NADP-
339 ME (Borba et al., 2018).

340

341 **Multiple genes encoding enzymes of the C₄ and Calvin-Benson-Bassham cycles** 342 **share the same occupied *cis*-elements**

343 To investigate whether genes involved in the C₄ phenotype are co-regulated, we
344 compared the number of instances where the same motifs were bound in multiple C₄,
345 Calvin-Benson-Bassham and C₂ cycle genes (Supplemental Table 12). While no single
346 *cis*-element was found in all genes that are preferentially expressed in mesophyll or
347 bundle-sheath cells, the number of genes possessing the same occupied motif ranged
348 from nine in *S. bicolor* and *S. italica* to four in *S. bicolor* and *Z. mays* whole leaves
349 respectively (Supplemental Table 9 & 12). These data support a model where the

350 combinatorial action of multiple transcription factors controls groups of C₄ genes to
351 produce the gene expression patterns required for C₄ photosynthesis.

352 We next performed comparative analysis of motifs bound by transcription factors to
353 determine whether the set of *cis*-elements found in C₄ genes of each species were
354 common, or whether C₄ genes are regulated differently in each species. In pairwise
355 comparisons, DGF fell into three categories: conserved and occupied by a transcription
356 factor, conserved but only occupied in one species, and not conserved (Figure 5A). Only a
357 small percentage of DGF were both conserved in sequence and bound by transcription
358 factors (Figure 5B, Supplemental Table 13). Consistent with this, the majority of C₄ gene
359 orthologs did not share DGF. Due to the lack of DGF saturation in maize, these estimates
360 likely set lower bounds for the extent of conservation. However, in several cases,
361 patterning of C₄ gene expression correlated with a set of motifs shared across species
362 (Figure 5C). In some cases, these shared *cis*-elements were present in the ancestral C₃
363 state. For instance, the *TRANSKETOLASE* (*TKL*) gene contains several conserved DGF
364 that are present in the bundle sheath of the C₄ species but also in whole leaves of C₃
365 *Brachypodium* (Figure 5). This finding is consistent with the notion that C₄ photosynthesis
366 makes use of existing regulatory architecture found in C₃ plants. Nevertheless, overall,
367 these data also indicate that the majority of C₄ gene expression appears to be associated
368 with species-specific regulatory networks.

369

370 **Hyper-conserved *cis*-regulators of C₄ genes**

371 To investigate the extent to which transcription factor binding sites associated with C₄
372 genes within a C₄ lineage are conserved, genes encoding the core C₄ cycle were
373 compared in *S. bicolor* and *Z. mays* (Figure 6A, Supplemental Table 14). 27 genes
374 associated with the C₄ and Calvin-Benson-Bassham Cycles contained a total of 379 DGF.
375 Although many of these transcription factor footprints were conserved in sequence within
376 orthologous genes, only nine were both conserved and bound by a transcription factor
377 (Figure 6A). Again, due to the lack of DGF saturation in maize, these data likely represent
378 minimum estimates of conservation.

379 Genome-wide, the number of DGF that were conserved in sequence and bound by a
380 transcription factor decayed in a non-linear manner with phylogenetic distance (Figure 6B,
381 Supplemental Table 15). For example, *Z. mays* and *S. bicolor* shared 5,775 DGF that were
382 both conserved and occupied. *S. italica* shared only 670 DGF with *Z. mays* and *S. bicolor*
383 (Figure 6B). Finally, comparison of these C₄ grasses with C₃ *B. distachyon* yielded 93 DGF
384 that have been conserved over >60Myr of evolution. Because nuclei from *B. distachyon*

385 were sampled later in the photoperiod than those from the C₄ grasses, and DGF may well
386 vary over the diel cycle, it is possible that this is an underestimate of DGF conservation.
387 However, 41 of these highly conserved DGF were present in whole leaf samples of the C₃
388 species, but in the C₄ species were restricted to the bundle sheath (Figure 6B). Gene
389 Ontology analysis did not detect enrichment of any specific terms for hyper-conserved
390 DGF associated with the bundle sheath, but for whole leaves detected over-representation
391 of “cell component” categories such as membrane bound organelles and the nucleus
392 (Supplemental Table 16). In whole leaves, this set of ancient and highly conserved DGF
393 were located predominantly in 5' UTRs and coding sequences, but in bundle sheath
394 strands over fifty percent of these hyper-conserved DGF were in coding sequences
395 (Figure 6B). Overall, these data indicate that certain duons are highly conserved across
396 deep evolutionary time. The frequent use of hyper-conserved duons in the bundle sheath
397 implies that this cell type uses an ancient and highly conserved regulatory code.

398 Discussion

399 Genome-wide transcription factor binding in grasses

400 The dataset provides insight into the regulation of gene expression in cereals in general,
401 and to C₄ photosynthesis in particular. Consistent with previous analysis ranging from *A.*
402 *thaliana* (Sullivan et al., 2014) to metazoans (Natarajan et al., 2012; Stergachis et al.,
403 2013, 2014), the majority of DGF detected in the four grasses were centred around
404 annotated transcription start sites. However, in these cereals it is noteworthy that
405 transcription factor binding was prevalent in genic sequence. Whilst we cannot rule out the
406 possibility that this distribution is in some way related to the methodology used in this
407 study, there is evidence that the exact distribution of transcription factor binding appears to
408 be species specific. For example, whilst in *A. thaliana* DNaseI-SEQ revealed enrichment of
409 DHS in sequence ~400 base pairs upstream of transcription start sites as well as 5' UTRs
410 (Sullivan et al., 2014) and ATAC-SEQ of *A. thaliana*, *Medicago truncatula* and *Oryza*
411 *sativa* detected most transposase hypersensitive sites upstream of genes, in *Solanum*
412 *lycopersicum* more were present in introns and exons than upstream of annotated
413 transcription start sites (Maher et al., 2017).

414 The prevalence of transcription factor binding to coding sequences is relevant to
415 approaches used to generate transgenic plants and test gene function and regulation.
416 First, consistent with the prevalence of DGF downstream of the annotated transcription
417 start sites that we detected, it is noteworthy that during cereal transformation, exon and
418 intron sequences are frequently used to achieve stable expression of transgenes (Maas et
419 al., 1991; Cornejo et al., 1993; Jeon et al., 2000). It is possible that this strategy is required
420 in grasses because of the high proportion of transcription factor binding downstream of
421 annotated transcription start sites. These transcription factor binding sites in coding
422 sequence also have implications for synthetic biology. Although technologies such as type-
423 IIS restriction endonuclease cloning methods allow high-throughput testing of many
424 transgenes, they rely on sequence domestication. Whilst routinely this would maintain
425 amino acid sequence, without analysis of transcription factor binding sites it could mutate
426 motifs bound by transcription factors and lead to unintended modifications to gene
427 expression.

428

429 The transcription factor landscape underpinning gene expression in specific tissues

430 The finding that so few transcription factor binding sites were shared between bundle
431 sheath tissue and whole leaves of *S. bicolor*, *Z. mays* and *S. italica* argues for the need to
432 isolate these cells when attempting to understand the control of gene expression. Although

433 separating bundle sheath strands from C₄ leaves is relatively trivial (Covshoff et al., 2013;
434 Furbank et al., 1985; Leegood, 1985) this is not the case for C₃ leaves. Approaches in
435 which nuclei from specific cell-types are labelled with an exogenous tag (Deal and
436 Henikoff, 2011) now allow their transcription factor landscapes to be defined. The
437 application of DNaseI-SEQ to specific cell types has recently been used in roots (Maher et
438 al., 2017) and so in the future, this approach of both C₃ and C₄ leaves should provide
439 insight into how the extent to which gene regulatory networks have been re-wired during
440 the evolution of the complex C₄ trait.

441 Given the central importance of cellular compartmentation to C₄ photosynthesis, there
442 have been significant efforts to identify *cis*-elements that restrict gene expression to either
443 mesophyll or bundle sheath cells of C₄ leaves (Hibberd and Covshoff, 2010; Sheen, 1999;
444 Wang et al., 2014). As previous studies of C₄ gene regulation have focused on individual
445 genes and have been performed in various species, it has not been possible to obtain a
446 coherent picture of regulation of the C₄ pathway and along with many other systems, initial
447 analysis focussed on regulatory elements located in promoters of C₄ genes (Sheen, 1999).
448 However, it has become increasingly apparent that the patterning of gene expression
449 between cells in the C₄ leaf can be mediated by elements in various parts of a gene. In
450 addition to promoter elements (Sheen, 1999; Gowik et al., 2004), this includes
451 untranslated regions (Kajala et al., 2011; Patel et al., 2004; Viret et al., 1994; Williams et
452 al., 2016; Xu et al., 2001) and coding sequences (Brown et al., 2011; Reyna-Llorens et al.,
453 2018). By providing data on *in vivo* transcription factor occupancy for the complete C₄
454 pathway in three C₄ grasses, the data presented here allow broad comparisons and
455 provide several insights into regulatory networks controlling C₄ genes.

456 The DNaseI dataset indicates that cell specific gene expression in C₄ leaves is not
457 strongly correlated with changes to large-scale accessibility of DNA as defined by DHS.
458 This implies that modifications to transcription factor accessibility around any one gene
459 does not impact on its expression between tissues in the leaf. Rather, as only 8-24% of
460 transcription factor binding sites detected in the bundle sheath were also found in whole
461 leaves, the data strongly implicate complex modifications to patterns of transcription factor
462 binding in controlling gene expression between cell types. These findings are consistent
463 with analogous analysis in roots where genes with clear spatial patterns of expression are
464 bound by multiple transcription factors (Sparks et al., 2017) and highly combinatorial
465 interactions between multiple activators and repressors tune the output (de Lucas et al.,
466 2016).

467 The data also provide insight into *cis*-elements that underpin the C₄ phenotype. No
468 single *cis*-element was found in all genes preferentially expressed in either mesophyll or
469 bundle sheath cells of one species. This finding is consistent with analysis of yeast where
470 the output of genetic circuits can be maintained despite rapid turnover of *cis*-regulatory
471 mechanisms underpinning them (Tsong et al., 2006). However, we did detect small
472 numbers of C₄ genes that shared common transcription factor footprints (Figure 5,
473 Supplemental Table 14 & 15), which is consistent with previous analysis that identified
474 shared *cis*-elements in *PPDK* and *CA*, or *NAD-ME1* and *NAD-ME2* in C₄ *Gynandropsis*
475 *gynandra* (Williams et al., 2016; Reyna-Llorens et al., 2018). Interspecific comparisons
476 further underlined the high rate of divergence in the *cis*-regulatory logic used to control C₄
477 genes. For example, although we detected highly similar transcription factor footprints in
478 the *OMT1* and *TKL* genes of the three C₄ species we assessed, this was not apparent for
479 any other C₄ genes. As a result of the apparent rapid rate of evolution in *cis*-regulatory
480 architecture in these C₄ species, attempts to engineer C₄ photosynthesis into C₃ crops to
481 increase yield (Hibberd et al., 2008) may benefit from using pre-existing regulatory
482 mechanisms controlling mesophyll or bundle sheath expression in ancestral C₃ species.

483

484 **Characteristics of the transcription factor binding in the ancestral C₃ state that have** 485 **impacted on evolution of the C₄ pathway**

486 Comparison of transcription factor binding in the C₃ grass *B. distachyon* with three C₄
487 species provides insight into mechanisms associated with the evolution of C₄
488 photosynthesis. For all four grasses, irrespective of whether they used C₃ or C₄
489 photosynthesis, the most abundant DNA motifs bound by transcription factors were similar.
490 Thus, motifs recognised by the AP2-EREBP and MYB classes of transcription factor were
491 most commonly bound across each genome. This indicates that during the evolution of C₄
492 photosynthesis, there has been relatively little alteration to the most abundant classes of
493 transcription factors that bind DNA.

494 The repeated evolution of the C₄ pathway has frequently been associated with
495 convergent evolution (Sage, 2004; Sage et al., 2012). However, parallel alterations to
496 amino acid and nucleotide sequence that allow altered kinetics of the C₄ enzymes (Christin
497 et al., 2014, 2007) and patterning of C₄ gene expression (Brown et al., 2011) respectively
498 have also been reported. The genome-wide analysis of transcription factor binding
499 reported here indicates that only a small proportion of the C₄ cistrome is associated with
500 parallel evolution. These estimates regarding conservation between C₄ and C₃ species
501 may represent underestimates as whilst nuclei were all sampled in the light, those from

502 *B. distachyon* were sampled later in the photoperiod. Moreover, when orthologous genes
503 were compared between the four grasses assessed here, the majority of transcription
504 factor binding sites were not conserved, and of the DGF that were conserved, position
505 within orthologous genes varied. This indicates that C₄ photosynthesis in grasses is
506 tolerant to a rapid turnover of the *cis*-code, and that when motifs are conserved in
507 sequence, their position and frequency within a gene can vary. It therefore appears that
508 the cell-specific accumulation patterns of C₄ proteins can be maintained despite
509 considerable modifications to the cistrome of C₄ leaves. It was also the case that some
510 conserved motifs bound by transcription factors in the C₄ species were present in *B.*
511 *distachyon*, which uses the ancestral C₃ pathway. Previous work has shown that *cis*-
512 elements used in C₄ photosynthesis can be found in gene orthologs from C₃ species
513 (Williams et al., 2016; Reyna-Llorens et al., 2018). However, these previous studies
514 identified *cis*-elements that were conserved in both sequence and position. As it is now
515 clear that such conserved motifs are mobile within a gene, it seems likely that many more
516 examples of ancient *cis*-elements important in C₄ photosynthesis will be found in C₃ plants.

517 Although we were able to detect a small number of transcription factor binding sites that
518 were conserved and occupied in all four species sampled, these ancient hyper-conserved
519 motifs appear to have played a role in the evolution of C₄ photosynthesis. Interestingly, a
520 large proportion of these motifs bound by transcription factors were found in coding
521 sequences, and this bias was particularly noticeable in bundle sheath cells. Due to the
522 amino acid code, the rate of mutation of coding sequence compared with the genome is
523 restricted. If such regions have a longer half-life than transcription factor binding sites in
524 other regions of the genome, then they may represent an excellent source of raw material
525 for the repeated evolution of complex traits (Martin and Orgogozo, 2013). Our data
526 documenting the frequent use of hyper-conserved DGF in the C₄ bundle sheath implies
527 that this tissue may use an ancient and highly conserved regulatory code. It appears that
528 during the evolution of the C₄ pathway, which relies on heavy use of the bundle sheath,
529 this ancient code has been co-opted to control photosynthesis gene expression.

530 In summary, the data provide a transcription factor binding atlas for leaves of grasses
531 using either C₃ or C₄ photosynthesis. Whilst we did not achieve DGF saturation in maize,
532 commonalities between the four species were apparent. Sequences bound by transcription
533 factors were found within genes as well as promoter regions, and many of these motifs
534 represent duons. In terms of the regulation of tissue specific gene expression, whilst
535 bundle sheath strands and whole C₄ leaves shared considerable similarity in regions of
536 DNA accessible to transcription factors, the short sequences actually bound by

537 transcription factors varied dramatically. We identified a small number of transcription
538 factor motifs that were conserved in these species. The data also provide insight into the
539 regulatory architecture associated with C₄ photosynthesis more specifically. Whilst we
540 found some evidence that multiple genes important for C₄ photosynthesis share common
541 *cis*-elements bound by transcription factors, this was not widespread. This may well relate
542 to the relatively rapid turnover in the *cis*-code, and so it is possible that transcription factors
543 interacting with these motifs are more conserved. Analysis of transcription factor footprints
544 in specific cell types from leaves of C₃ grasses should in the future provide insight into the
545 extent to which gene regulatory networks have altered during the transition from C₃ to C₄
546 photosynthesis.

547 **Methods**

548 **Growth conditions and isolation of nuclei**

549 *S. bicolor*, *S. italica* and *Z. mays* were grown under controlled conditions at the Plant
550 Growth Facilities of the Department of Plant Sciences at the University of Cambridge in a
551 chamber set to 12 h/12 h light/dark; 28 °C light/20 °C dark; 400 $\mu\text{mol m}^{-2} \text{s}^{-1}$ photon flux
552 density, 60% humidity. For germination, *S. bicolor* and *Z. mays* seeds were imbibed in
553 H₂O for 48 h, *S. italica* seeds were incubated on wet filter paper at 30 °C overnight in the
554 dark. *Z. mays*, *S. bicolor* and *S. italica* were grown on 3:1 (v/v) M3 compost to medium
555 vermiculite mixture, with a thin covering of soil. Seedlings were hand-watered.

556 *B. distachyon* plants were grown in a separate growth facility under controlled
557 conditions optimised for its growth at the Sainsbury Laboratory Cambridge University, first
558 under short day conditions 14 h/10 h, light/dark for 2 weeks and then shifted to long day 20
559 h/4 h, light/dark, for 1 week and harvested at ZT20. Temperature was set at 20 °C,
560 humidity 65% and light intensity 350 $\mu\text{mol m}^{-2} \text{s}^{-1}$. All tissue was harvested from August to
561 October 2015.

562 To isolate nuclei from *S. bicolor*, *Z. mays* and *S. italica* mature third and fourth leaves
563 with a fully developed ligule were harvested 4-6 h into the light cycle 18 days after
564 germination. Bundle sheath cells were mechanically isolated as described previously
565 (Markelz et al., 2003). At least 3 g of tissue was used for each extraction. Nuclei were
566 isolated using a sucrose gradient adapted and yield quantified using a haemocytometer.
567 For *B. distachyon* plants were flash frozen and material pulverised in a coffee grinder. 3 g
568 of plant material was added to 45 ml NIB buffer (10mM Tris-HCl, 0.2M sucrose, 0.01%
569 (v/v) Triton X-100, pH 5.3 containing protease inhibitors (Sigma-Aldrich)) and incubated at
570 4°C on a rotating wheel for 5 min, afterwards debris was removed by sieving through 2
571 layers of Miracloth (Millipore) into pre-cooled flasks. Nuclei were spun down 4,000 rpm, 4
572 °C for 20 min. Plastids were lysed by adding Triton to a final concentration of 0.3% (v/v)
573 and incubated for 15 min on ice. Nuclei were pelleted by centrifugation at 5000 rpm at 4 °C
574 for 15 min. Pellets were washed 3 times with chilled NIB buffer.

575

576 **Deproteinized DNA extraction**

577 For isolation of deproteinated DNA from *S. bicolor*, *Z. mays*, *B. distachyon* and *S. italica*
578 mature third and fourth leaves with a fully developed ligule were harvested 4 h into the
579 light cycle, 18 days after germination. 100 mg of tissue was used for each extraction.
580 Deproteinated DNA was extracted using a QIAGEN DNeasy Plant Mini Kit (QIAGEN, UK)
581 according to the manufacturer's instructions.

582

583 **DNaseI digestion, sequencing and library preparation**

584 To obtain sufficient DNA each biological replicate consisted of leaves from tens of
585 individuals and to conform to standards set by the Human Genome project at least two
586 biological replicates were sequenced for each sample. 2×10^8 of freshly extracted nuclei
587 were re-suspended at 4 °C in digestion buffer (15 mM Tris-HCl, 90 mM NaCl, 60 mM KCl,
588 6 mM CaCl₂, 0.5 mM Spermidine, 1 mM EDTA and 0.5 mM EGTA, pH 8.0). DNaseI
589 (Fermentas) at 7.5 U was added to each tube and incubated at 37 °C for 3 min. Digestion
590 was arrested with addition of 1:1 volume of stop buffer (50 mM Tris-HCl, 100 mM NaCl,
591 0.1% (w/v) SDS, 100 mM EDTA, pH 8.0, 1 mM Spermidine, 0.3 mM Spermine, RNase-A
592 40 µg/ml) and incubated at 55 °C for 15 min. 50 U of Proteinase K was added and
593 samples incubated at 55 °C for 1 h. DNA was isolated with 25:24:1
594 Phenol:Chloroform:Isoamyl Alcohol (Ambion) followed by ethanol precipitation. Fragments
595 from 50 to 550 bp were selected using agarose gel electrophoresis. The extracted DNA
596 samples were quantified fluorometrically with a Qubit 3.0 Fluorometer (Life Technologies),
597 and a total of 10 ng of digested DNA (200 pg l⁻¹) was used for library construction.

598 Initial sample quality control of pre-fragmented DNA was assessed using a TapeStation
599 DNA 1000 High Sensitivity Screen Tape (Agilent, Cheadle UK). Sequencing ready libraries
600 were prepared using the Hyper Prep DNA Library preparation kit (Kapa Biosystems,
601 London UK) selecting fragments from 70-350 bp for optimization (see He et al., 2014) and
602 indexed for pooling using NextFlex DNA barcoded adapters (Bioo Scientific, Austin TX
603 US). In order to reduce bias due to amplification of DNA fragments by the polymerase
604 chain reaction, as recommended by the manufacturers, a low number of cycles (17 cycles)
605 was used. Libraries were quantified using a TapeStation DNA 1000 Screen Tape and by
606 qPCR using an NGS Library Quantification Kit (KAPA Biosystems) on an AriaMx qPCR
607 system (Agilent) and then normalised, pooled, diluted and denatured for sequencing on
608 the NextSeq 500 (Illumina, Chesterford UK). The main library was spiked at 10% with the
609 PhiX control library (Illumina). Sequencing was performed using Illumina NextSeq in the
610 Departments of Biochemistry and Pathology at the University of Cambridge, UK, with 2x75
611 cycles of sequencing. For the deproteinized DNase I seq experiments 1 µg of
612 deproteinized DNA was resuspended in 1 ml of digestion buffer (15 mM Tris-HCl, 90 mM
613 NaCl, 60 mM KCl, 6 mM CaCl₂, 0.5 mM spermidine, 1 mM EDTA and 0.5 mM EGTA, pH
614 8.0). DNaseI (Fermentas) at 2.5 U was added to each tube and incubated at 37 °C for 2
615 min. Digestion was arrested with addition of 1:1 volume of stop buffer (50 mM Tris-HCl,
616 100 mM NaCl, 0.1% (w/v) SDS, 100 mM EDTA, pH 8.0, 1 mM Spermidine, 0.3 mM

617 Spermine, RNase A 40 µg/ml) and incubated at 55 °C for 15 min. 50 U of Proteinase K
618 was added and samples incubated at 55 °C for 1 h. DNA was isolated by mixing with 1 ml
619 25:24:1 Phenol:Chloroform:Isoamyl Alcohol (Ambion) and spun for 5 min at 13,00 rpm
620 followed by ethanol precipitation of the aqueous phase. Samples were then size-selected
621 (50-400 bp) using agarose gel electrophoresis. The extracted DNA samples were
622 quantified fluorometrically using Qubit 3.0 Fluorometer (Life technologies), and a total of 1
623 ng of digested DNA was used for library construction. Sequencing ready libraries were
624 prepared using a KAPA Hyper Prep Kit (KAPA Biosystems, London UK) according to the
625 manufacturer's instructions. In order to reduce bias due to amplification of DNA fragments
626 by the polymerase chain reaction, as recommended by the manufacturers, 17 cycles were
627 used. Quality of the libraries were checked using a Bioanalyzer High Sensitivity DNA Chip
628 (Agilent Technologies). Libraries were quantified by Qubit 3.0 Fluorometer (Life
629 Technologies) and qPCR using an NGS Library Quantification Kit (KAPA Biosystems) and
630 then normalised, pooled, diluted and denatured for paired end sequencing using High
631 Output 150 cycle run (2x 75 bp reads). Sequencing was performed using NextSeq 500
632 (Illumina, Chesterford UK) in the Sainsbury Laboratory University of Cambridge, UK, with
633 2x75 cycles of sequencing.

634

635 **DNaseI-SEQ Data processing**

636 Genome sequences were downloaded from Phytozome (v10) (Goodstein et al., 2012).
637 The following genome assemblies were used: Bdistachyon_283_assembly_v2.0;
638 Sbicolor_255_v2.0; Sitalica_164_v2; Zmays_284_AGPv3. Due to the lack of guidelines for
639 DNaseI-SEQ experiments in plants we followed the guidelines from the Encyclopaedia of
640 DNA Elements (ENCODE3). Reads were mapped to genomes using bowtie2 (Langmead
641 and Salzberg, 2012) and processed using samtools (Li et al., 2009) to remove those with a
642 MAPQ score <42. DHS were called using MACS2 (Feng et al., 2012) and the final set of
643 peak calls were determined using the irreproducible discovery rate (IDR) (Li and Dewey,
644 2011), calculated using the script batch_consistency_analysis.R
645 ([https://github.com/modENCODE-](https://github.com/modENCODE-DCC/Galaxy/blob/master/modENCODE_DCC_tools/idr/batch-consistency-analysis.r)
646 [DCC/Galaxy/blob/master/modENCODE_DCC_tools/idr/batch-consistency-analysis.r](https://github.com/modENCODE-DCC/Galaxy/blob/master/modENCODE_DCC_tools/idr/batch-consistency-analysis.r)). The
647 Irreproducible Discovery Rate framework adapted from the ENCODE 3 pipeline (Marinov
648 et al., 2014; <https://sites.google.com/site/anshulkundaje/projects/idr>) aims to measure the
649 reproducibility of findings by identifying the point (threshold) in which peaks are no longer
650 consistent across replicates.

651

652 **Quality metrics and identification of Digital Genomic Footprints (DGF)**

653 SPOT score (number of a subsample of mapped reads (5M) in DHS/Total number of
654 subsampled, mapped reads (5M) (John et al., 2011)) was calculated using BEDTools
655 (Quinlan and Hall, 2010) to determine the number of mapped reads possessing at least 1
656 bp overlap with a DHS site. Normalized Strand Cross-correlation coefficient (NSC) and
657 Relative Strand Cross-correlation coefficient (RSC) scores were calculated using SPP
658 (Kharchenko et al., 2008) and PCR bottleneck coefficient (PBC) was calculated using
659 BEDTools. To account for cutting bias associated with the DNaseI enzyme DNaseI-SEQ
660 on naked DNA was performed. These data were used to generate background signal
661 profiles and calculate the footprint log-likelihood ratio for each footprint using the R
662 package MixtureModel (Yardimci et al., 2014) such that those with low log likelihood ratios
663 (FLR <0) were removed. Digital Genomic Footprints (DGF) were identified using
664 Wellington (Piper et al., 2013) and differential DGF were identified using Wellington
665 bootstrap (Piper et al., 2015).

666

667 **Data visualisation**

668 DHS and DGF sequences were loaded into and visualized in the Integrative Genomics
669 Viewer (Thorvaldsdóttir et al., 2013) and figures produced in Inkscape, plots were
670 generated with R package ggplot2 (Wickham, 2010) and figures depicting conservation of
671 DGF or motifs between orthologous sequences were generated using genoplotR (Guy et
672 al., 2010). Word clouds were created with the wordcloud R package (Fellows, 2012).
673 TreeView images were produced by processing DGF data using
674 'dnase_to_javatreeview.py' from pyDNase (Piper et al., 2013, 2015) and loaded into
675 TreeView (Saldanha, 2004). Average cut density plots were generated using the script
676 'dnase_average_profile.py' from pyDNase. Genomic features were annotated and
677 distribution calculated using PAVIS (Huang et al., 2013) and plotted using ggplot2. Circular
678 plots showing the distribution of ChIP-SEQ peaks, DHS and DGF across the maize
679 genome was generated using the R package circlize (Gu, 2014).

680

681 **DNase cutting bias calculations and ChIP-SEQ analysis**

682 After sequencing, the number of DNA 6-mer centred at each DNase cleavage site
683 (between 3rd and 4th base) was counted and normalized by the total number of counts.
684 Next, DNA 6-mer frequencies were normalized by the frequencies of each DNA 6-mer in
685 the genome. The resulting background signal profile was used as input in the

686 FootprintMixture.R package (https://ohlerlab.mdc-berlin.de/software/FootprintMixture_109/) (Supplemental Figure 2).

688 ChIP-SEQ peaks from 117 transcription factors were obtained from the pre-release
689 maize cistrome data collection (<http://www.epigenome.cuhk.edu.hk/C3C4.html>).
690 Permutation tests between ChIP peaks and DHS or DGF were performed using regioneR
691 (Gel et al., 2016) using 100 permutations.

692

693 ***de novo* motif prediction, motif scanning and enrichment testing**

694 *de novo* motif prediction was performed using findMotifsGenome.pl script from the
695 HOMER suite (Heinz et al., 2010) using digital genomic footprints (DGF) as input together
696 with the reference genome sequence for each species. Motif scanning was performed
697 using FIMO (Grant et al., 2011) with default parameters. To determine overrepresentation
698 of TF family motifs in samples hypergeometric tests were performed using R. The
699 distribution of each motif across different genomic features was obtained for each
700 annotated motif by dividing the number of hits in a particular feature by the total number of
701 hits in the genome.

702

703 **Whole genome alignments, pairwise cross mapping of genomic features and variant** 704 **data processing**

705 To cross map genomic features between species, mapping files were generated
706 according to (http://genomewiki.ucsc.edu/index.php/Whole_genome_alignment_howto)
707 using tools from the UCSC Genome Browser, including trfBig, faToNib, faSize, lavToPsl,
708 faSplit, axtChain, chainNet (Kent et al., 2002) and LASTZ (Harris, 2007). Briefly, whole
709 genome alignment was performed with LASTZ, matching alignments next to each other
710 were chained together using axtChain, sorted with axtSort, then netted together to form
711 larger blocks with chainNet. Genomic features were then mapped between genomes
712 using bnMapper (Denas et al., 2015). For the variant analysis on duons, *Z. mays* variant
713 data (Bukowski et al., 2018) was downloaded from
714 <http://cbsusrv04.tc.cornell.edu/users/panzea/download.aspx?filegroupid=16> following
715 instructions. After downloading vcf files were annotated using SnpEff (Cingolani et al,
716 2012; <https://doi.org/10.4161/fly.19695>) with the B73_RefGen_v4 genome assembly
717 specifically to allow identification of non-synonymous sites. A custom script was used to
718 identify all four-fold degenerate sites (FFDS) in the *Z. mays* genome. This bed file in turn
719 was used to identify which of the synonymous polymorphic sites were FFDS. Each
720 polymorphic site had its allele frequencies calculated. Putative *Z. mays* duons were

721 identified by intersecting (with bedtools intersect) the final DGF identified with exonic
722 regions. These duons were then used to extract only those exons within which a duon was
723 found. These exons in turn had the duon regions themselves subtracted to leave the exon
724 region except the duon. This provided the surrounding exonic sequences with which to
725 compare to the duons. These two regions were then intersected with the polymorphism
726 data to identify both the number of occurrences and allelic frequencies of polymorphic
727 sites (FFDS and non-synonymous) within both the duons and their surrounding exonic
728 sequences.

729

730 **Accession numbers**

731 Methods for DNaseI digestion are on protocols.io
732 ([dx.doi.org/10.17504/protocols.io.hdfb23n](https://doi.org/10.17504/protocols.io.hdfb23n)). Raw sequencing data and processed files are
733 deposited in Gene Expression Omnibus (GSE97369). For full methods, commands and
734 scripts, please see github ([https://github.com/hibberd-lab/Burgess-Reyna_llorens-](https://github.com/hibberd-lab/Burgess-Reyna_llorens-monocot-DNase)
735 [monocot-DNase](https://github.com/hibberd-lab/Burgess-Reyna_llorens-monocot-DNase)) and Figshare 10.6084/m9.figshare.7649450.

736 **Supplemental Legends**

737 **Supplemental Figure 1:** DNaseI digestion of nuclei for sequencing. Representative
738 images of digested samples separated on 2% (w/v) agarose gels by electrophoresis. (A)
739 *S. bicolor* whole leaf (WL); (B) *S. bicolor* Bundle Sheath (BS); (C) *Z. mays* WL; (D) *Z.*
740 *mays* BS; (E) *B. distachyon* WL; (F) *S. italica* WL; (G) *S. italica* BS. Each gel represents a
741 separate biological replicate, and the units of DNaseI used are illustrated above. Samples
742 selected for sequencing are indicated in red.

743

744 **Supplemental Figure 2:** Bias in DNaseI-SEQ cleavage. (A) TreeView diagrams
745 illustrating cut density around individual digital genomic footprint (DGF) predicted from
746 performing DNaseI-SEQ on deproteinated genomic DNA from each species. Each row
747 represents an individual DGF, cuts are coloured according to whether they align to the
748 positive (red) or negative (blue) strand and indicate increased cutting in a 100 bp window
749 on either side of the DGF. (B) Pearson correlation coefficient of DNase I cleavage bias
750 between *Z. mays*, *S. bicolor*, *S. italica* and *B. distachyon*. (C) Schematic illustrating the
751 process adopted to determine DNaseI cutting bias and then normalize to allow digital
752 genomic footprinting.

753

754 **Supplemental Figure 3:** Saturation analysis of footprints. Digital genomic footprints were
755 predicted from subsets (12.5, 25, 50, 75 and 100%) of uniquely mapped reads obtained
756 from DNaseI-SEQ of whole leaf samples in each species.

757

758 **Supplemental Figure 4:** Genome wide comparison of DGF and ChIP-SEQ peaks from
759 117 maize transcription factors. (A) Density plot of DHS, DGF, ChIP-SEQ peaks and
760 intersecting DGF/ChIP-SEQ peaks across the maize genome. The center of the plot
761 shows a word cloud representing transcription factor families in the ChIP-SEQ dataset. (B)
762 Effect of shifting DHS and DGF features from their original position on local z-scores as
763 determined by permutation tests between ChIP-SEQ peaks derived from 117 transcription
764 factors. A total of 100 permutations were performed for each comparison. The sharper
765 peak derived from shifting the DGF indicates a higher sensitivity to position and therefore
766 strong overlap with ChIP-SEQ data. (C) Density plot depicting the distribution of ChIP-SEQ
767 signals per kilobase (kb) from the transcription start site (TSS) of *Z. mays*.

768

769 **Supplemental Figure 5:** Density plot depicting the distribution of DGF per kilobase (kb)
770 from the transcription start site (TSS) of *S. bicolor*, *Z. mays*, *S. italica* and *B. distachyon*
771 whole leaves.

772

773 **Supplemental Figure 6:** Nucleotide proportion of duons and surrounding exons used in
774 the substitution analysis for *Z. mays*. The frequency of each nucleotide was divided by the
775 total length to determine nucleotide proportions across duons as well as surrounding exon
776 sequences.

777

778 **Supplemental Figure 7:** Transcript abundance for genes in mesophyll and bundle sheath
779 cells associated with DHS and DGF in *S. bicolor*. (A) Cell preferential gene expression
780 profiles of highly abundant M and BS genes expressed as transcripts per million reads
781 (TPM). (B) Schematic representing DHS, DGF and DE DGF present in whole leaf (blue)
782 and BS (orange) of *S. bicolor*.

783

784 **Supplemental Figure 8:** Differential accessibility of broad regulatory regions in *S. bicolor*
785 is not sufficient for cell preferential gene expression. Percentage of differentially detected
786 DHS among BS and M specific genes in *S. bicolor* compared with randomly generated
787 gene samples (n=50).

788

789 **Supplemental Figure 9:** Representation of the C₄ pathway showing differentially
790 accessible DHS, DGF and cell specific DGF between whole leaf (blue) and bundle sheath
791 (orange) samples in *S. bicolor*. CA; Carbonic Anhydrase, PEPC; Phosphoenolpyruvate
792 carboxylase, PPDK; Pyruvate, orthophosphate dikinase, MDH; Malate dehydrogenase,
793 NADP-ME; NADP-dependent malic enzyme, RBCS1A; Ribulose bisphosphate
794 carboxylase small subunit1A, OAA; Oxaloacetate, Mal; Malate, PEP;
795 Phosphoenolpyruvate, Pyr; Pyruvate, Asp; Aspartate.

796

797 **Supplemental Figure 10:** Representation of the C₄ pathway showing differentially
798 accessible DHS, DGF and cell specific DGF between whole leaf (blue) and bundle sheath
799 (orange) samples in *S. italica*. CA; Carbonic Anhydrase, PEPC; Phosphoenolpyruvate
800 carboxylase, PPDK; Pyruvate, orthophosphate dikinase, MDH; Malate dehydrogenase,
801 NADP-ME; NADP-dependent malic enzyme, RBCS1A; Ribulose bisphosphate
802 carboxylase small subunit1A, OAA; Oxaloacetate, Mal; Malate, PEP;
803 Phosphoenolpyruvate, Pyr; Pyruvate, Asp; Aspartate.

804

805 **Supplemental Figure 11:** Representation of the C₄ pathway showing differentially
806 accessible DHS, DGF and cell specific DGF between whole leaf (blue) and bundle sheath
807 (orange) samples in *Z. mays*. CA; Carbonic Anhydrase, PEPC; Phosphoenolpyruvate
808 carboxylase, PPK; Pyruvate, orthophosphate dikinase, MDH; Malate dehydrogenase,
809 NADP-ME; NADP-dependent malic enzyme, RBCS1A; Ribulose bisphosphate
810 carboxylase small subunit1A, OAA; Oxaloacetate, Mal; Malate, PEP;
811 Phosphoenolpyruvate, Pyr; Pyruvate, Asp; Aspartate.

812

813

814 **Supplemental Table 1:** Summary of DNaseI-SEQ Quality Metrics. Values include the
815 number of uniquely mapped reads with MAPQ scores ≥ 42 (NMAP), PCR bottleneck
816 Coefficient (PBC), Normalized Strand Cross-correlation Coefficient (NSC), Relative Strand
817 Cross-correlation Coefficient (RSC) the optimal number of peaks calculated by IDR
818 method (PEAKS), and the Signal Portion Of Tags (SPOT) for whole leaf and bundle
819 sheath samples from *B. distachyon*, *S. italica*, *S. bicolor* and *Z. mays*.

820

821 **Supplemental Table 2:** Summary Statistics for Genomic Features Identified in Whole Leaf
822 Samples. Including information about the number of DHS and DGF identified per sample,
823 and the number of genes that could be annotated with at least one genomic feature for *B.*
824 *distachyon*, *S. italica*, *S. bicolor* and *Z. mays*.

825

826 **Supplemental Table 3:** DNaseI cutting bias calculation summary for whole leaf and
827 bundle sheath data.

828

829 **Supplemental Table 4:** Transcription factors included in the ChIP-SEQ data analysis.

830

831 **Supplemental Table 5:** Summary statistics of DNaseI-SEQ analysis of Bundle Sheath
832 samples

833

834 **Supplemental Table 6:** Summary Statistics of Overlap between DHS in Whole Leaf and
835 Bundle Sheath Samples.

836

837 **Supplemental Table 7:** Summary Statistics of Overlap between DGF in Whole Leaf and
838 Bundle Sheath Samples.

839

840 **Supplemental Table 8:** Summary Statistics for Differential Digital Genomic Footprint
841 Calling. Including the total number of differential DGF (DE DGF) and the number of DE
842 DGF in DE genes for *S. italica*, *S. bicolor* and *Z. mays*.

843

844 **Supplemental Table 9:** Motifs mapped to genes of the C₄, CBB and C₂ cycles in *Z. mays*,
845 *S. bicolor*, *S. italica* for whole leaf and bundle sheath samples and in *B. distachyon* for
846 whole leaf samples

847

848 **Supplemental Table 10:** Hypergeometric tests for enrichment of individual motifs in *Z.*
849 *mays*, *S. bicolor*, *S. italica* for whole leaf and bundle sheath samples.

850

851 **Supplemental Table 11:** Hypergeometric tests for enrichment of cell specific individual
852 motifs in *Z. mays*, *S. bicolor*, *S. italica* for whole leaf and bundle sheath samples.

853

854 **Supplemental Table 12:** Number of genes in C₄, CBB and C₂ cycles annotated with a
855 given motif in *Z. mays*, *S. italica*, *S. bicolor* and *B. distachyon*.

856

857 **Supplemental Table 12:** Statistics for Cross Mapping of genomic features between *S.*
858 *bicolor*, *S. italica*, *Z. mays* and *B. distachyon*.

859

860 **Supplemental Table 13:** DGF conserved and occupied in *Z. mays*, *S. bicolor*, *S. italica* for
861 whole leaf and bundle sheath samples and in *B. distachyon* for whole leaf samples.

862

863 **Supplemental Table 14:** DGF in C₄ genes that are conserved between *Z. mays* and *S.*
864 *bicolor*.

865

866 **Supplemental Table 15:** DGF conserved in all four species.

867

868 **Supplemental Table 16:** Gene Ontology analysis on hyper-conserved DGF in whole leaf
869 samples of *S. italica*, *S. bicolor*, *Z. mays* and *B. distachyon*.

870

871

872 **Acknowledgements:** SJB was supported by the 3to4 grant from the EU and BB/I002243
873 from the BBSRC, IRL by CONACyT and BBSRC grant BB/L014130, SRS and PS by an

874 Advanced ERC grant 694733 Revolution to JMH, and KJ by a Gatsby Career
875 Development Fellowship. We would like to acknowledge Aslihan Karabacak for support in
876 implementing the FootprintMixture package.

877

878

879 **Contributions:** SJB, I-RL and JMH conceptualised the experiments. SJB and I-RL grew
880 and harvested nuclei from *S. bicolor*, *S. italica* and *Z. mays*. KJ provided the nuclei from *B.*
881 *distachyon*. SJB and I-RL performed DNase I experiments and data analysis. PS extracted
882 nuclei and performed DNaseI experiments on naked DNA, SRS undertook the variant
883 analysis. SJB, I-RL, SRS and JMH wrote the manuscript and prepared the figures.

884 **References**

- 885 Bauwe, H., Hagemann, M., and Fernie, A.R. (2010). Photorespiration: players, partners
886 and origin. *Trends Plant Sci.* 15: 330–6.
- 887 Bolduc, N. et al. (2012a). A Genome-Wide Regulatory Framework Identifies Maize
888 Pericarp Color1 Controlled Genes. *Plant Cell* 27: 543–553.
- 889 Bolduc, N. and Hake, S. (2009). The Maize Transcription Factor KNOTTED1 Directly
890 Regulates the Gibberellin Catabolism Gene *ga2ox1*. *Plant Cell* 21: 1647–1658.
- 891 Bolduc, N., Yilmaz, A., Mejia-Guerra, M.K., Morohashi, K., O'Connor, D., Grotewold, E.,
892 and Hake, S. (2012b). Unraveling the KNOTTED1 regulatory network in maize
893 meristems. *Genes Dev.* 26: 1685–1690.
- 894 Borba, A. R., Serra, T. S., Górska, A., Gouveia, P., Cordeiro, A. M., Reyna-Llorens, I.,
895 Kneřová, J., Barros, P.M., Abreu, I.A., Oliveira, Hibberd, J.M., Saibo, N. (2018).
896 Synergistic binding of bHLH transcription factors to the promoter of the maize NADP-
897 ME gene used in C4 photosynthesis is based on an ancient code found in the
898 ancestral C3 state. *Mol. Bio and Evolution* 35: 1690-1705.
- 899 Bowes, G., Ogren, W.L., and Hageman, R.H. (1971). Phosphoglycolate production
900 catalyzed by ribulose diphosphate carboxylase. *Biochem. Biophys. Res. Commun.* 45:
901 716–722.
- 902 Brown, N.J., Newell, C.A., Stanley, S., Chen, J.E., Perrin, A.J., Kajala, K., and Hibberd,
903 J.M. (2011). Independent and Parallel Recruitment of Preexisting Mechanisms
904 Underlying C₄ Photosynthesis. *Science* 331: 1436–1439.
- 905 Bukowski, R., Guo, X., Lu, Y., Zou, C., He, B., Rong, Z., Wang, B., Xu, D., Yang, B., Xie,
906 C. and Fan, L. (2017). Construction of the third-generation *Zea mays* haplotype map.
907 *GigaScience* 7: gix134.
- 908 Chang YM, Liu WY, Shih ACC, Shen MN, Lu CH, Lu MYJ, Yang HW, Wang TY, Chen SC,
909 Chen SM, et al. (2012) Characterizing regulatory and functional differentiation between
910 maize mesophyll and bundle sheath cells by transcriptomic analysis. *Plant Physiol* 160:
911 165–177.
- 912 Christin, P.A., Salamin, N., Savolainen, V., Duvall, M.R., and Besnard, G. (2007). C₄
913 photosynthesis evolved in grasses via parallel adaptive genetic changes. *Curr. Biol.*
914 17: 1241–1247.
- 915 Christin, P. A., & Osborne, C. P. (2014). The evolutionary ecology of C₄ plants. *New*
916 *Phytologist* 204: 765-781.
- 917 Cingolani, P., Platts, A., Wang, L.L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X.
918 and Ruden, D.M. (2012). A program for annotating and predicting the effects of single
919 nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster*
920 strain w1118; iso-2; iso-3. *Fly* 6: 80-92.
- 921 Cornejo, M.-J., Luth, D., Blankenship, K.M., Anderson, O.D., and Blechl, A.E. (1993).
922 Activity of a maize ubiquitin promoter in transgenic rice. *Plant Mol. Biol.* 23: 567–581.
- 923 Covshoff, S., Furbank, R.T., Leegood, R.C., and Hibberd, J.M. (2013). Leaf rolling allows
924 quantification of mRNA abundance in mesophyll cells of sorghum. *J. Exp. Bot.* 64:
925 807–813.
- 926 Deal, R.B. and Henikoff, S. (2011). The INTACT method for cell type-specific gene
927 expression and chromatin profiling in *Arabidopsis thaliana*. *Nat. Protoc.* 6: 56–68.
- 928 Denas, O., Sandstrom, R., Cheng, Y., Beal, K., Herrero, J., Hardison, R.C., and Taylor, J.
929 (2015). Genome-wide comparative analysis reveals human-mouse regulatory
930 landscape and evolution. *BMC Genomics* 16: 87.
- 931 Emms, D. M., Covshoff, S., Hibberd, J. M., & Kelly, S. (2016). Independent and parallel
932 evolution of new genes by gene duplication in two origins of C₄ photosynthesis
933 provides new insight into the mechanism of phloem loading in C₄ species. *Mol Bio*
934 *and Evol.* 33(7), 1796-1806.

935 Eveland, A.L., Goldshmidt, A., Pautler, M., Morohashi, K., Liseron-Monfils, C., Lewis,
936 M.W., Kumari, S., Hiraga, S., Yang, F., Unger-Wallace, E. and Olson, A. (2014).
937 Regulatory modules controlling maize inflorescence architecture. *Genome Research*
938 24: 431-443.

939 Fellows, I. (2012). wordcloud: Word clouds. R Packag. version 2: 109.

940 Feng, J., Liu, T., Qin, B., Zhang, Y., and Liu, X.S. (2012). Identifying ChIP-SEQ
941 enrichment using MACS. *Nat. Protoc.* 7: 1728–1740.

942 Furbank, R.T. (2011). Evolution of the C₄ photosynthetic mechanism: are there really three
943 C₄ acid decarboxylation types? *J. Exp. Bot.* 62: 3103–3108.

944 Furbank, R.T., Stitt, M., and Foyer, C.H. (1985). Intercellular compartmentation of sucrose
945 synthesis in leaves of *Zea mays* L. *Planta* 164: 172–178.

946 Gel, B., Díez-Villanueva, A., Serra, E., Buschbeck, M., Peinado, M. A., & Malinverni, R.
947 (2015). regioneR: an R/Bioconductor package for the association analysis of genomic
948 regions based on permutation tests. *Bioinformatics* 32: 289-291.

949 Giuliano, G., Pichersky, E., Malik, V.S., Timko, M.P., Scolnik, P.A., and Cashmore, A.R.
950 (1988). An evolutionarily conserved protein binding sequence upstream of a plant light
951 regulated gene. *Proc. Natl. Acad. Sci.* 85: 7089–7093.

952 Goodstein, D.M., Shu, S., Howson, R., Neupane, R., Hayes, R.D., Fazo, J., Mitros, T.,
953 Dirks, W., Hellsten, U., Putnam, N., and Rokhsar, D.S. (2012). Phytozome: a
954 comparative platform for green plant genomics. *Nucleic Acids Res.* 40: D1178-86.

955 Gowik, U., Burscheidt, J., Akyildiz, M., Schlue, U., Koczor, M., Streubel, M., and Westhoff,
956 P. (2004). *cis*-regulatory elements for mesophyll-specific gene expression in the C₄
957 plant *Flaveria trinervia*, the promoter of the C₄ PHOSPHOENOLPYRUVATE
958 CARBOXYLASE gene. *Plant Cell* 16: 1077–1090.

959 Grant, C.E., Bailey, T.L., and Noble, W.S. (2011). FIMO: scanning for occurrences of a
960 given motif. *Bioinformatics* 27: 1017–1018.

961 Gu, Z., Gu, L., Eils, R., Schlesner, M., & Brors, B. (2014). circlize implements and
962 enhances circular visualization in R. *Bioinformatics* 30: 2811-2812.

963 Guy, L., Roat Kultima, J., and Andersson, S.G.E. (2010). genoPlotR: comparative gene
964 and genome visualization in R. *Bioinformatics* 26: 2334–2335.

965 Harris, R.S. (2007). Improved pairwise alignment of genomic DNA. The Pennsylvania
966 State University. (PhD Thesis).

967 Hatch, M.D. (1987). C₄ photosynthesis: a unique elend of modified biochemistry, anatomy
968 and ultrastructure. *Biochim. Biophys. Acta - Rev. Bioenerg.* 895: 81–106.

969 He, H.H., Meyer, C.A., Chen, M.W., Zang, C., Liu, Y., Rao, P.K., Fei, T., Xu, H., Long, H.,
970 Liu, X.S. and Brown, M. (2014). Refined DNase-seq protocol and data analysis
971 reveals intrinsic bias in transcription factor footprint identification. *Nature Methods* 11:
972 73.

973 Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y.C., and Laslo, P. (2010). Simple
974 combinations of lineage-determining transcription factors prime cis-regulatory
975 elements required for macrophage and B cell identities. *Mol. Cell.* 38: 576-89

976 Hesselberth, J.R., Chen, X., Zhang, Z., Sabo, P.J., Sandstrom, R., Reynolds, A.P.,
977 Thurman, R.E., Neph, S., Kuehn, M.S., Noble, W.S., Fields, S., and
978 Stamatoyannopoulos, J.A. (2009). Global mapping of protein-DNA interactions in vivo
979 by digital genomic footprinting. *Nat. Methods* 6: 283–9.

980 Hibberd, J.M. and Covshoff, S. (2010). The Regulation of Gene Expression Required for
981 C₄ Photosynthesis. *Annu. Rev. Plant Biol.* 61: 181–207.

982 Huang, W., Loganantharaj, R., Schroeder, B., Fargo, D., & Li, L. (2013). PAVIS: a tool for
983 Peak Annotation and Visualization. *Bioinformatics* 29: 3097-3099.

984 Jeon, J.-S., Lee, S., Jung, K.-H., Jun, S.-H., Kim, C., and An, G. (2000). Tissue-
985 Preferential Expression of a Rice α -Tubulin Gene, *OsTubA1*, Mediated by the First
986 Intron. *Plant Physiol.* 123: 1005–1014.

- 987 John, S., Sabo, P.J., Thurman, R.E., Sung, M.-H., Biddie, S.C., Johnson, T.A., Hager,
988 G.L., and Stamatoyannopoulos, J.A. (2011). Chromatin accessibility pre-determines
989 glucocorticoid receptor binding patterns. *Nat Genet.* 43: 264–268.
- 990 John, C. R., Smith-Unna, R. D., Woodfield, H., Covshoff, S., & Hibberd, J. M. (2014).
991 Evolutionary convergence of cell-specific gene expression in independent lineages of
992 C4 grasses. *Plant Phys.*, 165(1), 62-75.
- 993 Jordan, D.B. and Ogren, W.L. (1984). The CO₂/O₂ specificity of Ribulose 1,5-Bisphosphate
994 Carboxylase Oxygenase - dependence on Ribulose bisphosphate concentration, pH
995 and temperature. *Planta* 161: 308–313.
- 996 Kajala, K., Williams, B.P., Brown, N.J., Taylor, L.E., and Hibberd, J.M. (2011). Multiple
997 Arabidopsis genes primed for direct recruitment into C₄ photosynthesis. *Plant J.* 69:
998 47–56.
- 999 Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and
1000 Haussler, and D. (2002). The Human Genome Browser at UCSC. *Genome Res.* 12:
1001 996–1006.
- 1002 Khan, A., Fornes, O., Stigliani, A., Gheorghe, M., Castro-Mondragon, J.A., van der Lee,
1003 R., Bessy, A., Cheneby, J., Kulkarni, S.R., Tan, G. and Baranasic, D.,. (2017).
1004 JASPAR 2018: update of the open-access database of transcription factor binding
1005 profiles and its web framework. *Nucleic Acids Research* 46 (D1): D260-D266.
- 1006 Kharchenko, P. V, Tolstorukov, M.Y., and Park, P.J. (2008). Design and analysis of ChIP-
1007 SEQ experiments for DNA-binding proteins. *Nat Biotech.* 26: 1351–1359.
- 1008 Kozaki, A., Hake, S., & Colasanti, J. (2004). The maize ID1 flowering time regulator is a
1009 zinc finger protein with novel DNA binding properties. *Nucleic Acids Research*, 32:
1010 1710-1720.
- 1011 Langmead, B. and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2.
1012 *Nature Methods* 9: 357–359.
- 1013 Leegood, R.C. (1985). The intercellular compartmentation of metabolites in leaves of *Zea*
1014 *mays* L. *Planta* 164: 163–171.
- 1015 Li, B. and Dewey, C.N. (2011). RSEM: accurate transcript quantification from RNA-Seq
1016 data with or without a reference genome. *BMC Bioinformatics* 12: 323.
- 1017 Li, C., Qiao, Z., Qi, W., Wang, Q., Yuan, Y., Yang, X., Tang, Y., Mei, B., Lv, Y., Zhao, H.
1018 and Xiao, H. (2015). Genome-wide characterization of cis-acting DNA targets reveals
1019 the transcriptional regulatory framework of opaque2 in maize. *The Plant Cell* 27: 532–
1020 545.
- 1021 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis,
1022 G., Durbin, R., and Subgroup, 1000 Genome Project Data Processing (2009). The
1023 Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25: 2078–2079.
- 1024 de Lucas, M., Pu, L., Turco, G., Gaudinier, A., Morao, A.K., Harashima, H., Kim, D., Ron,
1025 M., Sugimoto, K., Roudier, F., and Brady, S.M. (2016). Transcriptional Regulation of
1026 Arabidopsis Polycomb Repressive Complex 2 Coordinates Cell-Type Proliferation and
1027 Differentiation. *Plant Cell* 28: 2616 -2631.
- 1028 Maas, C., Laufs, J., Grant, S., Korfhage, C., and Werr, W. (1991). The combination of a
1029 novel stimulatory element in the first exon of the maize *SHRUNKEN-1* gene with the
1030 following intron 1 enhances reporter gene expression up to 1000-fold. *Plant Mol. Biol.*
1031 16: 199–207.
- 1032 Maher, K.A., Bajic, M., Kajala, K., Reynoso, M., Pauluzzi, G., West, D.A., Zumstein, K.,
1033 Woodhouse, M., Bubb, K., Dorrity, M.W. and Queitsch, C. (2018). Profiling of
1034 accessible chromatin regions across multiple plant species and cell types reveals
1035 common gene regulatory principles and new control modules. *The Plant Cell* 30: 15-
1036 36.
- 1037 Marinov, G.K., Kundaje, A., Park, P.J., and Wold, B.J. (2014). Large-Scale Quality
1038 Analysis of Published ChIP-SEQ Data. *G3* 4: 209–223.

- 1039 Markelz, N. H., Costich, D. E., & Brutnell, T. P. (2003). Photomorphogenic responses in
1040 maize seedling development. *Plant Physiology*, 133(4), 1578-1591.
- 1041 Martin, A. and Orgogozo, V. (2013). The Loci of repeated evolution: a catalog of genetic
1042 hotspots of phenotypic variation. *Evolution* 67: 1235–50.
- 1043 Natarajan, A., Yardimci, G.G., Sheffield, N.C., Crawford, G.E., and Ohler, U. (2012).
1044 Predicting cell-type-specific gene expression from regions of open chromatin.
1045 *Genome Research* 22: 1711–1722.
- 1046 Neph, S., Vierstra, J., Stergachis, A.B., Reynolds, A.P., Haugen, E., Vernet, B., Thurman,
1047 R.E., John, S., Sandstrom, R., Johnson, A.K. and Maurano, M.T. (2012). An
1048 expansive human regulatory lexicon encoded in transcription factor footprints. *Nature*
1049 489: 83-90.
- 1050 Niu, X., Helentjaris, T., & Bate, N. J. (2002). Maize ABI4 binds coupling element1 in
1051 abscisic acid and sugar response genes. *The Plant Cell*, 14: 2565-2575.
- 1052 O'Malley, R.C., Huang, S.C., Song, L., Lewsey, M.G., Bartlett, A., Nery, J.R., Galli, M.,
1053 Gallavotti, A., and Ecker, J.R. (2016). Cistrome and Epicistrome Features Shape the
1054 Regulatory DNA Landscape. *Cell* 165: 1280–1292.
- 1055 Pajoro, A., Madrigal, P., Muiño, J.M., Matus, J.T., Jin, J., Mecchia, M.A., Debernardi, J.M.,
1056 Palatnik, J.F., Balazadeh, S., Arif, M. and Ó'Maoiléidigh, D.S. (2014). Dynamics of
1057 chromatin accessibility and gene regulation by MADS-domain transcription factors in
1058 flower development. *Genome Biology* 15: R41.
- 1059 Patel, M., Corey, A.C., Yin, L.P., Ali, S.J., Taylor, W.C., and Berry, J.O. (2004).
1060 Untranslated regions from C₄ Amaranth *AhRbcS1* mRNAs confer translational
1061 enhancement and preferential bundle sheath cell expression in transgenic C₄ *Flaveria*
1062 *bidentis*. *Plant Physiol.* 136: 3550–3561.
- 1063 Pautler, M., Eveland, A.L., LaRue, T., Yang, F., Weeks, R., Lunde, C., Je, B. II, Meeley,
1064 R., Komatsu, M., Vollbrecht, E., Sakai, H., and Jackson, D. (2015). *FASCIATED*
1065 *EAR4* Encodes a bZIP Transcription Factor That Regulates Shoot Meristem Size in
1066 Maize. *Plant Cell* 27: 104-120.
- 1067 Piper, J., Assi, S.A., Cauchy, P., Ladroue, C., Cockerill, P.N., Bonifer, C., and Ott, S.
1068 (2015). Wellington-bootstrap: differential DNase-seq footprinting identifies cell-type
1069 determining transcription factors. *BMC Genomics* 16: 1000.
- 1070 Piper, J., Elze, M.C., Cauchy, P., Cockerill, P.N., Bonifer, C., and Ott, S. (2013).
1071 Wellington: a novel method for the accurate identification of digital genomic footprints
1072 from DNase-seq data. *Nucleic Acids Research* 41: e201–e201.
- 1073 Quinlan, A.R. and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing
1074 genomic features. *Bioinformatics* 26: 841–842.
- 1075 Reyna-Llorens, I., Burgess, S.J., Reeves, G., Singh, P., Stevenson, S.R., Williams, B.P.,
1076 Stanley, S., and Hibberd, J.M. (2018). Ancient duons may underpin spatial patterning
1077 of gene expression in C₄ leaves. *Proc. Natl. Acad. Sci.* 115:1931-1936.
- 1078 Sage, R. (2004). The evolution of C₄ photosynthesis. *New Phytol.* 161: 341–370.
- 1079 Sage, R.F., Christin, P.-A., and Edwards, E.J. (2011). The C₄ plant lineages of planet
1080 Earth. *J. Exp. Bot.* 62: 3171–3181.
- 1081 Sage, R.F., Sage, T.L., and Kocacinar, F. (2012). Photorespiration and the evolution of C₄
1082 photosynthesis. *Annu. Rev. Plant Biol.* 63: 19–47.
- 1083 Sage, R.F. and Zhu, X.G. (2011). Exploiting the engine of C₄ photosynthesis. *J. Exp. Bot.*
1084 62: 2989–3000.
- 1085 Saldanha, A.J. (2004). Java Treeview--extensible visualization of microarray data.
1086 *Bioinformatics.* 20: 3246-3248.
- 1087 Schäffner, A. R., & Sheen, J. (1991). Maize *rbcS* promoter activity depends on sequence
1088 elements not found in dicot *rbcS* promoters. *The Plant Cell* 3: 997-1012.

1089 Sharwood, R.E., Ghannoum, O., Kapralov, M. V, Gunn, L.H., and Whitney, S.M. (2016).
1090 Temperature responses of Rubisco from Paniceae grasses provide opportunities for
1091 improving C₃ photosynthesis. *Nat. Plants* 2: 16186.
1092 Sheen, J. (1999). C₄ gene expression. *Ann. Rev Plant Physiol. Plant Mol Biol* 50: 187–
1093 217.
1094 Sparks, E.E. et al. (2017). Establishment of Expression in the SHORTROOT-
1095 SCARECROW Transcriptional Cascade through Opposing Activities of Both
1096 Activators and Repressors. *Dev. Cell* 39: 585–596.
1097 Stergachis, A.B., Neph, S., Sandstrom, R., Haugen, E., Reynolds, A.P., Zhang, M., Byron,
1098 R., Canfield, T., Stelhing-Sun, S., Lee, K. and Thurman, R.E. (2014). Conservation of
1099 trans-acting circuitry during mammalian regulatory evolution. *Nature*, 515: 365-370.
1100 Stergachis, A.B., Haugen, E., Shafer, A., Fu, W., Vernot, B., Reynolds, A., Raubitschek,
1101 A., Ziegler, S., LeProust, E.M., Akey, J.M., Stamatoyannopoulos, J.A. (2013). Exonic
1102 transcription factor binding directs codon choice and affects protein evolution. *Science*
1103 342: 1367–72.
1104 Sullivan, A.M., Arsovski, A.A., Lempe, J., Bubb, K.L., Weirauch, M.T., Sabo, P.J.,
1105 Sandstrom, R., Thurman, R.E., Neph, S., Reynolds, A.P. and Stergachis, A.B. (2014).
1106 Mapping and dynamics of regulatory DNA and transcription factor networks in *A.*
1107 *thaliana*. *Cell Reports* 8: 2015-2030.
1108 Taniguchi, M., Izawa, K., Ku, M.S.B., Lin, J.H., Saito, H., Ishida, Y., Ohta, S., Komari, T.,
1109 Matsuoka, M., and Sugiyama, T. (2000). Binding of cell type-specific nuclear proteins
1110 to the 5' flanking region of maize C₄ phosphoenolpyruvate carboxylase gene confers
1111 its differential transcription in mesophyll cells. *Plant Mol. Biol.* 44: 543–557.
1112 Thorvaldsdóttir, H., Robinson, J.T., and Mesirov, J.P. (2013). Integrative Genomics Viewer
1113 (IGV): high-performance genomics data visualization and exploration. *Briefings in*
1114 *Bioinformatics* 14: 178–192.
1115 Thurman, R.E., Rynes, E., Humbert, R., Vierstra, J., Maurano, M.T., Haugen, E., Sheffield,
1116 N.C., Stergachis, A.B., Wang, H., Vernot, B. and Garg, K. (2012). The accessible
1117 chromatin landscape of the human genome. *Nature* 489: 75-82.
1118 Tolbert, N.E. (1971). Microbodies - peroxisomes and glyoxysomes. *Annu. Rev. Plant*
1119 *Physiol.* 22: 45–74.
1120 Tsong, A.E., Tuch, B.B., Li, H., and Johnson, A.D. (2006). Evolution of alternative
1121 transcriptional circuits with identical logic. *Nature* 443: 415–420.
1122 Vollbrecht, E., Springer, P. S., Goh, L., Buckler IV, E. S., & Martienssen, R. (2005).
1123 Architecture of floral branch systems in maize and related grasses. *Nature* 436: 1119-
1124 1126.
1125 Viret, J.F., Mabrouk, Y., and Bogorad, L. (1994). Transcriptional photoregulation of cell-
1126 type preferred expression of maize *RbcS-m3*: 3' and 5' sequences are involved. *Proc.*
1127 *Natl. Acad. Sci.* 91: 8577–8581.
1128 Wickham, H. (2010). ggplot2: elegant graphics for data analysis. *J. Stat. Softw.* 35: 65-68.
1129 Williams, B.P., Burgess, S.J., Reyna-Llorens, I., Knerova, J., Aubry, S., Stanley, S., and
1130 Hibberd, J.M. (2016). An untranslated *cis*-element regulates the accumulation of
1131 multiple C₄ enzymes in *Gynandropsis gynandra* mesophyll cells. *Plant Cell* 28: 454–
1132 465.
1133 Xing, K. and He, X. (2015). Reassessing the “Duon” Hypothesis of Protein Evolution. *Mol.*
1134 *Biol. Evol.* 32: 1056–1062.
1135 Xu, T., Purcell, M., Zucchi, P., Helentjaris, T., and Bogorad, L. (2001). TRM1, a YY1-like
1136 suppressor of *RbcS-m3* expression in maize mesophyll cells. *Proc. Natl. Acad. Sci.*
1137 98: 2295–2300.
1138 Yardimci, G. G., Frank, C. L., Crawford, G. E., & Ohler, U. (2014). Explicit DNase
1139 sequence bias modelling enables high-resolution transcription factor footprint
1140 detection. *Nucleic acids research*, 42: 11865-11878.

1141 Yu, C.P., Chen, S.C.C., Chang, Y.M., Liu, W.Y., Lin, H.H., Lin, J.J., Chen, H.J., Lu, Y.J.,
1142 Wu, Y.H., Lu, M.Y.J. and Lu, C.H. (2015b). Transcriptome dynamics of developing
1143 maize leaves and genome wide prediction of *cis*-elements and their cognate
1144 transcription factors. *Proc. Natl. Acad. Sci.* 112: E2477–E2486.

1145 Zentner, G.E. and Henikoff, S. (2014). High-resolution digital profiling of the epigenome.
1146 *Nat Rev Genet* 15: 814–827.

1147 Zhang, T., Marand, A.P., and Jiang, J. (2016). PlantDHS: a database for DNase I
1148 hypersensitive sites in plants. *Nucleic Acids Res.* 44: D1148–D1153.

1149 Zhang, W., Wu, Y., Schnable, J.C., Zeng, Z., Freeling, M., Crawford, G.E., and Jiang, J.
1150 (2012a). High-resolution mapping of open chromatin in the rice genome. *Genome*
1151 *Res.* 22: 151–162.

1152 Zhang, W., Zhang, T., Wu, Y., and Jiang, J. (2012b). Genome-Wide Identification of
1153 Regulatory DNA Elements and Protein-Binding Footprints Using Signatures of Open
1154 Chromatin in Arabidopsis. *The Plant Cell* 24: 2719–2731.

1155 **Figure Legends:**

1156 **Figure 1: Transcription factor binding atlas for whole leaf samples of four grasses.**

1157 (A) Schematic of phylogenetic relationship between species analysed. The two
1158 independent origins of C₄ photosynthesis are highlighted with black and white circles
1159 (figure not drawn to scale). (B) Summary of sampling and the total number of DNaseI-
1160 hypersensitive sites (DHS) and Digital Genomic Footprints (DGF) identified across all four
1161 species. (C) TreeView diagrams illustrating cut density around individual digital genomic
1162 footprint (DGF). Each row represents an individual DGF, cuts are coloured according to
1163 whether they align to the positive (red) or negative (blue) strand and indicate increased
1164 cutting in a 50 bp window on either side of the DGF. The total number of DGF per sample
1165 is shown at the bottom. (D) Representation of DNaseI-SEQ data from *S. bicolor*, depicting
1166 gene (grey), DHS (light blue), DGF (orange) and DNaseI cut density (dark blue) at three
1167 scales: genome wide, with chromosome number and position indicated (top),
1168 chromosomal (second level) and kilobase genomic region (third level). Between each level
1169 the expanded area is denoted by dashed lines. (E) Pie-chart representing the distribution
1170 of DHS among genomic features. Promoters are defined as sequence up to 2000 base
1171 pairs (bp) upstream of the transcriptional start site, downstream represent regions 1000
1172 downstream the transcription termination site while intergenic represent > 1000 bp
1173 downstream the transcription termination site until the next promoter region. (F) Bar chart
1174 representing number of DHS in genic and promoter regions.

1175

1176 **Figure 2: Digital genomic footprints in whole leaves of four grasses.**

1177 (A) DNA motifs from previous studies in maize (¹Yu et al., 2015; ²Kozaki et al., 2004; ³Niu et al., 2002;
1178 ⁴Vollbrecht et al., 2005; ⁵Eveland et al., 2014; ⁶Li et al., 2015; ⁷Bolduc et al., 2012) were
1179 detected in whole leaves and bundle sheath strands from maize. (B) Pie-chart
1180 representing the distribution of DGF among genomic features. Promoters are defined as
1181 sequence up to 2000 base pairs (bp) upstream of the transcriptional start site, downstream
1182 represent regions 1000 downstream the transcription termination site while intergenic
1183 represent > 1000 bp downstream the transcription termination site until the next promoter
1184 region. (C) Bar chart representing number of DGF in genic and promoter regions. (D)
1185 Polymorphic sites per kb in duons and surrounding exons at FourFold Degenerate Sites
1186 (FFDS) and non-synonymous sites. Chi-squared tests indicate reduced rates of mutation
1187 at FFDS than expected by chance.

1188

1189 **Figure 3: Characterisation of the DNA binding landscape in the C₄ Bundle Sheath.**
1190 (A) Schematic showing that due to their negative distribution below the background signal
1191 derived from reads mapping to the genome, footprints associated with low abundance
1192 cells such as the Bundle Sheath (BS) are unlikely to be detected from whole leaf (WL)
1193 samples. (B) Bundle sheath isolation for DNase I-SEQ experiments, with phylogeny (left)
1194 and workflow (right). (C) DGF identified in the maize *ZmRBCS3* gene coincide with I- and
1195 HOMO-boxes known to regulate gene expression. The gene model is annotated with
1196 whole leaf (blue) and BS (orange) DGF, and the I- and HOMO-boxes are indicated below.
1197 (D) DHS distribution across the maize *PEPC* gene in BS and WL samples. (E) Transcript
1198 abundance expressed as transcripts per million reads (TPM) of three co-linear genes on
1199 chromosome seven of sorghum - C₄ *MDH* (Sobic.007G166300), non C₄ MDH (non C₄,
1200 Sobic.007G166200) and an uncharacterised NAC domain protein (Sobic.007G166100) in
1201 bundle sheath and mesophyll cells. Schematic of these co-linear genes from *S. bicolor*,
1202 depicting three classes of alterations to DNA accessibility and transcription factor binding
1203 to genes that are differentially expressed between whole leaf and bundle sheath cells. (F)
1204 Whole leaf (blue) and bundle sheath (orange) DHS, DGF and differentially (DE) enriched
1205 DGF, as determined by the Wellington Bootstrap algorithm, are depicted. Regions where a
1206 DHS was identified in one sample but not another are indicated by dashed boxes.

1207

1208 **Figure 4: Cistromes associated with cell specific gene expression in C₄ grasses.** (A)
1209 Number of previously reported motifs as well as those defined *de novo* in the grasses. (B)
1210 Density plots depicting average DNaseI activity on positive (red) and negative (blue)
1211 strands centred around a *de novo* motif. (C) Bar chart depicting percentage of DGF
1212 annotated with known or *de novo* motifs. (D) Comparison of transcription factor motif
1213 prevalence in Whole Leaf (WL) samples from *S. italica*, *Z. mays*, *B. distachyon* compared
1214 with *S. bicolor*. Word clouds depict frequency of motifs associated with transcription factor
1215 families, with larger names more abundant. Scatter plots compare frequency of
1216 transcription factor motifs within DGF, ranked from low (most abundant) to high (least
1217 abundant). Correlation between samples is indicated as Kendall's Tau coefficient (τ). (E)
1218 Comparison of transcription factor motif prevalence in BS enriched and WL enriched DGF
1219 from *S. bicolor*, as in (D), word clouds depict frequency of motifs associated with
1220 transcription factor families and plots compare frequency of transcription factor motifs
1221 within DGF ranked from low to high. Similarly scatter plots compare transcription factor
1222 motif prevalence in BS enriched and whole leaf enriched DGF from *S. bicolor*. (F) Venn
1223 diagram showing enriched motifs for each cell type in all three C₄ species.

1224

1225 **Figure 5: *Cis*-elements show high rates of turnover and mobility in grasses.** (A)

1226 Scenarios for DGF conservation between species. Reads derived from DNaseI cuts are
1227 depicted in grey, DGF that are both conserved and occupied between species by red, and
1228 DGF that are conserved but unoccupied by blue shading. (B) Bar-plot representing
1229 pairwise comparisons of DGF occupancy. (C) Schematic depicting the position of a
1230 transcription factor motif consistently associated with the bundle sheath enriched
1231 *TRANSKETOLASE (TKL)* gene in *S. bicolor*, *Z. mays*, *S. italica* and C₃ *B. distachyon*. The
1232 position of motifs conserved between orthologous genes is depicted by solid lines and
1233 orange) and varies between species.

1234

1235 **Figure 6: Hyper-conserved *cis*-elements in grasses recruited into C₄ photosynthesis.**

1236 (A) Conservation of regulation in C₄ and Calvin Benson Bassham cycle genes following
1237 the divergence of *Z. mays* and *S. bicolor*. The number of carbon atoms (red dots) and
1238 metabolite flow (red dashed line) between mesophyll (grey) and bundle sheath (orange)
1239 cells are illustrated along with the degree of conservation of DGF associated with BS
1240 strands. (B) Conservation of DGF occupancy in grasses across evolutionary time. Results
1241 are depicted for whole leaf (WL - blue) and bundle sheath (BS - orange) DGF. The asterisk
1242 indicates 41 DGF that are conserved in the BS of the C₄ species but are also found in
1243 whole leaves of *B. distachyon*). Pie-charts display the distribution of conserved and
1244 occupied DGF for whole leaf and BS strands. Promoters are defined as sequence up to
1245 2000 base pairs (bp) upstream of the transcriptional start site, downstream represent
1246 regions 1000 downstream the transcription termination site while intergenic represent >
1247 1000 bp downstream the transcription termination site until the next promoter region.

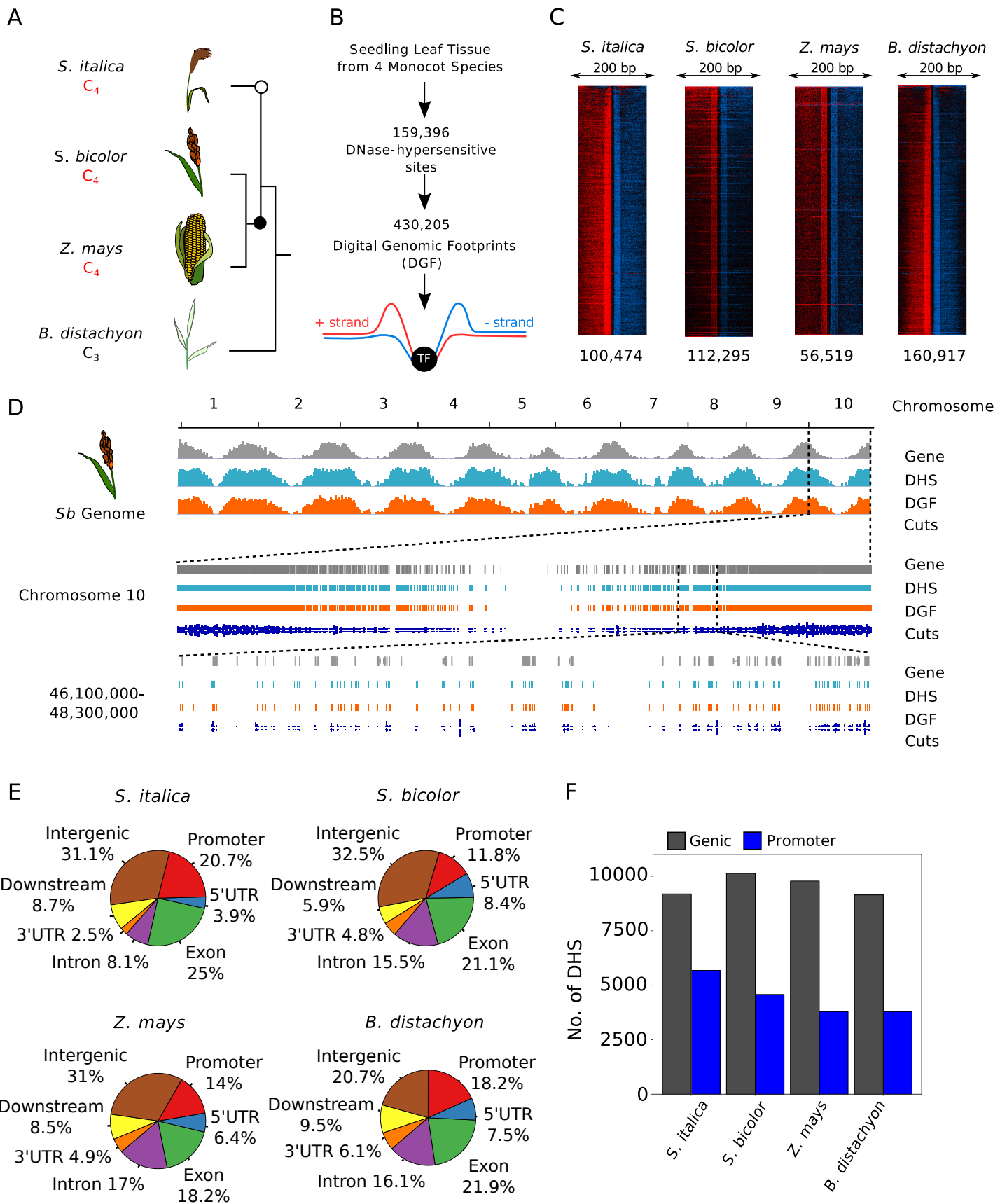


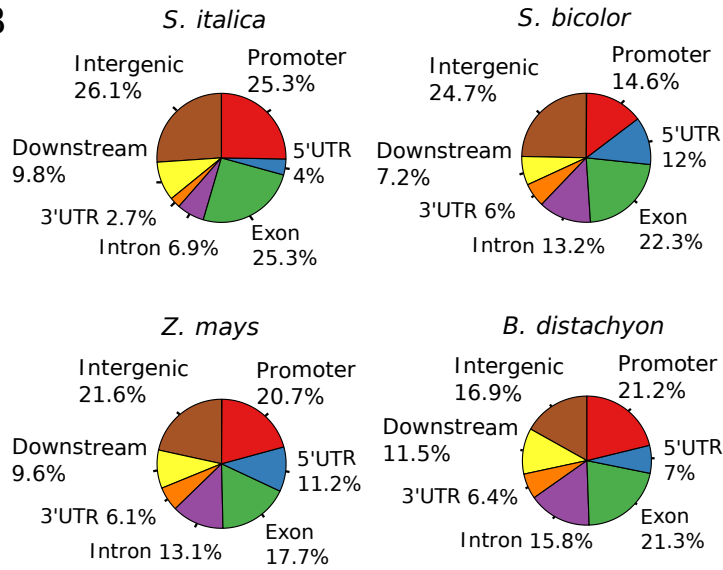
Figure 1

Figure 1: Transcription factor binding atlas for whole leaf samples of four grasses. (A) Schematic of phylogenetic relationship between species analysed. The two independent origins of C4 photosynthesis are highlighted with black and white circles (figure not drawn to scale). (B) Summary of sampling and the total number of DNaseI-hypersensitive sites (DHS) and Digital Genomic Footprints (DGF) identified across all four species. (C) TreeView diagrams illustrating cut density around individual digital genomic footprint (DGF). Each row represents an individual DGF, cuts are coloured according to whether they align to the positive (red) or negative (blue) strand and indicate increased cutting in a 50 bp window on either side of the DGF. The total number of DGF per sample is shown at the bottom. (D) Representation of DNaseI-SEQ data from *S. bicolor*, depicting gene (grey), DHS (light blue), DGF (orange) and DNase-I cut density (dark blue) at three scales: genome wide, with chromosome number and position indicated (top), chromosomal (second level) and kb genomic region (third level). Between each level the expanded area is denoted by dashed lines. (E) Pie-chart representing the distribution of DHS among genomic features. Promoters are defined as sequence up to 2000 base pairs (bp) upstream of the transcriptional start site, downstream represent regions 1000 downstream the transcription termination site while intergenic represent > 1000 bp downstream the transcription termination site until the next promoter region. (F) Bar chart representing number of DHS in genic and promoter regions.

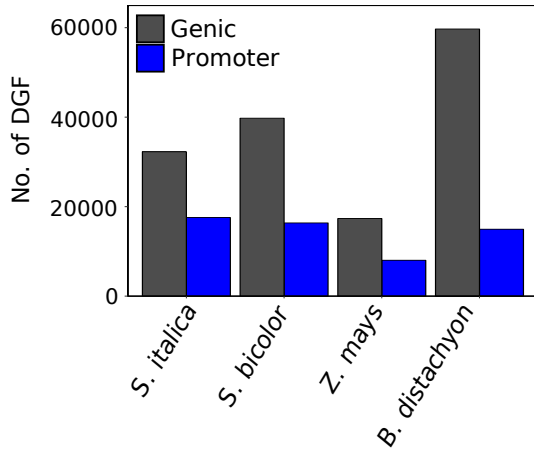
A

Motif ID	Transcription factor	Validation	Number of motif hits in DGF	
			Whole leaf	Bundle sheath
ID0008	EF2	EMSA ¹	359	1024
ID0009	ABI4/ERF11	EMSA ¹	1741	5425
ID0026	CIB1	EMSA ¹	621	1685
ID0037	WRI1	EMSA ¹	3635	13445
ID0061	HOX16	EMSA ¹	118	259
ID0062	RS2	EMSA ¹	261	581
ID0063	GBF6	EMSA ¹	1083	3476
ID0064	GLK1/GLK2	EMSA ¹	140	285
ID0439	MYB59	EMSA ¹	452	897
MA0120.1	ID1	SELEX ²	340	802
MA0123.1	abi4	SELEX ³	2505	9610
MA1416.1	RAMOSA1	ChIP-seq ^{4,5}	705	1253
MA1417.1	O2	ChIP-seq ⁶	300	666
NA	KNT1	ChIP-seq ⁷	2207	5327

B



C



D

Polymorphic sites per kb		Duon	Chi Squared
Surrounding exon	Four Fold Degenerate Sites (FFDS)		
Observed	33.41	30.21	33.52, p=7.04e-9*
Expected	33.16	33.16	
Non-synonymous		Duon	Chi Squared
Observed	20.96		
Expected	20.90	20.90	

Figure 2

Figure 2: Digital genomic footprints in whole leaves of four grasses. (A) DNA motifs from previous studies in maize (1Yu et al., 2015; 2Kozaki et al., 2004; 3Niu et al., 2002; 4Vollbrecht et al., 2005; 5Eveland et al., 2014; 6Li et al., 2015; 7Bolduc et al., 2012) were detected in whole leaves and bundle sheath strands from maize. (B) Pie-chart representing the distribution of DGF among genomic features. Promoters are defined as sequence up to 2000 base pairs (bp) upstream of the transcriptional start site, downstream represent regions 1000 downstream the transcription termination site while intergenic represent > 1000 bp downstream the transcription termination site until the next promoter region. (C) Bar chart representing number of DGF in genic and promoter regions. (D) Polymorphic sites per kb in duons and surrounding exons at FourFold Degenerate Sites (FFDS) and non-synonymous sites. Chi-squared tests indicate reduced rates of mutation at FFDS than expected by chance.

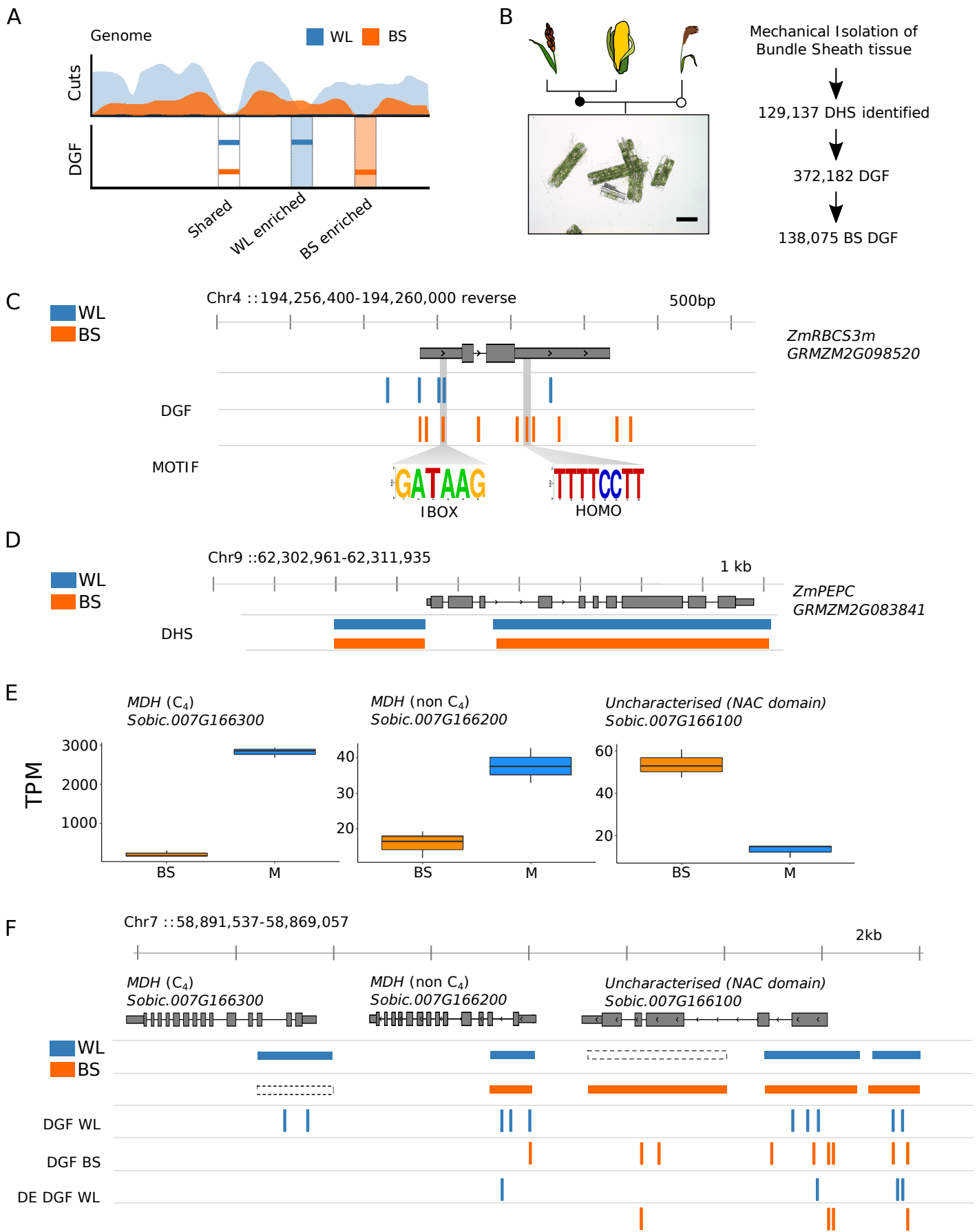


Figure 3

Figure 3: Characterisation of the DNA binding landscape in the C₄ Bundle Sheath. (A) Schematic showing that due to their negative distribution below the background signal derived from reads mapping to the genome, footprints associated with low abundance cells such as the Bundle Sheath (BS) are unlikely to be detected from whole leaf (WL) samples. (B) Bundle sheath isolation for DNaseI-SEQ experiments, with phylogeny (left) and workflow (right). (C) DGF identified in the maize *ZmRBCS3* gene coincide with I- and HOMO-boxes known to regulate gene expression. The gene model is annotated with whole leaf (blue) and BS (orange) DGF, and the I- and HOMO-boxes are indicated below. (D) DHS distribution across the maize PEPC gene in BS and WL samples. (E) Transcript abundance expressed as transcripts per million reads (TPM) of three co-linear genes on chromosome seven of sorghum - C₄ *MDH* (Sobic.007G166300), non C₄ *MDH* (non C₄, Sobic.007G166200) and an uncharacterised NAC domain protein (Sobic.007G166100) in bundle sheath and mesophyll cells. Schematic of these co-linear genes from *S. bicolor*, depicting three classes of alterations to DNA accessibility and transcription factor binding to genes that are differentially expressed between whole leaf and bundle sheath cells. (F) Whole leaf (blue) and bundle sheath (orange) DHS, DGF and differentially (DE) enriched DGF, as determined by the Wellington Bootstrap algorithm, are depicted. Regions where a DHS was identified in one sample but not another are indicated by dashed boxes.

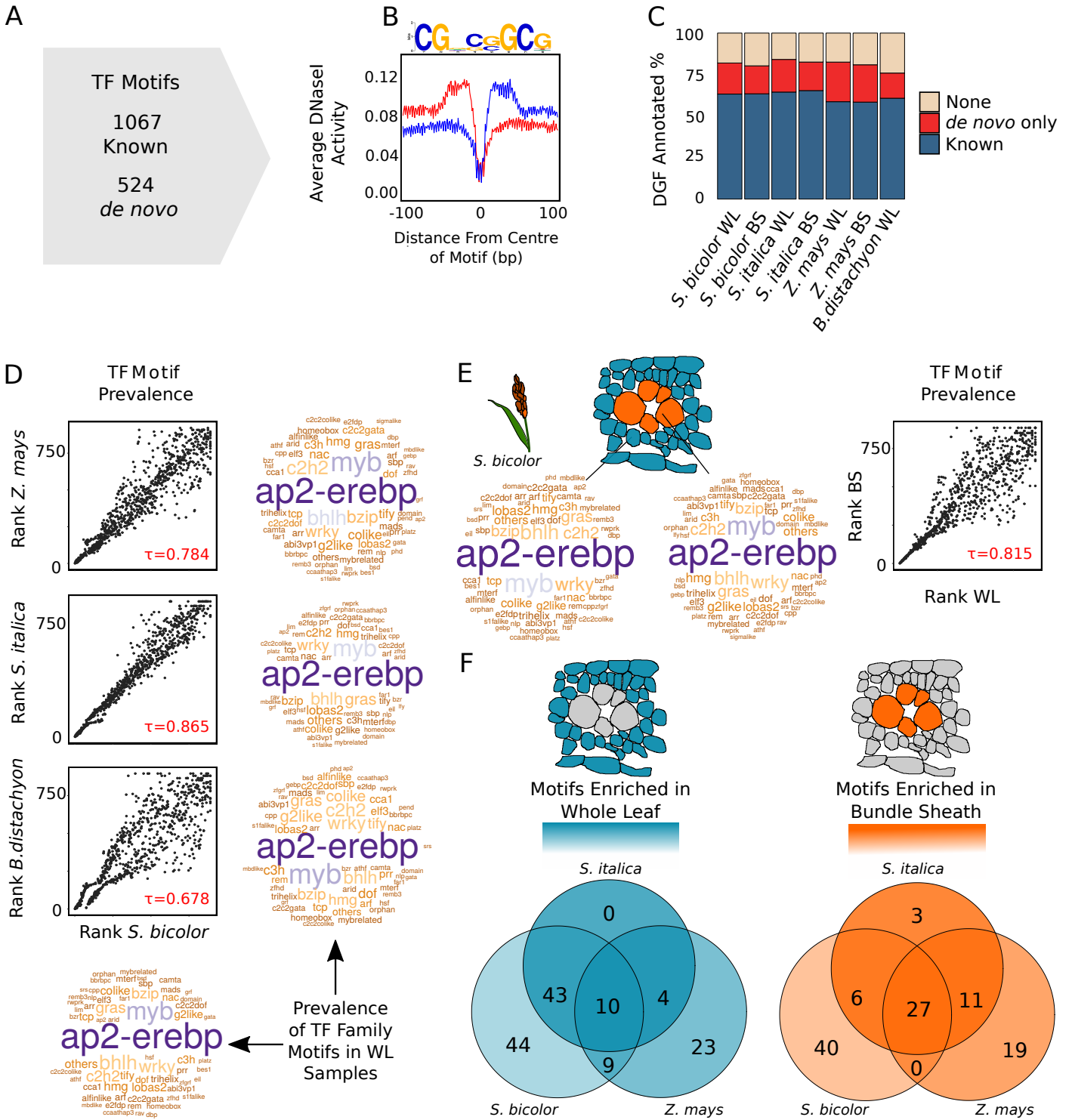


Figure 4: Cistromes associated with cell specific gene expression in C₄ grasses. (A) Number of previously reported motifs as well as those defined *de novo* in the grasses. (B) Density plots depicting average DNase-I activity on positive (red) and negative (blue) strands centred around a *de novo* motif. (C) Bar chart depicting percentage of DGF annotated with known or *de novo* motifs. (D) Comparison of transcription factor motif prevalence in Whole Leaf (WL) samples from *S. italica*, *Z. mays*, *B. distachyon* compared with *S. bicolor*. Word clouds depict frequency of motifs associated with transcription factor families, with larger names more abundant. Scatter plots compare frequency of transcription factor motifs within DGF, ranked from low (most abundant) to high (least abundant). Correlation between samples is indicated as Kendall's Tau coefficient (τ). (E) Comparison of transcription factor motif prevalence in BS enriched and WL enriched DGF from *S. bicolor*, as in (D), word clouds depict frequency of motifs associated with transcription factor families and plots compare frequency of transcription factor motifs within DGF ranked from low to high. Similarly scatter plots compare transcription factor motif prevalence in BS enriched and whole leaf enriched DGF from *S. bicolor*. (F) Venn diagram showing enriched motifs for each cell type in all three C₄ species.

Figure 5: *Cis*-elements show high rates of turnover and mobility in grasses. (A) Scenarios for DGF conservation between species. Reads derived from DNase-I cuts are depicted in grey, DGF that are both conserved and occupied between species by red, and DGF that are conserved but unoccupied by blue shading. (B) Bar-plot representing pairwise comparisons of DGF occupancy. (C) Schematic depicting the position of transcription factor motif consistently associated with the bundle sheath enriched *TRANSKETOLASE (TKL)* gene in *S. bicolor*, *Z. mays*, *S. italica* and C₃ *B. distachyon*. The position of motifs conserved between orthologous genes is depicted by solid lines and varies between species.

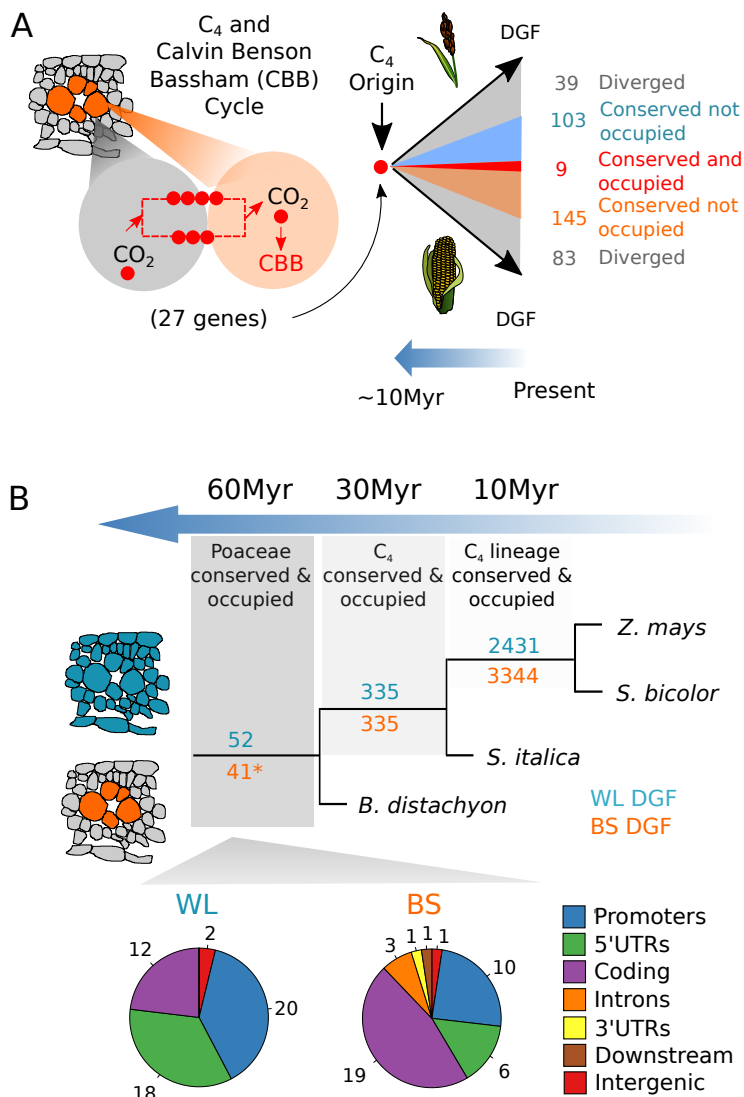


Figure 6: Hyper-conserved *cis*-elements in grasses recruited into C₄ photosynthesis. (A) Conservation of regulation in C₄ and Calvin Benson Bassham cycle genes following the divergence of *Z. mays* and *S. bicolor*. The number of carbon atoms (red dots) and metabolite flow (red dashed line) between mesophyll (grey) and bundle sheath (orange) cells are illustrated along with the degree of conservation of DGF associated with BS strands. (B) Conservation of DGF occupancy in grasses across evolutionary time. Results are depicted for whole leaf (WL - blue) and bundle sheath (BS - orange) DGF. The asterisk indicates 41 DGF that are conserved in the BS of the C₄ species but are also found in whole leaves of *B. distachyon*. Pie-charts display the distribution of conserved and occupied DGF for whole leaf and BS strands. Promoters are defined as sequence up to 2000 base pairs (bp) upstream of the transcriptional start site, downstream represent regions 1000 downstream the transcription termination site while intergenic represent > 1000 bp downstream the transcription termination site until the next promoter region. (D) Bar chart representing DGF number in genic versus promoter regions.

Parsed Citations

Bauwe, H., Hagemann, M., and Fernie, A.R. (2010). Photorespiration: players, partners and origin. *Trends Plant Sci.* 15: 330–6.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Bolduc, N. et al. (2012a). A Genome-Wide Regulatory Framework Identifies Maize Pericarp Color1 Controlled Genes. *Plant Cell* 27: 543–553.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Bolduc, N. and Hake, S. (2009). The Maize Transcription Factor KNOTTED1 Directly Regulates the Gibberellin Catabolism Gene *ga2ox1*. *Plant Cell* 21: 1647–1658.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Bolduc, N., Yilmaz, A., Mejia-Guerra, M.K., Morohashi, K., O'Connor, D., Grotewold, E., and Hake, S. (2012b). Unraveling the KNOTTED1 regulatory network in maize meristems. *Genes Dev.* 26: 1685–1690.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Borba, A R., Serra, T. S., Górska, A., Gouveia, P., Cordeiro, A. M., Reyna-Llorens, I., Kneřová, J., Barros, P.M., Abreu, I.A., Oliveira, Hibberd, J.M., Saibo, N. (2018). Synergistic binding of bHLH transcription factors to the promoter of the maize NADP-ME gene used in C4 photosynthesis is based on an ancient code found in the ancestral C3 state. *Mol. Bio and Evolution* 35: 1690-1705.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Bowes, G., Ogren, W.L., and Hageman, R.H. (1971). Phosphoglycolate production catalyzed by ribulose diphosphate carboxylase. *Biochem. Biophys. Res. Commun.* 45: 716–722.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Brown, N.J., Newell, C.A., Stanley, S., Chen, J.E., Perrin, A.J., Kajala, K., and Hibberd, J.M. (2011). Independent and Parallel Recruitment of Preexisting Mechanisms Underlying C4 Photosynthesis. *Science* 331: 1436–1439.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Bukowski, R., Guo, X., Lu, Y., Zou, C., He, B., Rong, Z., Wang, B., Xu, D., Yang, B., Xie, C. and Fan, L. (2017). Construction of the third-generation *Zea mays* haplotype map. *GigaScience* 7: gix134.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Chang YM, Liu WY, Shih ACC, Shen MN, Lu CH, Lu MYJ, Yang HW, Wang TY, Chen SC, Chen SM, et al. (2012) Characterizing regulatory and functional differentiation between maize mesophyll and bundle sheath cells by transcriptomic analysis. *Plant Physiol* 160: 165–177.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Christin, P.A., Salamin, N., Savolainen, V., Duvall, M.R., and Besnard, G. (2007). C4 photosynthesis evolved in grasses via parallel adaptive genetic changes. *Curr. Biol.* 17: 1241–1247.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Christin, P. A., & Osborne, C. P. (2014). The evolutionary ecology of C4 plants. *New Phytologist* 204: 765-781.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Cingolani, P., Platts, A., Wang, L.L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X. and Ruden, D.M. (2012). A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly* 6: 80-92.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Cornejo, M.-J., Luth, D., Blankenship, K.M., Anderson, O.D., and Blechl, A.E. (1993). Activity of a maize ubiquitin promoter in transgenic rice. *Plant Mol. Biol.* 23: 567–581.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Covshoff, S., Furbank, R.T., Leegood, R.C., and Hibberd, J.M. (2013). Leaf rolling allows quantification of mRNA abundance in mesophyll cells of sorghum. *J. Exp. Bot.* 64: 807–813.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Deal, R.B. and Henikoff, S. (2011). The INTACT method for cell type-specific gene expression and chromatin profiling in *Arabidopsis*

thaliana. Nat. Protoc. 6: 56–68.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Denas, O., Sandstrom, R., Cheng, Y., Beal, K., Herrero, J., Hardison, R.C., and Taylor, J. (2015). Genome-wide comparative analysis reveals human-mouse regulatory landscape and evolution. BMC Genomics 16: 87.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Emms, D. M., Covshoff, S., Hibberd, J. M., & Kelly, S. (2016). Independent and parallel evolution of new genes by gene duplication in two origins of C4 photosynthesis provides new insight into the mechanism of phloem loading in C4 species. Mol Bio and Evol, 33(7), 1796-1806.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Eveland, A.L., Goldshmidt, A., Pautler, M., Morohashi, K., Liseron-Monfils, C., Lewis, M.W., Kumari, S., Hiraga, S., Yang, F., Unger-Wallace, E. and Olson, A. (2014). Regulatory modules controlling maize inflorescence architecture. Genome Research 24: 431-443.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Fellows, I. (2012). wordcloud: Word clouds. R Packag. version 2: 109.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Feng, J., Liu, T., Qin, B., Zhang, Y., and Liu, X.S. (2012). Identifying ChIP-SEQ enrichment using MACS. Nat. Protoc. 7: 1728–1740.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Furbank, R.T. (2011). Evolution of the C4 photosynthetic mechanism: are there really three C4 acid decarboxylation types? J. Exp. Bot. 62: 3103–3108.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Furbank, R.T., Stitt, M., and Foyer, C.H. (1985). Intercellular compartmentation of sucrose synthesis in leaves of Zea mays L. Planta 164: 172–178.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Gel, B., Díez-Villanueva, A., Serra, E., Buschbeck, M., Peinado, M. A., & Malinverni, R. (2015). regioneR: an R/Bioconductor package for the association analysis of genomic regions based on permutation tests. Bioinformatics 32: 289-291.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Giuliano, G., Pichersky, E., Malik, V.S., Timko, M.P., Scolnik, P.A., and Cashmore, A.R. (1988). An evolutionarily conserved protein binding sequence upstream of a plant light regulated gene. Proc. Natl. Acad. Sci. 85: 7089–7093.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Goodstein, D.M., Shu, S., Howson, R., Neupane, R., Hayes, R.D., Fazo, J., Mitros, T., Dirks, W., Hellsten, U., Putnam, N., and Rokhsar, D.S. (2012). Phytozome: a comparative platform for green plant genomics. Nucleic Acids Res. 40: D1178-86.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Gowik, U., Burscheidt, J., Akyildiz, M., Schlue, U., Koczor, M., Streubel, M., and Westhoff, P. (2004). cis-regulatory elements for mesophyll-specific gene expression in the C4 plant Flaveria trinervia, the promoter of the C4 PHOSPHOENOLPYRUVATE CARBOXYLASE gene. Plant Cell 16: 1077–1090.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Grant, C.E., Bailey, T.L., and Noble, W.S. (2011). FIMO: scanning for occurrences of a given motif. Bioinformatics 27: 1017–1018.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Gu, Z., Gu, L., Eils, R., Schlesner, M., & Brors, B. (2014). circlize implements and enhances circular visualization in R. Bioinformatics 30: 2811-2812.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Guy, L., Roat Kultima, J., and Andersson, S.G.E. (2010). genoPlotR: comparative gene and genome visualization in R. Bioinformatics 26: 2334–2335.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Harris, R.S. (2007). Improved pairwise alignment of genomic DNA. The Pennsylvania State University. (PhD Thesis).

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Hatch, M.D. (1987). C4 photosynthesis: a unique elend of modified biochemistry, anatomy and ultrastructure. *Biochim. Biophys. Acta - Rev. Bioenerg.* 895: 81–106.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

He, H.H., Meyer, C.A, Chen, M.W., Zang, C., Liu, Y., Rao, P.K., Fei, T., Xu, H., Long, H., Liu, X.S. and Brown, M. (2014). Refined DNase-seq protocol and data analysis reveals intrinsic bias in transcription factor footprint identification. *Nature Methods* 11: 73.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y.C., and Laslo, P. (2010). Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol. Cell.* 38: 576-89

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Hesselberth, J.R., Chen, X., Zhang, Z., Sabo, P.J., Sandstrom, R., Reynolds, A.P., Thurman, R.E., Neph, S., Kuehn, M.S., Noble, W.S., Fields, S., and Stamatoyannopoulos, J.A (2009). Global mapping of protein-DNA interactions in vivo by digital genomic footprinting. *Nat. Methods* 6: 283–9.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Hibberd, J.M. and Covshoff, S. (2010). The Regulation of Gene Expression Required for C4 Photosynthesis. *Annu. Rev. Plant Biol.* 61: 181–207.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Huang, W., Loganantharaj, R., Schroeder, B., Fargo, D., & Li, L. (2013). PAMS: a tool for Peak Annotation and Visualization. *Bioinformatics* 29: 3097-3099.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Jeon, J.-S., Lee, S., Jung, K.-H., Jun, S.-H., Kim, C., and An, G. (2000). Tissue-Preferential Expression of a Rice α -Tubulin Gene, OsTubA1, Mediated by the First Intron. *Plant Physiol.* 123: 1005–1014.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

John, S., Sabo, P.J., Thurman, R.E., Sung, M.-H., Biddie, S.C., Johnson, T.A, Hager, G.L., and Stamatoyannopoulos, J.A (2011). Chromatin accessibility pre-determines glucocorticoid receptor binding patterns. *Nat Genet.* 43: 264–268.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

John, C. R., Smith-Unna, R. D., Woodfield, H., Covshoff, S., & Hibberd, J. M. (2014). Evolutionary convergence of cell-specific gene expression in independent lineages of C4 grasses. *Plant Phys.*, 165(1), 62-75.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Jordan, D.B. and Ogren, W.L. (1984). The CO₂/O₂ specificity of Ribulose 1,5-Bisphosphate Carboxylase Oxygenase - dependence on Ribulose bisphosphate concentration, pH and temperature. *Planta* 161: 308–313.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Kajala, K., Williams, B.P., Brown, N.J., Taylor, L.E., and Hibberd, J.M. (2011). Multiple Arabidopsis genes primed for direct recruitment into C4 photosynthesis. *Plant J.* 69: 47–56.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and Haussler, and D. (2002). The Human Genome Browser at UCSC. *Genome Res.* 12: 996–1006.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Khan, A, Fornes, O., Stigliani, A, Gheorghe, M., Castro-Mondragon, J.A, van der Lee, R., Bessy, A, Cheneby, J., Kulkarni, S.R., Tan, G. and Baranasic, D., (2017). JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Research* 46 (D1): D260-D266.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Kharchenko, P. V, Tolstorukov, M.Y., and Park, P.J. (2008). Design and analysis of ChIP-SEQ experiments for DNA-binding proteins. *Nat Biotech.* 26: 1351–1359.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Kozaki, A, Hake, S., & Colasanti, J. (2004). The maize ID1 flowering time regulator is a zinc finger protein with novel DNA binding

properties. *Nucleic Acids Research*, 32: 1710-1720.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Langmead, B. and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods* 9: 357–359.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Leegood, R.C. (1985). The intercellular compartmentation of metabolites in leaves of *Zea mays* L. *Planta* 164: 163–171.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

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.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Li, C., Qiao, Z., Qi, W., Wang, Q., Yuan, Y., Yang, X., Tang, Y., Mei, B., Lv, Y., Zhao, H. and Xiao, H. (2015). Genome-wide characterization of cis-acting DNA targets reveals the transcriptional regulatory framework of *opaque2* in maize. *The Plant Cell* 27: 532–545.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., and Subgroup, 1000 Genome Project Data Processing (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25: 2078–2079.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

de Lucas, M., Pu, L., Turco, G., Gaudinier, A., Morao, A.K., Harashima, H., Kim, D., Ron, M., Sugimoto, K., Roudier, F., and Brady, S.M. (2016). Transcriptional Regulation of Arabidopsis Polycomb Repressive Complex 2 Coordinates Cell-Type Proliferation and Differentiation. *Plant Cell* 28: 2616 -2631.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Maas, C., Laufs, J., Grant, S., Korfhage, C., and Werr, W. (1991). The combination of a novel stimulatory element in the first exon of the maize *SHRUNKEN-1* gene with the following intron 1 enhances reporter gene expression up to 1000-fold. *Plant Mol. Biol.* 16: 199–207.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Maher, K.A., Bajic, M., Kajala, K., Reynoso, M., Pauluzzi, G., West, D.A., Zumstein, K., Woodhouse, M., Bubb, K., Dorrity, M.W. and Queitsch, C. (2018). Profiling of accessible chromatin regions across multiple plant species and cell types reveals common gene regulatory principles and new control modules. *The Plant Cell* 30: 15-36.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Marinov, G.K., Kundaje, A., Park, P.J., and Wold, B.J. (2014). Large-Scale Quality Analysis of Published ChIP-SEQ Data. *G3* 4: 209–223.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Markelz, N. H., Costich, D. E., & Brutnell, T. P. (2003). Photomorphogenic responses in maize seedling development. *Plant Physiology*, 133(4), 1578-1591.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Martin, A and Orgogozo, V. (2013). The Loci of repeated evolution: a catalog of genetic hotspots of phenotypic variation. *Evolution* 67: 1235–50.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Natarajan, A, Yardimci, G.G., Sheffield, N.C., Crawford, G.E., and Ohler, U. (2012). Predicting cell-type-specific gene expression from regions of open chromatin. *Genome Research* 22: 1711–1722.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Neph, S., Vierstra, J., Stergachis, A.B., Reynolds, A.P., Haugen, E., Vernot, B., Thurman, R.E., John, S., Sandstrom, R., Johnson, A.K. and Maurano, M.T. (2012). An expansive human regulatory lexicon encoded in transcription factor footprints. *Nature* 489: 83-90.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

Niu, X., Helentjaris, T., & Bate, N. J. (2002). Maize *ABI4* binds coupling element1 in abscisic acid and sugar response genes. *The Plant Cell*, 14: 2565-2575.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only Title Only Author and Title](#)

O'Malley, R.C., Huang, S.C., Song, L., Lewsey, M.G., Bartlett, A., Nery, J.R., Galli, M., Gallavotti, A., and Ecker, J.R. (2016). Cistrome and Epicistrome Features Shape the Regulatory DNA Landscape. *Cell* 165: 1280–1292.

- Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Pajoro, A., Madrigal, P., Muiño, J.M., Matus, J.T., Jin, J., Mecchia, M.A., Debernardi, J.M., Palatnik, J.F., Balazadeh, S., Arif, M. and O'Maoiléidigh, D.S. (2014). Dynamics of chromatin accessibility and gene regulation by MADS-domain transcription factors in flower development. *Genome Biology* 15: R41.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Patel, M., Corey, A.C., Yin, L.P., Ali, S.J., Taylor, W.C., and Berry, J.O. (2004). Untranslated regions from C4 Amaranth AhRbcS1 mRNAs confer translational enhancement and preferential bundle sheath cell expression in transgenic C4 *Flaveria bidentis*. *Plant Physiol.* 136: 3550–3561.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Pautler, M., Eveland, A.L., LaRue, T., Yang, F., Weeks, R., Lunde, C., Je, B. II, Meeley, R., Komatsu, M., Vollbrecht, E., Sakai, H., and Jackson, D. (2015). FASCIATED EAR4 Encodes a bZIP Transcription Factor That Regulates Shoot Meristem Size in Maize. *Plant Cell* 27: 104-120.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Piper, J., Assi, S.A., Cauchy, P., Ladroue, C., Cockerill, P.N., Bonifer, C., and Ott, S. (2015). Wellington-bootstrap: differential DNase-seq footprinting identifies cell-type determining transcription factors. *BMC Genomics* 16: 1000.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Piper, J., Elze, M.C., Cauchy, P., Cockerill, P.N., Bonifer, C., and Ott, S. (2013). Wellington: a novel method for the accurate identification of digital genomic footprints from DNase-seq data. *Nucleic Acids Research* 41: e201–e201.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Quinlan, A.R. and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26: 841–842.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Reyna-Llorens, I., Burgess, S.J., Reeves, G., Singh, P., Stevenson, S.R., Williams, B.P., Stanley, S., and Hibberd, J.M. (2018). Ancient duons may underpin spatial patterning of gene expression in C4 leaves. *Proc. Natl. Acad. Sci.* 115:1931-1936.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Sage, R. (2004). The evolution of C4 photosynthesis. *New Phytol.* 161: 341–370.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Sage, R.F., Christin, P.-A., and Edwards, E.J. (2011). The C4 plant lineages of planet Earth. *J. Exp. Bot.* 62: 3171–3181.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Sage, R.F., Sage, T.L., and Kocacinar, F. (2012). Photorespiration and the evolution of C4 photosynthesis. *Annu. Rev. Plant Biol.* 63: 19–47.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Sage, R.F. and Zhu, X.G. (2011). Exploiting the engine of C4 photosynthesis. *J. Exp. Bot.* 62: 2989–3000.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Saldanha, A.J. (2004). Java Treeview—extensible visualization of microarray data. *Bioinformatics.* 20: 3246-3248.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Schäffner, A R., & Sheen, J. (1991). Maize rbcS promoter activity depends on sequence elements not found in dicot rbcS promoters. *The Plant Cell* 3: 997-1012.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Sharwood, R.E., Ghannoum, O., Kapralov, M. V, Gunn, L.H., and Whitney, S.M. (2016). Temperature responses of Rubisco from Paniceae grasses provide opportunities for improving C3 photosynthesis. *Nat. Plants* 2: 16186.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)
- Sheen, J. (1999). C4 gene expression. *Ann. Rev Plant Physiol. Plant Mol Biol* 50: 187–217.**
Pubmed: [Author and Title](#)
Google Scholar: [Author Only Title Only Author and Title](#)

Sparks, E.E. et al. (2017). Establishment of Expression in the SHORTROOT-SCARECROW Transcriptional Cascade through Opposing Activities of Both Activators and Repressors. Dev. Cell 39: 585–596.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Stergachis, A.B., Neph, S., Sandstrom, R., Haugen, E., Reynolds, A.P., Zhang, M., Byron, R., Canfield, T., Stelhing-Sun, S., Lee, K. and Thurman, R.E. (2014). Conservation of trans-acting circuitry during mammalian regulatory evolution. Nature, 515: 365-370.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Stergachis, A.B., Haugen, E., Shafer, A, Fu, W., Vernot, B., Reynolds, A, Raubitschek, A, Ziegler, S., LeProust, E.M., Akey, J.M., Stamatoyannopoulos, J.A. (2013). Exonic transcription factor binding directs codon choice and affects protein evolution. Science 342: 1367–72.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Sullivan, A.M., Arsovski, A.A, Lempe, J., Bubb, K.L., Weirauch, M.T., Sabo, P.J., Sandstrom, R., Thurman, R.E., Neph, S., Reynolds, A.P. and Stergachis, A.B. (2014). Mapping and dynamics of regulatory DNA and transcription factor networks in *A. thaliana*. Cell Reports 8: 2015-2030.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Taniguchi, M., Izawa, K., Ku, M.S.B., Lin, J.H., Saito, H., Ishida, Y., Ohta, S., Komari, T., Matsuoka, M., and Sugiyama, T. (2000). Binding of cell type-specific nuclear proteins to the 5' flanking region of maize C4 phosphoenolpyruvate carboxylase gene confers its differential transcription in mesophyll cells. Plant Mol. Biol. 44: 543–557.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Thorvaldsdóttir, H., Robinson, J.T., and Mesirov, J.P. (2013). Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration. Briefings in Bioinformatics 14: 178–192.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Thurman, R.E., Rynes, E., Humbert, R., Vierstra, J., Maurano, M.T., Haugen, E., Sheffield, N.C., Stergachis, A.B., Wang, H., Vernot, B. and Garg, K. (2012). The accessible chromatin landscape of the human genome. Nature 489: 75-82.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Tolbert, N.E. (1971). Microbodies - peroxisomes and glyoxysomes. Annu. Rev. Plant Physiol. 22: 45–74.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Tsong, A.E., Tuch, B.B., Li, H., and Johnson, A.D. (2006). Evolution of alternative transcriptional circuits with identical logic. Nature 443: 415–420.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Vollbrecht, E., Springer, P. S., Goh, L., Buckler IV, E. S., & Martienssen, R. (2005). Architecture of floral branch systems in maize and related grasses. Nature 436: 1119-1126.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Viret, J.F., Mabrouk, Y., and Bogorad, L. (1994). Transcriptional photoregulation of cell-type preferred expression of maize RbcS-m3: 3' and 5' sequences are involved. Proc. Natl. Acad. Sci. 91: 8577–8581.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Wickham, H. (2010). ggplot2: elegant graphics for data analysis. J. Stat. Softw. 35: 65-68.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Williams, B.P., Burgess, S.J., Reyna-Llorens, I., Knerova, J., Aubry, S., Stanley, S., and Hibberd, J.M. (2016). An untranslated cis-element regulates the accumulation of multiple C4 enzymes in Gynandropsis gynandra mesophyll cells. Plant Cell 28: 454–465.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Xing, K. and He, X. (2015). Reassessing the "Duon" Hypothesis of Protein Evolution. Mol. Biol. Evol. 32: 1056–1062.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Xu, T., Purcell, M., Zucchi, P., Helentjaris, T., and Bogorad, L. (2001). TRM1, a YY1-like suppressor of RbcS-m3 expression in maize mesophyll cells. Proc. Natl. Acad. Sci. 98: 2295–2300.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Yardimci, G. G., Frank, C. L., Crawford, G. E., & Ohler, U. (2014). Explicit DNase sequence bias modelling enables high-resolution transcription factor footprint detection. *Nucleic acids research*, 42: 11865-11878.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Yu, C.P., Chen, S.C.C., Chang, Y.M., Liu, W.Y., Lin, H.H., Lin, J.J., Chen, H.J., Lu, Y.J., Wu, Y.H., Lu, M.Y.J. and Lu, C.H. (2015b). Transcriptome dynamics of developing maize leaves and genome wide prediction of cis-elements and their cognate transcription factors. *Proc. Natl. Acad. Sci.* 112: E2477–E2486.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Zentner, G.E. and Henikoff, S. (2014). High-resolution digital profiling of the epigenome. *Nat Rev Genet* 15: 814–827.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Zhang, T., Marand, A.P., and Jiang, J. (2016). PlantDHS: a database for DNase I hypersensitive sites in plants. *Nucleic Acids Res.* 44: D1148–D1153.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Zhang, W., Wu, Y., Schnable, J.C., Zeng, Z., Freeling, M., Crawford, G.E., and Jiang, J. (2012a). High-resolution mapping of open chromatin in the rice genome. *Genome Res.* 22: 151–162.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)

Zhang, W., Zhang, T., Wu, Y., and Jiang, J. (2012b). Genome-Wide Identification of Regulatory DNA Elements and Protein-Binding Footprints Using Signatures of Open Chromatin in Arabidopsis. *The Plant Cell* 24: 2719–2731.

Pubmed: [Author and Title](#)

Google Scholar: [Author Only](#) [Title Only](#) [Author and Title](#)