

# IsoSeq transcriptome assembly of C<sub>3</sub> panicoid grasses provides tools to study evolutionary change in the Panicoideae

Daniel S. Carvalho<sup>1</sup>, Aime V. Nishimwe<sup>1</sup> & James C. Schnable<sup>1</sup>

1. Center for Plant Science Innovation, Department of Agronomy and Horticulture, University of Nebraska-Lincoln, Lincoln, NE, USA.

7 Abstract

8 The number of plant species with genomic and transcriptomic data has been increasing rapidly. The grasses – Poaceae – have been well represented among species with published reference genomes. However, as a result the genomes of wild grasses are less frequently targeted by sequencing efforts. Sequence data from wild relatives of crop species in the grasses can aid the study of domestication, gene discovery for breeding and crop improvement, and improve our understanding of the evolution of C<sub>4</sub> photosynthesis. Here we used long read sequencing technology to characterize the transcriptomes of three C<sub>3</sub> panicoid grass species: *Dichanthelium oligosanthes*, *Chasmanthium laxum*, and *Hymenachne amplexicaulis*. Based on alignments to the sorghum genome we estimate that assembled consensus transcripts from each species capture between 54.2 and 65.7% of the conserved syntenic gene space in grasses. Genes co-opted into C<sub>4</sub> were also well represented in this dataset, despite concerns that, because these genes might play roles unrelated to photosynthesis in the target species, they would be expressed at low levels and missed by transcript-based sequencing. A combined analysis using syntenic orthologous genes from grasses with published reference genomes and consensus long read sequences from these wild species was consistent with previously published phylogenies. It is hoped that this data, targeting under represented classes of species within the PACMAD grasses – wild species and species utilizing C<sub>3</sub> photosynthesis – will aid in future studies of domestication and C<sub>4</sub> evolution by decreasing the evolutionary distance between C<sub>4</sub> and C<sub>3</sub> species within this clade, enabling more accurate comparisons associated with evolution of the C<sub>4</sub> pathway.

## 28 Introduction

29 The pace of plant genome sequencing has accelerated in recent years. However, despite decreases in sequencing costs and improvements in genome assembly quality, species selected for whole genome sequencing often meet one or more of the following criteria: A) agricultural importance, B) status as a genetic model system or C) ecological importance. Sequence data from species which lack direct economic, ecological, or genetic model importance can enable comparative analyses to address biological questions in crops and model species (Michael and Jackson, 2013; Ellegren, 2014). C<sub>4</sub> photosynthesis has evolved multiple times in the grasses (GPWG II, 2012), making it particularly amenable to study through comparative genetic approaches (Wang et al., 2009; Huang et al., 2016). C<sub>4</sub> photosynthesis requires both substantial biochemical and anatomical changes (Kellogg, 2013). All grasses which utilize the C<sub>4</sub> pathway belong to the PACMAD clade, a group of grass subfamilies and tribes which includes substantial numbers of both C<sub>3</sub> and C<sub>4</sub> species (GPWG II, 2012). Substantial new insights into both the genes involved in producing the biochemical and anatomical changes required for C<sub>4</sub> photosynthesis, as well as the potential function of individual amino acid residues can be obtained from comparative analysis of individual gene families across species utilizing either C<sub>3</sub> or C<sub>4</sub> photosynthesis within the PACMAD clade (Christin et al., 2007, 2015; Moreno-Villena et al., 2017). However, assembling sequence data for a single gene family from a large enough set of species through PCR amplification and individual Sanger sequencing remains a time and labor intensive process.

47 Many domesticated grasses belong to the PACMAD clade, including such as maize (*Zea mays*), sugar cane (*Saccharum spp.*), sorghum (*Sorghum bicolor*), and foxtail millet (*Setaria italica*). However every domesticated grass in the PACMAD clade with a sequenced genome utilizes one or more variants of the C<sub>4</sub> photosynthetic pathway (Schnable et al., 2009; Garsmeur et al., 2018; Paterson et al., 2009; Bennetzen et al., 2012). As a result, while published whole genome sequence assemblies exist for at least 14 grasses within the PACMAD clade (Table 1), only one of these (*Dichanthelium oligosanthes*, a wild species) (Studer et al., 2016) utilizes C<sub>3</sub> photosynthesis. Long-read sequencing can effectively generate sequence for large numbers of full length cDNAs even in species lacking reference genome assemblies (An et al., 2018; Zhang et al., 2019). One concern with utilizing this technology for

comparative genetic studies is that the higher error rate, particularly the frequencies<sup>57</sup> of insertion and deletion errors, make data from long read based sequencing of non-model species<sup>58</sup> unsuitable for use in comparative evolutionary analyses (Gonzalez-Garay, 2016). However, we previously<sup>59</sup> found that observed synonymous substitution rates calculated from consensus sequences<sup>60</sup> constructed using PacBio IsoSeq pipeline were not elevated relative to a sister lineage where gene<sup>61</sup> sequences were taken from a sanger-based whole genome assembly, indicating sequence data obtained<sup>62</sup> in this manner may indeed be suitable for comparative evolutionary analyses (Yan et al., 2019).

Species	Relevance	C <sub>3</sub> /C <sub>4</sub>	Genome Publication
<i>Dichanthelium oligosanthes</i>	Wild Species	C <sub>3</sub>	Studer et al. (2016)
<i>Eleusine coracana</i> <sup>a</sup>	Grain Crop	C <sub>4</sub>	Hittalmani et al. (2017)
<i>Eragrostis tef</i> <sup>a</sup>	Grain Crop	C <sub>4</sub>	Cannarozzi et al. (2014); VanBuren et al. (2019)
<i>Miscanthus x giganteus</i> <sup>b</sup>	Biomass Crop	C <sub>4</sub>	Swaminathan et al. (2010) VanBuren et al. (2015, 2018)
<i>Oropetium thomaeum</i> <sup>a</sup>	Wild Species	C <sub>4</sub>	Lovell et al. (2018)
<i>Panicum hallii</i> <sup>c</sup>	Wild Species	C <sub>4</sub>	Zou et al. (2019)
<i>Panicum miliaceum</i> <sup>c</sup>	Grain Crop	C <sub>4</sub>	Casler et al. (2011)
<i>Panicum virgatum</i> <sup>c</sup>	Biomass Crop	C <sub>4</sub>	Varshney et al. (2017b)
<i>Pennisetum glaucum</i> <sup>c</sup>	Grain Crop	C <sub>4</sub>	Garsmeur et al. (2018)
<i>Saccharum spp.</i> <sup>b</sup>	Sugar Crop	C <sub>4</sub>	Bennetzen et al. (2012)
<i>Setaria italica</i> <sup>c</sup>	Grain Crop	C <sub>4</sub>	Brutnell et al. (2010)
<i>Setaria viridis</i> <sup>c</sup>	Genetic Model	C <sub>4</sub>	Paterson et al. (2009)
<i>Sorghum bicolor</i> <sup>b</sup>	Grain/Biomass/ Sugar Crop	C <sub>4</sub>	Schnable et al. (2009)
<i>Zea mays</i> <sup>b</sup>	Grain Crop & Genetic Model	C <sub>4</sub>	

Table 1: Published reference genomes for grass species within the PACMAD clade. Species sharing a common inferred evolutionary origin of C<sub>4</sub> photosynthesis as reported in (?) are indicated by superscript letters.

Here we report the sequencing and characterization of IsoSeq based transcriptomes for three<sup>65</sup> additional PACMAD grasses, selecting to enable wider scale studies of protein sequence changes<sup>66</sup> associated with the many parallel origins of C<sub>4</sub> photosynthesis within that clade (Figure 1). These<sup>67</sup> species were specifically selected to augment C<sub>3</sub>/C<sub>4</sub> comparisons: *Hymenachne amplexicaulis*,<sup>68</sup> *Chasmanthium laxum*, and *D. oligosanthes*. *H. amplexicaulis* is a member of the grass tribe Pas-

paleae which contains a mixture of C<sub>3</sub> and C<sub>4</sub> species. The Paspaleae are sister to exclusively C<sub>4</sub><sup>69</sup> clade. This C<sub>4</sub> clade is variously considered to either consider of two tribes the Andropogoneae<sup>71</sup> and Arundinelleae, or a single expanded Andropogoneae including those species otherwise included<sup>72</sup> in the Arundinelleae. In either nomenclature this clade includes both maize and sorghum, two<sup>73</sup> species with extensive genomic, genetic, and phenotypic resources. *H. amplexicaulis* is found in<sup>74</sup> moist habitats and thrives under flooded conditions (Kibbler and Bahnisch, 1999). *Chasmanthium*<sup>75</sup> *laxum* belongs to the grass tribe Chasmanthieae (7 species). The Chasmanthieae all appear to utilize<sup>76</sup> C<sub>3</sub> photosynthesis (Kellogg, 2015) and are generally placed as early diverging lineage within the<sup>77</sup> Panicoideae, the grass sub-family containing maize, sorghum, sugar cane, miscanthus, switchgrass,<sup>78</sup> foxtail millet, and proso millet (GPWG II, 2012). *C. laxum* can occur in a variety of environments<sup>79</sup> such as: woods, meadows and swamps (Yates, 1966). The final species targeted for transcriptome<sup>80</sup> sequences was *Dichanthelium oligosanthes*. *D. oligosanthes* is the only PACMAD species exclusively<sup>81</sup> utilizing C<sub>3</sub> photosynthesis with a published genome sequence to date (Studer et al., 2016).

It is a member of the grass tribe Paniceae, a group which also includes foxtail millet, proso millet,<sup>83</sup> and switchgrass, but is an outgroup to the MPC C<sub>4</sub> subclade of exclusively C<sub>4</sub>-utilizing species<sup>84</sup> within that tribe (Giussani et al., 2001; GPWG II, 2012; Washburn et al., 2015, 2017). As the<sup>85</sup> published *D. oligosanthes* reference genome was constructed utilizing short read sequencing, the<sup>86</sup> inclusion of *D. oligosanthes* provided an opportunity to improve the proportion of genes with full<sup>87</sup> length sequences from this lineage available for comparative analyses. *D. oligosanthes* is present<sup>88</sup> in small glades on

the edge of woods (A.J. Studer, personal communication, April 08, 2019). The <sup>89</sup> placement of *C. laxum* as an outgroup to other panicoid grasses with sequenced reference genomes <sup>90</sup> and *D. oligosanthos* as sister to other members of the Paniceae with sequenced reference genomes <sup>91</sup> were recovered in a preliminary analysis of our long read dataset. Support of the placement of <sup>92</sup> *H. amplexicaulis* as a sister group to Andropogoneae (sorghum and maize) was strong but not <sup>93</sup> unambiguous.

## <sup>94</sup> Methods

### <sup>95</sup> Plant material, RNA extraction, and sequencing

<sup>96</sup> For all three species, young leaf tissue was harvested from mature plants growing in the greenhouses <sup>97</sup> of the University of Nebraska's Beadle Center, latitude: 40.8190, longitude: -96.6932, on October <sup>98</sup> 05 2017. Young leaves were harvested from a *C. laxum* plant germinated from seed collected <sup>99</sup> with accession Kellogg 1268 in Corkwood Conservation Area, just outside of Neelyville, MO, USA.

<sup>100</sup> Full details of this collection are published on Tropicos: <https://www.tropicos.org/Specimen/> <sup>101</sup> 100877982. Leaf tissue from *D. oligosanthos* was harvested from a plant descended from Kellogg <sup>102</sup> 1175, which was collected in Shaw Nature Reserve, west of St. Louis, MO, USA. Full details of <sup>103</sup> this collection are published on Tropicos: <http://www.tropicos.org/Specimen/100315254>. The <sup>104</sup> specific *D. oligosanthos* plant used as a tissue donor had experienced at least three generations of <sup>105</sup> selfing relative to the originally collected plant. This selfing occurred via an independent lineage <sup>106</sup> from the F2 plant derived from the same collection which was used to generate the DNA for <sup>107</sup> the *D. oligosanthos* reference genome Studer et al. (2016). Young leaves were harvested from *H.*

<sup>108</sup> *amplexicaulis* which had been clonally propagated from collection PH2016. PH2016 was originally <sup>109</sup> collected by Pu Huang in Myakka River state park in Florida, USA on March 22nd, 2016. A <sup>110</sup> clone of this same accession, grown in the same greenhouse, is deposited at the University of <sup>111</sup> Nebraska-Lincoln Herbarium with index number NEB-328848.

<sup>112</sup> Tissue samples were ground in liquid N<sub>2</sub> and then approximately 200mg of powdered tissue was <sup>113</sup> added to 2 μL of TriPure isolation reagent (Roche Life Science, catalog number #11667157001). <sup>114</sup> The RNA samples mixed with TriPure were then separated using chloroform, precipitated using <sup>115</sup> isopropanol, and RNA pellets were washed using 75% ethanol. The samples were air-dried and <sup>116</sup> diluted in RNeasy (Ambion). Total RNA concentration was measured using a NanoDrop 1000 <sup>117</sup> spectrophotometer and the integrity was assessed based on electrophoresis on a 1% agarose gel. <sup>118</sup> 10 μL of total RNA for each species was shipped to the Duke Center for Genomic and Computational <sup>119</sup> Biology (GCB), Duke University, USA. Concentrations at the time of shipment ranging from 226.07 <sup>120</sup> to 1,374 ng/μL. OD260/280 ratios for RNA samples were of: 1.93 (*H. amplexicaulis*), 1.92 (*C.*

<sup>121</sup> *Laxum*) and 2.03 (*D. oligosanthos*) within the recommended range for IsoSeq library construction <sup>122</sup> of 1.8-2.2 provided by PacBio. One IsoSeq library was constructed per species and each library <sup>123</sup> was sequenced using a single SMRT cell on a PacBio Sequel.

### <sup>124</sup> Consensus reads and transcriptome assembly

<sup>125</sup> Two separate sequence datasets were produced per library: full length (FL) transcripts and non <sup>126</sup> full length (NFL) transcripts. A given transcript was considered FL if the sequence read contained <sup>127</sup> both 5' and 3' adapters as well as poly-A tail and are not redundant to other transcripts. The tran <sup>128</sup> scripts lacking the poly-A tail or one of the adapters are instead included in the non-full length <sup>129</sup> dataset. Sequence reads from both files were used to assemble consensus transcriptomes using

<sup>130</sup> the software pbtranscript to cluster redundant sequences, part of the SMRT pipe package (ver- <sup>131</sup> sion 5.1) with default parameters ([https://www.pacb.com/wp-content/uploads/SMRT\\_Tools\\_](https://www.pacb.com/wp-content/uploads/SMRT_Tools_Reference_Guide_v600.pdf) <sup>132</sup> [Reference\\_Guide\\_v600.pdf](https://www.pacb.com/wp-content/uploads/SMRT_Tools_Reference_Guide_v600.pdf)). For each final consensus transcript, the single longest ORF present <sup>133</sup> within that transcript was selected as the CDS sequence for downstream analyses. Consensus <sup>134</sup> transcripts were obtained from full-length and non-full length transcripts generated from oligo-dT <sup>135</sup> purified mRNAs. As oligo-dT purification can capture the 3' ends of partially fragmented mRNA <sup>136</sup> molecules it was anticipated that some transcripts may be less than full length and missing start <sup>137</sup> codons. Therefore, ORFs were required to include an in frame stop codon but were not required to <sup>138</sup> include an

in-frame "ATG" which may result in additional non-translated codons being appended<sup>139</sup> to the 5' end of the putative CDS but avoids CDS truncation when the 5' end of the sequence was<sup>140</sup> not recovered.

#### <sup>141</sup>Sequence data from species with published reference genomes

<sup>142</sup> CDS file containing only one primary transcript per gene downloaded from Phytozome 12 (<https://phytozome.jgi.doe.gov/pz/portal.html>) was used from *Brachypodium distachyon* (Initiative<sup>144</sup> et al., 2010), *Oryza sativa* (rice) (Yu et al., 2005; Kawahara et al., 2013), *Sorghum bicolor* (sorghum)<sup>145</sup> (Paterson et al., 2009) and *Setaria italica* (foxtail millet) (Bennetzen et al., 2012). CDS sequences<sup>146</sup> for version 2 of the *Oropetium thomaeum* (oropetium) genome (GenomeID 51527) (VanBuren<sup>147</sup> et al., 2018), and the draft *Eragrostis tef* genome (GenomeID 50954) (VanBuren et al., 2019)<sup>148</sup> were downloaded from CoGe (Lyons and Freeling, 2008). CDS sequences for the initial release<sup>149</sup> of the *Pennisetum glaucum* (pearl millet) genome where downloaded from GigaDB. (Varshney<sup>150</sup> et al., 2017b,a). CDS sequences for B73\_RefGenV4 of the *Zea mays* (maize) reference genome was<sup>151</sup> retrieved from Ensembl (Jiao et al., 2017). In cases where only a complete set of CDS sequences<sup>152</sup> was released for a given species, we arbitrarily selected the longest annotated transcript from a<sup>153</sup> given locus to be the single representative transcript for downstream analyses. mRNA sequences<sup>154</sup> from *D. oligosanthos* were obtained CoGe (Genome ID 35847) (Studer et al., 2016) for comparison<sup>155</sup> with IsoSeq *D. oligosanthos* transcripts. PEPC and PPK gene families from sorghum for further<sup>156</sup> manual curation and phylogeny analysis were obtained from Christin et al. (2007) and Wang et al. (2009), respectively.<sup>157</sup>

#### <sup>158</sup>Putative orthology assignments

<sup>159</sup> CDS sequences obtained from *H. amplexicaulis*, *C. laxum* and *D. oligosanthos* as described above<sup>160</sup> were compared to the primary CDS sequences of each annotated gene in the sorghum genome using<sup>161</sup> LASTZ version 1.04.00 (Harris, 2007) with the following parameters: -identity=70 -coverage=50<sup>162</sup> -ambiguous=iupac, -notransition, and -seed=match12. CDS sequences from the three target<sup>163</sup> species were presumed to belong to an orthologous group as a given sorghum gene if the sorghum<sup>164</sup> CDS sequence and target species CDS sequence were reciprocally identified as each others high<sup>165</sup> scoring hit in the LASTZ analysis. The comparison of *D. oligosanthos* consensus transcripts and<sup>166</sup> annotated mRNA sequences from the published reference genome for *D. oligosanthos* defined equiv<sup>167</sup> alence as sequences which were reciprocally identified as highest scoring LASTZ matches.

<sup>168</sup> Orthologous relationships between sorghum genes and genes in other species with sequenced<sup>169</sup> reference genomes were inferred based on syntenic orthology. For each combination of sorghum<sup>170</sup> and rice, brachypodium, oropetium, teff, foxtail millet, pearl millet, sorghum, and maize all by all<sup>171</sup> LASTZ comparisons were performed using the same parameters described above. The resulting<sup>172</sup> LASTZ output was employed to identify initial syntenic genomic blocks using QuotaAlign with<sup>173</sup> the parameters -tandemNmax=10, cscore=0.5, -merge and -Dm=20 (Tang et al., 2011). The<sup>174</sup> quota was set to -quota=1:2 for maize and teff, and -quota=1:1 for all other species. Pairwise<sup>175</sup> syntenic block data was merged and polished using the methodology previously described in (Zhang<sup>176</sup> et al., 2017) to obtain the final set of high confidence syntenic ortholog groups employed for all<sup>177</sup> downstream analysis.

<sup>178</sup> Orthology was treated as a transitive property, thus each *H. amplexicaulis*, *C. laxum* or *D. oligosanthos* gene identified as putatively orthologous to a given sorghum gene based on reciprocal<sup>179</sup> best LASTZ hit analysis, was also considered to be putatively orthologous to syntenic orthologs<sup>181</sup> of that sorghum gene identified in each of the other species described above. The final sets of<sup>182</sup> putatively orthologous gene groups including both sequences from published reference genomes<sup>183</sup> and the long read sequencing described here is provided in Supplemental Material 1.

#### <sup>184</sup>Sequence alignment, QC, and phylogenetic analysis

<sup>185</sup> Kalign (v2.04) was used create a multiple sequence alignment from protein sequences obtained by<sup>186</sup> translating CDS sequences from all genes in a give putatively orthologous gene group. This gapped<sup>187</sup> protein alignment was in turn employed to create a codon-level DNA alignment of the original CDS<sup>188</sup> sequences. GBLOCKS version 0.91 was run with default parameters to identify high quality portions<sup>189</sup> of the sequence alignment and remove those portions of the alignment not meeting specified quality<sup>190</sup> thresholds (Talavera

and Castresana, 2007). Alignments including only those portions passing 191 GBlocks filtering were then used as input for RAxML version 8, using the GTRGAMMA model 192 and with a clade of rice and brachypodium specified as an outgroup, to obtain a phylogenetic tree 193 for each group of putatively orthologous genes (Stamatakis, 2014). When RAxML was unable to 194 construct a phylogeny in which rice and brachypodium formed monophyletic clade sister to other 195 other taxa the trees were omitted from downstream visualization. To plot all phylogenies, we used 196 Densitree, part of the BEAST2 package, was used to create combined blots of large numbers of trees 197 (Bouckaert et al., 2014). We performed bootstrap analyses of 100 randomly chosen trees (from 198 all trees generated without any filtering step) with 100 replicates to obtain the branch support 199 values of the most common tree topologies of the Densitree plot. For the trees built for PPDK 200 and PEPC, RAxML bootstrap analyses were performed with 1000 replicates. For visualization 201 purposes only, all branches were treated as having equal length in order to improve the ease of 202 visually comparing differences in topology. Data on the consistency or inconsistency of individual 203 portions of the phylogeny was judged from comparisons of the single best trees generated using the 204 sequence of each separate gene. However, it should be noted that bootstrapping was not performed 205 for the individual best trees for each gene, hence some disagreements in topology may simply reflect 206 poorly resolved nodes with limited support.

207 As a result of the separate whole genome duplications in the maize and teff lineages, in many 208 cases gene and species trees would contain different numbers of leaf nodes. For gene groups where 209 maize and teff had each fractionated back to single copy status, only a single alignment file was 210 created. If fractionation had already occurred in one lineage, but not the other, two separate 211 alignments were created, each sampling one of the two co-orthologous gene copies from the species 212 with a retained whole genome duplication derived gene pair. When fractionation had not occurred 213 in either lineage, four total alignments were generated per gene group, capturing all possible 214 pairwise combinations of the two teff gene copies and two maize gene copies.

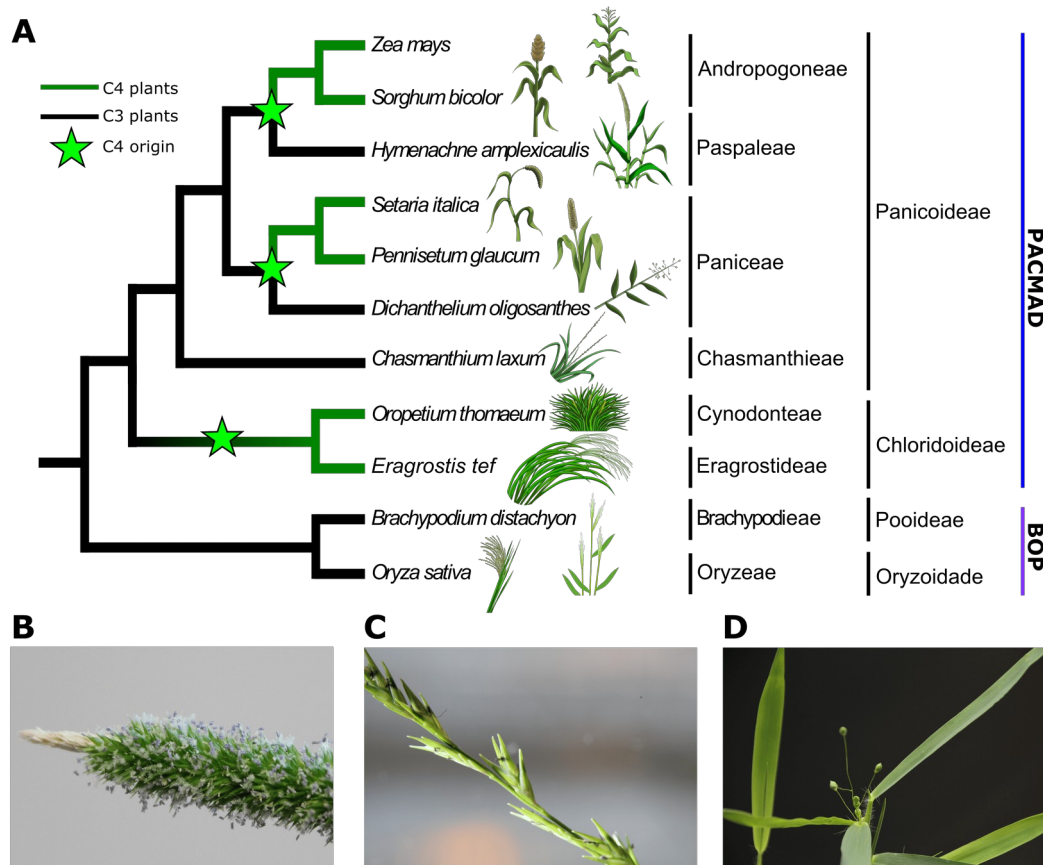


Figure 1: A) Current literature consensus phylogeny of the relationships between the grass species studied here. Lineages in green utilize C<sub>4</sub> photosynthesis, while lineages in black utilize C<sub>3</sub> photosynthesis. The green stars indicate apparent independent origins of C<sub>4</sub> photosynthesis. B) Inflorescence of *H. amplexicaulis*. C) Inflorescence of *C. laxum*. D) Inflorescence of *D. oligosanthes*.

## Results and Discussion

The number of raw reads generated per species was largely consistent and ranged from 708,681 to 734,932 (Table 2). After clustering both full length and non-full length transcripts to obtain a set of polished consensus transcripts, the number of sequences per species dropped to 164,640

to 193,422 (Table 2). The average length of consensus sequences ranged from 925 bp to 1,438 kb (Figure S1). A comparison between the mRNA sequences and IsoSeq transcripts from *D. oligosanthes* was performed to assess the improvement of sequence coverage. Out of the 13,847 reciprocal best LASTZ hits between the mRNA and IsoSeq data, 12,347 transcripts were longer than the mRNA sequences, while the remaining 1,500 sequences were either shorter or the same length in both datasets (Figure S2). The longer sequences recovered in this experiment may be the result

of a combination of fragmentation within the *D. oligosanthes* reference genome and improved capture of 5' and 3' untranslated regions, which are often missed by homology-based annotation of

genomic sequence. While *H. amplexicaulis* exhibited the shortest consensus transcript length, this was not reflected in a reduced number of complete ORFs – those containing both an in frame ATG and stop codon and occupying at least 60% of the total transcript length. The number of consensus transcripts significantly exceeded the expected number of expressed genes, however, this is consistent with other reference genome-free IsoSeq analyses (Li et al., 2017; Kuang et al., 2019; Yan et al., 2019). Inflated numbers of consensus transcripts can result from sequencing of multiple alternatively spliced isoforms of the same gene, sequencing of incompletely processed mRNA molecules (Martin et al., 2014), high sequence error rates preventing multiple sequences from the same transcript being collapsed into a consensus, divergent haplotypes of the same locus present in our clonally propagated, wild collected, or

partially inbred starting material, or contamination <sup>237</sup> of the original samples with mRNA from non-target organisms.

of the original samples with mRNA from non-

Species	Total reads	Raw data	CCS reads	FL reads	Average FL length	Consensus transcripts	Average consensus transcript length	Transcripts containing start codon
<i>Hymenachne amplexicaulis</i>	734,932 reads	5.8 GB	732,158 reads	284,027 reads	963 bp	193,422	925 bp	34,016/193,422 (17.5%)
<i>Dichanthelium oligosanthes</i>	708,681 reads	10.1 GB	701,802 reads	380,381 reads	1,460 bp	190,632	1,438 bp	6,055/190,632 (3.2%)
<i>Chasmanthium laxum</i>	729,710 reads	12.5 GB	649,149 reads	306,566 reads	1,294 bp	164,640	1,236 bp	

Table 2: Summary statistics for raw and processed long read sequence data generated from each of the three target species.

<sup>238</sup> Alignment of final consensus reads to the sorghum reference genome was employed to estimate <sup>239</sup> coverage of the shared grass gene space for data collected from each of our target species, as well as <sup>240</sup> to assist in further collapsing multiple redundant sequences originating from alternative splicing, <sup>241</sup> incomplete processing, or divergent haplotypes of transcripts originating from a single genetic <sup>242</sup> locus. In all three cases the majority of consensus transcripts could be aligned to known genes in <sup>243</sup> the sorghum genome, with an average of between 9.3 and 12.1 consensus transcripts aligning to <sup>244</sup> each sorghum gene represented in the transcriptome data (Table 3). Each of these three target <sup>245</sup> species is predicted to be diploid based on either flow cytometry based estimates of genome size <sup>246</sup> and/or imaging of chromosomes, thus a maximum of two transcripts per locus can be explained by <sup>247</sup> divergent haplotypes. The high number of consensus sequences aligned per represented sorghum <sup>248</sup> locus suggests that a large proportion of the overall inflation in consensus transcript number from <sup>249</sup> this dataset may result from alternative splice isoforms or sequencing of incompletely processed <sup>250</sup> mRNA molecules. It should also be noted that this analysis will confound lineage specific gene <sup>251</sup> duplications with divergent haplotypes and splice isoforms, however this bias will be consistent <sup>252</sup> across all three species.

Species	Sorghum gene space coverage	Sorghum syntenic gene space coverage	Transcript alignment rate
<i>H. amplexicaulis</i>	11,485 genes/34,211 genes (33.5%)	6,402 transcripts/11,800 genes (54.2%)	115,361 transcripts/193,422 transcripts (59.6%)
<i>D. oligosanthes</i>	14,159 genes/34,211 genes (41.3%)	7,760 transcripts/11,800 genes (65.7%)	171,465 transcripts/190,632 transcripts (89.9%)
<i>C. laxum</i>	13,446 genes/34,211 genes (39.3%)	7,418 transcripts/11,800 genes (62.8%)	125,357 transcripts/164,640 transcripts (76.1%)

Table 3: Alignment rates of consensus transcripts generated from each of the three target species to the sorghum gene space.

<sup>253</sup> For each sorghum gene which aligned to two or more consensus transcripts from the same <sup>254</sup> target species, a single representative transcript was selected for further downstream analysis (See <sup>255</sup> Methods). Between 11,485 and 14,159 sorghum genes had a corresponding representative transcript <sup>256</sup> in a given target species (Table 3). Here we were: 1) using a single library constructed per species, <sup>257</sup> rather than multiple libraries constructed using different size fractions, 2) using RNA from a single <sup>258</sup> tissue rather than pooled RNA from multiple tissue types, and 3) conducting comparisons between <sup>259</sup> more distantly related species. However, the total proportion of sorghum genes represented in each <sup>260</sup> transcriptome dataset was not substantially lower than the 14,401 *T. dactyloides*-maize gene pairs <sup>261</sup> identified in a previous study which implemented all of these best practices (Yan et al., 2019). This <sup>262</sup> may in part be explained by both sequencing and library preparation improvements between the <sup>263</sup> RSII and Sequel iterations of this sequencing technology.

<sup>264</sup> Manual curation was used to access the coverage and quality of sequences retrieved from these <sup>265</sup> three *C<sub>3</sub>* photosynthesis-utilizing PACMAD species for five genes known to be involved in *C<sub>4</sub>* <sup>266</sup> photosynthesis: PPKK, PEPC, NADP-MDH, NAD-ME and DCT2 in *C<sub>4</sub>* photosynthesis-utilizing <sup>267</sup> PACMAD species. In four cases, the representative transcript identified from each of the three <sup>268</sup> target species spanned every annotated codon in sorghum. The one exception was PPKK where <sup>269</sup> the representative transcript identified for *H. amplexicaulis* lacked the first annotated exon of the <sup>270</sup> annotated gene model in sorghum (Figure 2). Multiple isoforms of the PPKK gene have been <sup>271</sup> described in both maize and sorghum, with the shorter isoform, lacking the same exon absent in <sup>272</sup> *H. amplexicaulis* (Sheen, 1991; Wang et al., 2009). This shorter isoform lacks the chloroplast transit <sup>273</sup> peptide and encodes cytosolic PPKK protein not thought to be associated with *C<sub>4</sub>* photosynthesis <sup>274</sup> (Glackin and Grula, 1990; Sheen, 1991; Wang et al., 2009). However, these results would also be <sup>275</sup> consistent with sequencing of an incomplete transcript from *H. amplexicaulis* with break point <sup>276</sup> occurring at the same location as the 3' junction of the first exon.

<sup>277</sup> Phylogenetic consistency was assessed using a small subset of genes with high confidence syn <sup>278</sup> tenic orthologs identified in species with published reference genomes and representative transcripts <sup>279</sup> identified in each of the three target species. A subset of PACMAD species with sequenced refer <sup>280</sup> ence genomes were

included in these analyses. Excluded species included those with fragmented <sup>281</sup> genome assemblies at the time of this analyses as well as many, but not all, species with inde-

<sup>282</sup> pendent whole genome duplications (i.e. *Panicum virgatum*, *Miscanthus x giganteus* and *Eleusine*  
<sup>283</sup> *coracana*) as these increased the complexity of downstream analyses. Two well characterized C<sub>4</sub>  
<sup>284</sup> related proteins PPKK and PEPC were evaluated in detail. PPKK is a member of a small gene <sup>285</sup> family, while genes encoding PEPC proteins are more numerous in many grass species. Copies of <sup>286</sup> each of these two genes in sorghum previously identified as involved in the C<sub>4</sub> cycle were identified <sup>287</sup> from the literature [Christin et al. \(2007\)](#); [Wang et al. \(2009\)](#). Phylogenies constructed using known <sup>288</sup> gene copies from species with published reference genomes and isoseq transcripts of both PPKK <sup>289</sup> and PEPC clustered all C<sub>4</sub> gene copies together (Figure S3). A total of 11,800 genes were identified <sup>290</sup> at syntenic orthologous locations across the genomes of rice, brachypodium, teff, Oropetium, pearl <sup>291</sup> millet, foxtail millet, sorghum, and maize. Of these in 2,774 cases no representative transcripts

<sup>292</sup> were retrieved from *C. laxum*, *D. oligoanthes*, or *H. amplexicaulis*. These cases likely represent <sup>293</sup> conserved genes that are not expressed in developing photosynthetic tissue. In 1,611 cases, a rep<sup>294</sup> resentative transcript was identified in only one of the three target species, and in 2,276 cases, <sup>295</sup> representative transcripts were identified in two of the three target species. In the remaining 5,139 <sup>296</sup> cases representative transcripts were retrieved for all three target species. The complete lists of <sup>297</sup> each of these sets of conserved syntenic genes and corresponding transcripts from 0, 1, 2, or 3 of <sup>298</sup> the target species is provided as part of Supplemental Material 1.

<sup>299</sup> One potential concern is using transcriptome data from species utilizing C<sub>3</sub> photosynthesis to <sup>300</sup> provide sequence data for comparative genetic and evolutionary analyses of C<sub>4</sub> is that enzymes <sup>301</sup> involved in the C<sub>4</sub> cycle will likely different functions unrelated to photosynthesis in C<sub>3</sub> plants <sup>302</sup> ([Aubry et al., 2011](#)), and therefore may not be expressed in photosynthetic tissue and hence be <sup>303</sup> missing from from datasets derived from sequencing cDNAs. Of 31 core C<sub>4</sub> genes enumerated in <sup>304</sup> ([Huang et al., 2016](#)), 20 were part of the set of 11,800 sorghum genes with conserved syntenic <sup>305</sup> orthologs identified in each of the tested grass species with a published reference genome. Hence,

<sup>306</sup> these genes are almost certainly present within the genomes of *C. laxum*, *D. oligoanthes*, and *H.* <sup>307</sup> *amplexicaulis* as well, whether or not they were expressed to sufficient levels to be detected in this <sup>308</sup> analysis. Of these 20 syntenically-conserved C<sub>4</sub> related genes, sequence data was obtained from <sup>309</sup> all three target C<sub>3</sub> utilizing panicoid species in 16 cases. In the remaining four cases – DCT4c, <sup>310</sup> GLR, NADP-ME and SCL – no putatively orthologous transcript was identified in any of the three <sup>311</sup> species. There were no cases where a syntenically conserved gene linked to C<sub>4</sub> photosynthesis was <sup>312</sup> detected in some, but not all, of the three C<sub>3</sub> utilizing species evaluated.

<sup>313</sup> From the list containing a total of 5,139 conserved orthologous gene groups present in all species <sup>314</sup> 231 were discarded for one of several reasons, listed from most common to least common. 1) In <sup>315</sup> 113 cases the CDS sequence for the *O. thomaeum* genome included one or more in-frame stop codons. <sup>316</sup> 2) In 61 cases in at least one species represented by isoseq data no stop codon was present in any <sup>317</sup> of the 6 possible open reading frames, indicating either a sequencing error or incomplete <sup>318</sup> 3 prime coverage. 3) In 56 cases a syntenic orthologous gene present in version 2.1 of the *B. distachyon* <sup>319</sup> genome had been removed or renamed in version 3.1 of the *B. distachyon* genome. 4) One *O.*

<sup>320</sup> *thomaeum de novo* predicted gene region was not present in the CDS data.

<sup>321</sup> The remaining set of 4,908 conserved orthologous gene groups were used to generate protein <sup>322</sup> guided codon multiple sequence alignments (See Methods). A subset of these alignments containing <sup>323</sup> at least 900 nucleotides (300 codons) alignment scored as "high quality" by GBlocks were employed <sup>324</sup> to construct individual gene-level trees (Figure S4). In total 733 trees, representing 267 putatively <sup>325</sup> orthologous gene groups were constructed. Multiple trees resulted from retained duplicate gene <sup>326</sup> pairs resulting from lineage specific whole genome duplications in maize and teff. Each duplication <sup>327</sup> had the potential to create a retained syntenic gene pair which were each co-orthologous to single <sup>328</sup> gene copies in other grass species within the analysis. In order to maintain a consistent number <sup>329</sup> of final nodes, when a retained gene pair was observed in one or both species, multiple sampled <sup>330</sup> trees were generated (see Methods). A modest bias towards over representation of retained – <sup>331</sup> rather than fractionated – genes was observed in the set of genes which were represented in the <sup>332</sup> transcriptome assemblies from all three target species: 37% (1,843/4,908) of maize genes in this set <sup>333</sup> were retained as duplicate pairs vs 30% of all syntenic maize genes, and 100% of teff genes in this <sup>334</sup> set were retained as duplicate pairs vs 91% of all syntenic teff genes. Dating multiple whole genome <sup>335</sup> duplication events

back to the Rho event in grasses is complex (McKain et al., 2016). Therefore, the lack of consideration of these events back to the Rho event could cause a limitation/ambiguity in the gene trees. The rice and brachypodium clade represented a known outgroup as these two species belong to the BOP clade which diverged from the PACMAD clade of grasses early in the evolution of this family (GPWG II, 2012). In 49 cases, RAxML was unable to place the rice-brachypodium clade as an outgroup suggested broader issues with orthology assignment, correct ORF identification, or alignment. These trees were not included in downstream analyses.

Among the 684 remaining gene trees, 291 (42.5%) produced a single topology consistent with the prior literature on the relationship of these species (Figure 3). The second and third most common topologies were each represented by less than 7% of all calculated trees, 43 and 28 cases respectively.

The second and third most common topologies differed from prior published phylogenies regarding the placement *H. amplexicaulis*. In the second most common topology *H. amplexicaulis* was placed sister to all other panicoid grass species other than *C. laxum*. In the third most common topology *H. amplexicaulis* was placed sister to the Paniceae. Parallel analysis was conducted using all 4,908 conserved orthologous gene groups, including many cases with substantially shorter regions of high quality multiple sequence alignment. The pattern of trees recovered were largely consistent with those in (Figure 3). In the "all genes" analysis, the same first and second most common topologies

were retrieved as in the long alignment only analysis. The third most common in the "all genes" topology places *C. laxum* as sister to the combined Chloridoideae and Panicoideae (Supplemental Figure S5).

A bootstrap analysis was performed and values were retrieved for each of the three most common topologies. In the analysis, 100 randomly trees were chosen to perform the bootstrap analysis (see methods). The most common topology of the 100 trees was the same as the one observed in figure 3, appearing seven times. Both second and third common topologies in figure 3 only appeared once. The bootstrap values observed for the most common topology is the average of all seven bootstrap values of each branch (Figure S6). The bootstrap values of the internal branches in the second and third most common topologies are substantially smaller than the values obtained in the most common topology. Overall, the most common topology exhibits stronger bootstrap support values compared to the other two topologies.

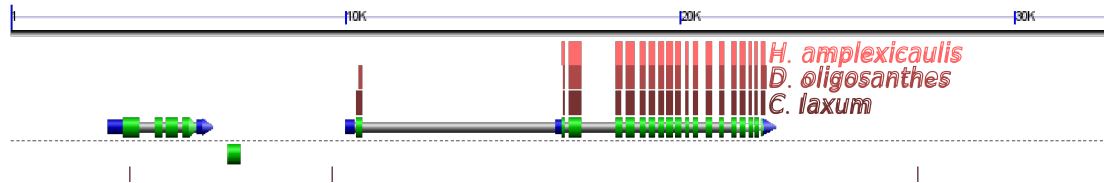


Figure 2: A GEvo panel showing transcript coverage of the C<sub>4</sub> PPK gene in *S. bicolor* Sobic.009G132900 in each of the three species texted. Red-brown boxes represent regions of similar sequence identified by BLASTN between the sorghum genome and consensus transcript sequences retrieved from *H. amplexicaulis*, *D. oligosanthos*, *C. laxum* (from top most to bottom most). The bottom track indicates the annotated gene structure, with intronic sequence indicated in gray and exonic sequence indicated in either blue (5' or 3' untranslated regions) or green (coding sequence). Top y-axis indicates scale of the displayed genomic region in kilobases. (Lyons and Freeling, 2008).

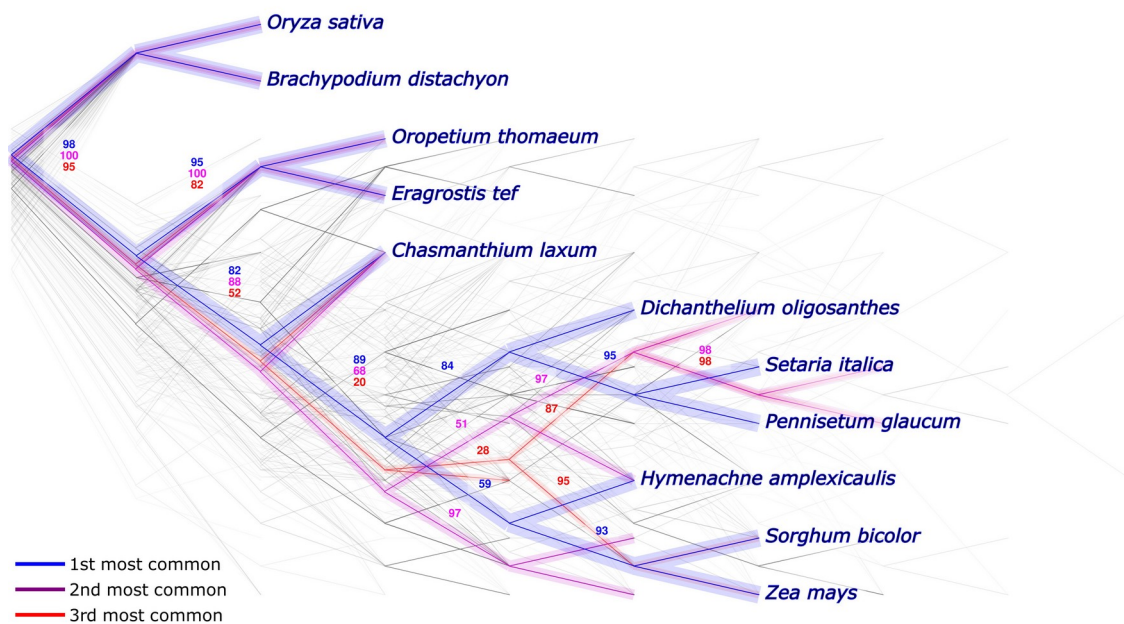


Figure 3: Seven hundred distinct phylogenetic trees calculated from separate multiple sequence alignments of 267 putatively orthologous gene groups with large regions of alignment scored as high quality. Blue indicates the most commonly observed topology (291 trees (42.5% of the total), purple and red indicate the second (43 trees (6.2%) and third most commonly observed topologies (28 trees (4%)), respectively. Numerical labels of branches for each topology indicate average bootstrap support from separately calculated bootstrap trees for 100 randomly selected gene groups, considering data only from those gene trees which supported that particular topology.

### 364 Data availability

365 Raw sequence data for *C. laxum*, *H. amplexicaulis* and *D. oligosanthes* have been deposited in the 366 NCBI SRA under accessions numbers: SRR7632721 (*C. laxum*), SRR7632716 (*H. amplexicaulis*) 367 and SRR9603193 (*D. oligosanthes*). Processed consensus transcript sequences generated for each 368 of these three species have been deposited in Zenodo DOI [10.5281/zenodo.3253206](https://doi.org/10.5281/zenodo.3253206). Syntenic gene 369 sets and putatively orthologous relationships are provided as Supplemental Material 1.

### 370 Competing interests

371 The authors declare that they have no competing interests.

### 372 Author's contributions

373 DSC and JSC wrote the paper and designed the experiments, DSC generated and analyzed the 374 transcriptome data. All authors have reviewed and approved the manuscript.

### 375 Acknowledgements

376 This work was supported by a CNPq Science Without Borders scholarship (214038/2014-9) to DSC. 377 This material is based upon work supported by the National Science Foundation under Grant No. 378 MCB-183830 to JCS. This project was supported by the Agriculture and Food Research Initiative 379 Grant number 2016-67013-2461 from the USDA National Institute of Food and Agriculture to 380 JCS. In addition, portions of this project were supported by Robert B. Daugherty Water for Food 381 Institute research support award to JCS. The authors would like to thank Lindsay Erndwein who 382 drew the illustrations of various grass species presented in Figure 1A.

383 Supplemental figures

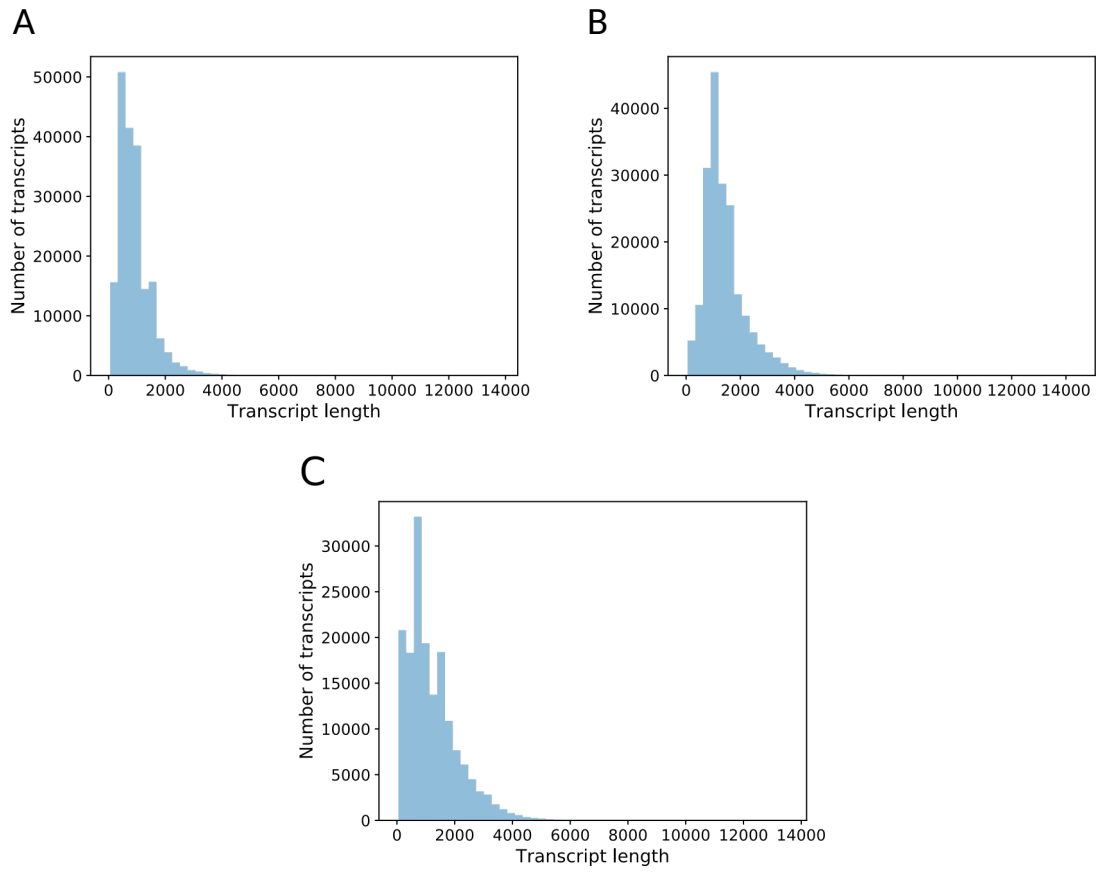


Figure S1: Distribution of lengths for polished transcript sequences: A) *H. amplexicaulis*, B) *D. oligosanthos*, C) *C. laxum*.

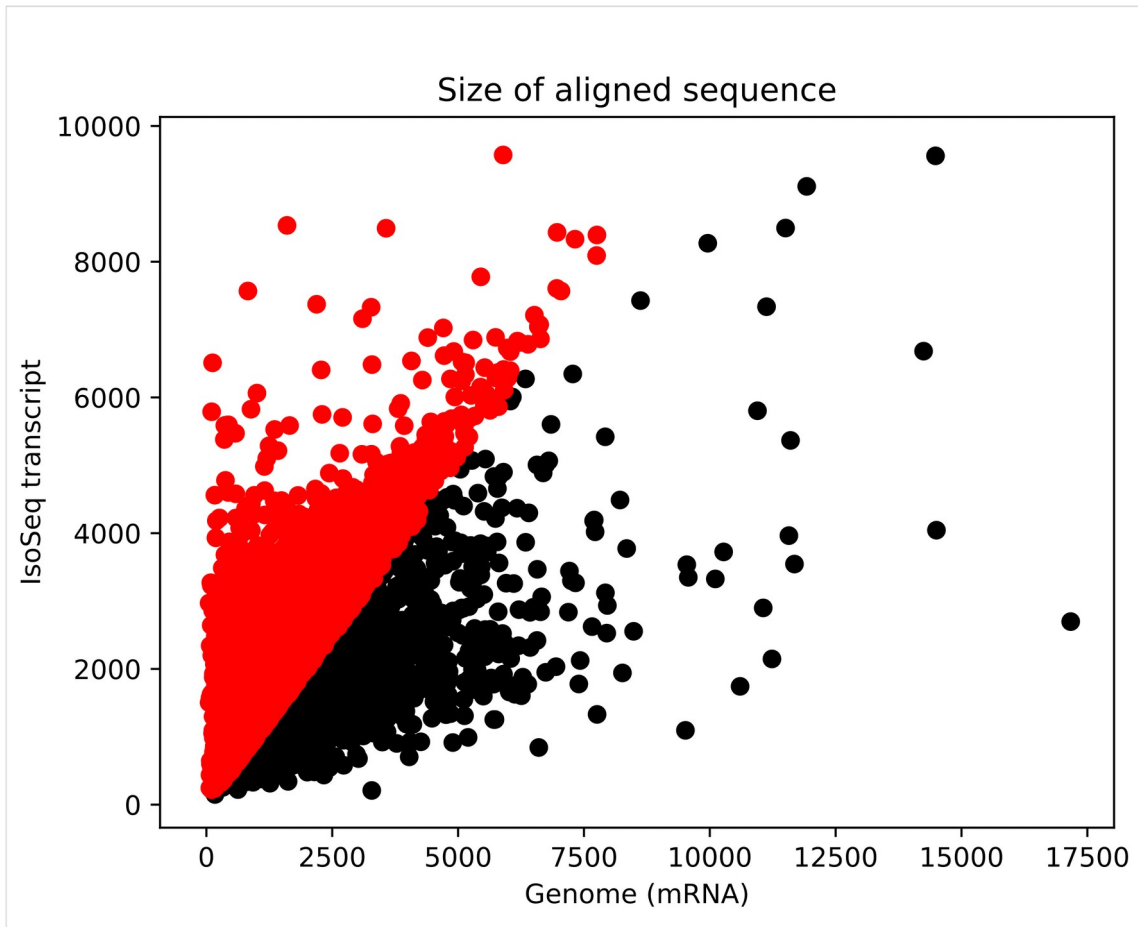


Figure S2: Relationship between aligned mRNA and IsoSeq transcript sizes. Red points indicate sequences where the transcript identified using IsoSeq was longer than the mRNA annotated in the published reference genome for *D. oligosanthos*. Black points indicate sequences where the annotated mRNA sequence was longer than the corresponding sequenced transcript.



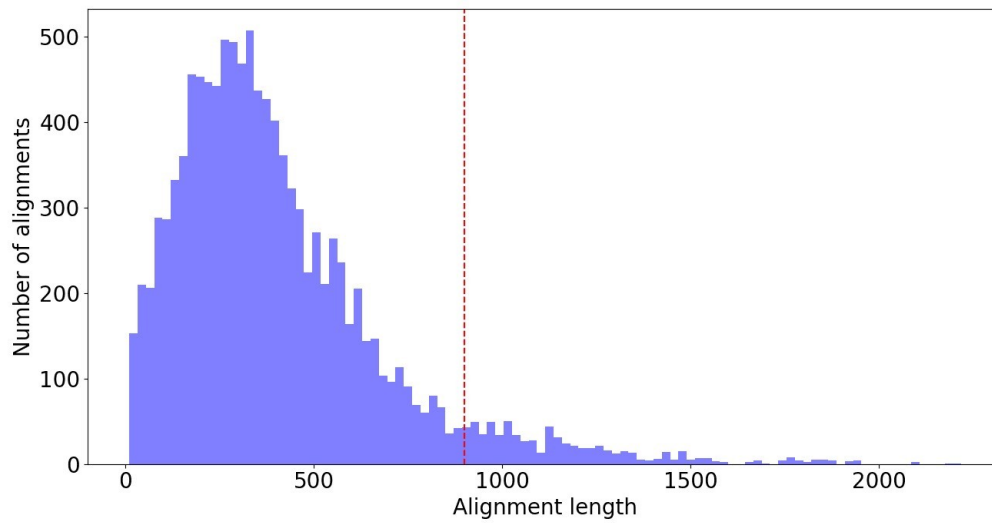


Figure S4: Distribution of alignment lengths after GBlocks cleaning. Dashed line represents the threshold of 900 nucleotides long sequences employed for Figure 3. Sequences represented on the right side of the histogram were analyzed.

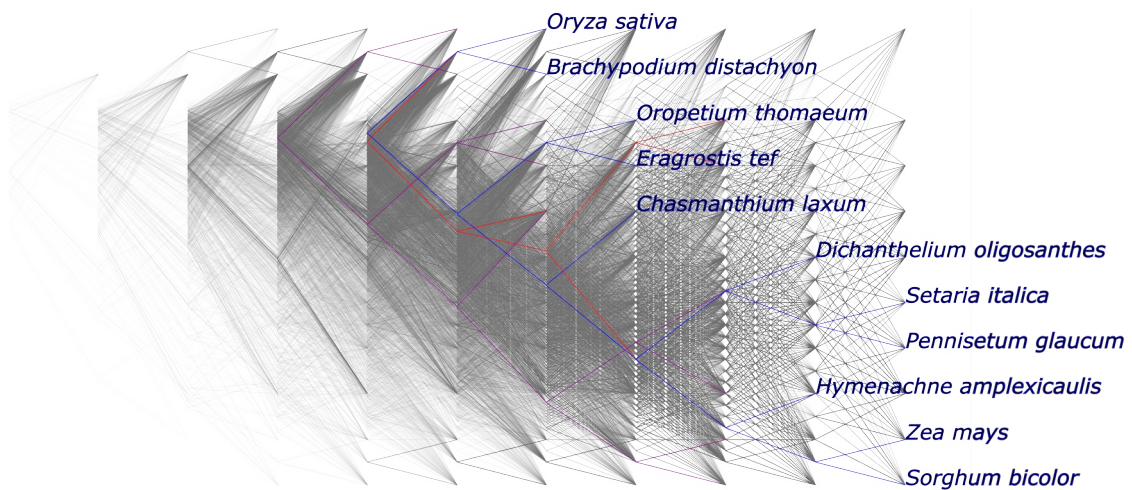


Figure S5: Plot of 11,371 phylogenetic trees. The blue branches represent the most common topology, purple and red branches represent second and third most common topologies, respectively. Figure generated using Densitree.

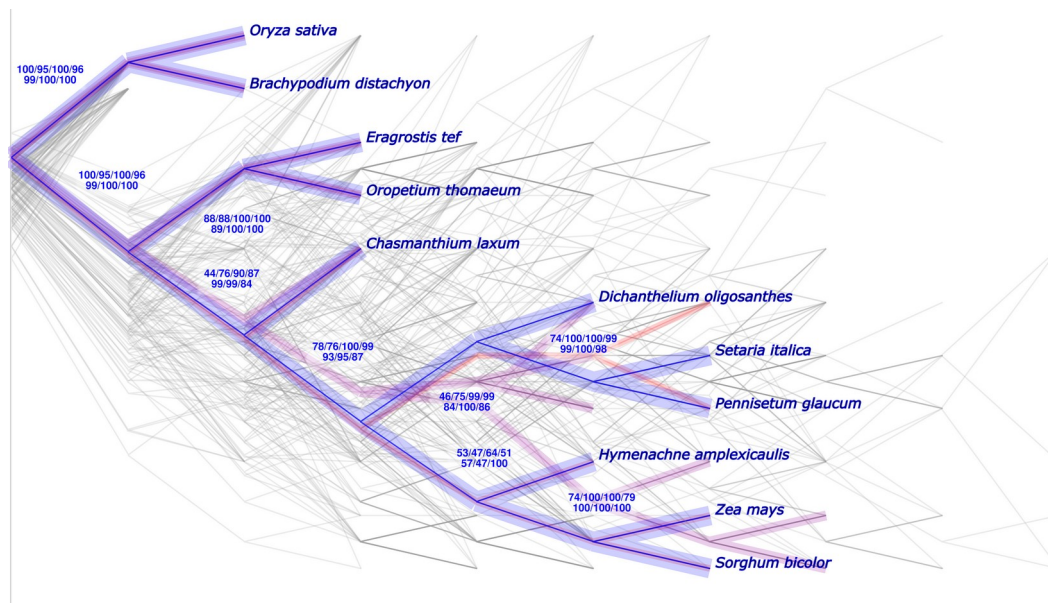


Figure S6: Plot of 100 phylogenetic trees randomly selected. The blue branches represent the most common topology (7 trees), purple and red branches represent the other two most common topologies (2 trees each). Bootstrap values (in blue) obtained from each of the seven trees of the common topology. The bootstrap values in the same position are from the same tree. Figure generated using Densitree.

## References

- An, D., Cao, H., Li, C., Humbeck, K. and Wang, W. (2018). Isoform sequencing and state-of-art applications for unravelling complexity of plant transcriptomes. *Genes* 9, 43.
- Aubry, S., Brown, N. J. and Hibberd, J. M. (2011). The role of proteins in c3 plants prior to their recruitment into the c4 pathway. *Journal of experimental botany* 62, 3049–3059.
- Bennetzen, J. L., Schmutz, J., Wang, H., Percifield, R., Hawkins, J., Pontaroli, A. C., Estep, M., Feng, L., Vaughn, J. N., Grimwood, J. et al. (2012). Reference genome sequence of the model plant setaria. *Nature biotechnology* 30, 555.
- Bouckaert, R., Heled, J., Kühnert, D., Vaughan, T., Wu, C.-H., Xie, D., Suchard, M. A., Rambaut, A. and Drummond, A. J. (2014). Beast 2: a software platform for bayesian evolutionary analysis. *PLoS computational biology* 10, e1003537.
- Brutnell, T. P., Wang, L., Swartwood, K., Goldschmidt, A., Jackson, D., Zhu, X.-G., Kellogg, E. and Van Eck, J. (2010). *Setaria viridis*: a model for c4 photosynthesis. *The Plant Cell* 22, 2537–2544.
- Cannarozzi, G., Plaza-Wüthrich, S., Esfeld, K., Larti, S., Wilson, Y. S., Girma, D., de Castro, E., Chanyalew, S., Blösch, R., Farinelli, L. et al. (2014). Genome and transcriptome sequencing identifies breeding targets in the orphan crop tef (*eragrostis tef*). *BMC genomics* 15, 581.
- Casler, M. D., Tobias, C. M., Kaeppler, S. M., Buell, C. R., Wang, Z.-Y., Cao, P., Schmutz, J. and Ronald, P. (2011). The switchgrass genome: tools and strategies. *The Plant Genome* 4, 273–282.
- Christin, P.-A., Arakaki, M., Osborne, C. P. and Edwards, E. J. (2015). Genetic enablers underlying the clustered evolutionary origins of c4 photosynthesis in angiosperms. *Molecular biology and evolution* 32, 846–858.
- Christin, P.-A., Salamin, N., Savolainen, V., Duvall, M. R. and Besnard, G. (2007). C4 photosynthesis evolved in grasses via parallel adaptive genetic changes. *Current Biology* 17, 1241–1247.

- <sup>411</sup> Ellegren, H. (2014). Genome sequencing and population genomics in non-model organisms. <sup>412</sup> *Trends in ecology & evolution* 29, 51–63.
- <sup>413</sup> Garsmeur, O., Droc, G., Antonise, R., Grimwood, J., Potier, B., Aitken, K., Jenkins, <sup>414</sup> J., Martin, G., Charron, C., Hervouet, C. et al. (2018). A mosaic monoploid reference <sup>415</sup> sequence for the highly complex genome of sugarcane. *Nature communications* 9, 2638.
- <sup>416</sup> Giussani, L. M., Cota-Sánchez, J. H., Zuloaga, F. O. and Kellogg, E. A. (2001). A <sup>417</sup> molecular phylogeny of the grass subfamily panicoideae (poaceae) shows multiple origins of c4 <sup>418</sup> photosynthesis. *American Journal of Botany* 88, 1993–2012.
- <sup>419</sup> Glackin, C. A. and Grula, J. W. (1990). Organ-specific transcripts of different size and <sup>420</sup> abundance derive from the same pyruvate, orthophosphate dikinase gene in maize. *Proceedings* <sup>421</sup> of the National Academy of Sciences 87, 3004–3008.
- <sup>422</sup> Gonzalez-Garay, M. L. (2016). Introduction to isoform sequencing using pacific biosciences <sup>423</sup> technology (iso-seq). In *Transcriptomics and Gene Regulation*, pp. 141–160. Springer. <sup>424</sup> GPWG II (2012). New grass phylogeny resolves deep evolutionary relationships and discovers c4 <sup>425</sup> origins. *New Phytologist* 193, 304–312.
- <sup>426</sup> Harris, R. S. (2007). *Improved pairwise alignment of genomic DNA*. The Pennsylvania State <sup>427</sup> University.
- <sup>428</sup> Hittalmani, S., Mahesh, H., Shirke, M. D., Biradar, H., Uday, G., Aruna, Y., Lo<sup>429</sup> hithaswa, H. and Mohanrao, A. (2017). Genome and transcriptome sequence of finger <sup>430</sup> millet (eleusine coracana (l.) gaertn.) provides insights into drought tolerance and nutraceutical <sup>431</sup> properties. *BMC genomics* 18, 465.
- <sup>432</sup> Huang, P., Studer, A. J., Schnable, J. C., Kellogg, E. A. and Brutnell, T. P. (2016). <sup>433</sup> Cross species selection scans identify components of c4 photosynthesis in the grasses. *Journal of* <sup>434</sup> *Experimental Botany* 68, 127–135.
- <sup>435</sup> Initiative, I. B. et al. (2010). Genome sequencing and analysis of the model grass brachypodium <sup>436</sup> distachyon. *Nature* 463, 763.
- <sup>437</sup> Jiao, Y., Peluso, P., Shi, J., Liang, T., Stitzer, M. C., Wang, B., Campbell, M. S., <sup>438</sup> Stein, J. C., Wei, X., Chin, C.-S. et al. (2017). Improved maize reference genome with <sup>439</sup> single-molecule technologies. *Nature* 546, 524.
- <sup>440</sup> Kawahara, Y., de la Bastide, M., Hamilton, J. P., Kanamori, H., McCombie, W. R.,  
<sup>441</sup> Ouyang, S., Schwartz, D. C., Tanaka, T., Wu, J., Zhou, S. et al. (2013). Improvement <sup>442</sup> of the oryza sativa nipponbare reference genome using next generation sequence and optical map <sup>443</sup> data. *Rice* 6, 4.
- <sup>444</sup> Kellogg, E. A. (2013). C4 photosynthesis. *Current Biology* 23, R594–R599.
- <sup>445</sup> Kellogg, E. A. (2015). Viii. subfamily panicoideae link (1827). In *Flowering Plants. Monocots*, <sup>446</sup> pp. 271–345. Springer.
- <sup>447</sup> Kibbler, H. and Bahnisch, L. (1999). Physiological adaptations of hymenachne amplexicaulis <sup>448</sup> to flooding. *Australian journal of experimental agriculture* 39, 429–435.
- <sup>449</sup> Kuang, X., Sun, S., Wei, J., Li, Y. and Sun, C. (2019). Iso-seq analysis of the taxus <sup>450</sup> cuspidata transcriptome reveals the complexity of taxol biosynthesis. *BMC plant biology* 19, <sup>451</sup> 210.
- <sup>452</sup> Li, J., Harata-Lee, Y., Denton, M. D., Feng, Q., Rathjen, J. R., Qu, Z. and Adelson, <sup>453</sup> D. L. (2017). Long read reference genome-free reconstruction of a full-length transcriptome from <sup>454</sup> astragalus membranaceus reveals transcript variants involved in bioactive compound biosynthe<sup>455</sup> sis. *Cell discovery* 3, 17031.
- <sup>456</sup> Lovell, J. T., Jenkins, J., Lowry, D. B., Mamidi, S., Sreedasyam, A., Weng, X., Barry, <sup>457</sup> K., Bonnette, J., Campitelli, B., Daum, C. et al. (2018). The genomic landscape of <sup>458</sup> molecular responses to natural drought stress in panicum hallii. *Nature communications* 9, <sup>459</sup> 5213.
- <sup>460</sup> Lyons, E. and Freeling, M. (2008). How to usefully compare homologous plant genes and <sup>461</sup> chromosomes as dna sequences. *The Plant Journal* 53, 661–673.

462 Martin, J. A., Johnson, N. V., Gross, S. M., Schnable, J., Meng, X., Wang, M.,  
463 Coleman-Derr, D., Lindquist, E., Wei, C.-L., Kaepler, S. et al. (2014). A near 464 complete snapshot of the  
zea mays seedling transcriptome revealed from ultra-deep sequencing. 465 *Scientific reports* 4, 4519.

466 McKain, M. R., Tang, H., McNeal, J. R., Ayyampalayam, S., Davis, J. I., dePamphilis, 467 C. W., Givnish, T. J.,  
Pires, J. C., Stevenson, D. W. and Leebens-Mack, J. H. 468 (2016). A phylogenomic assessment of ancient  
polyploidy and genome evolution across the 469 poales. *Genome biology and evolution* 8, 1150–1164.

470 Michael, T. P. and Jackson, S. (2013). The first 50 plant genomes. *The plant genome* 6.

471 Moreno-Villena, J. J., Dunning, L. T., Osborne, C. P. and Christin, P.-A. (2017). 472 Highly expressed genes  
are preferentially co-opted for c4 photosynthesis. *Molecular biology and 473 evolution* 35, 94–106.

474 Paterson, A. H., Bowers, J. E., Bruggmann, R., Dubchak, I., Grimwood, J., Gundlach, 475 H., Haberer, G.,  
Hellsten, U., Mitros, T., Poliakov, A. et al. (2009). The sorghum 476 bicolor genome and the diversification of  
grasses. *Nature* 457, 551.

477 Schnable, P. S., Ware, D., Fulton, R. S., Stein, J. C., Wei, F., Pasternak, S., Liang, 478 C., Zhang, J., Fulton, L.,  
Graves, T. A. et al. (2009). The b73 maize genome: complexity, 479 diversity, and dynamics. *science* 326, 1112–  
1115.

480 Sheen, J. (1991). Molecular mechanisms underlying the differential expression of maize pyruvate, 481  
orthophosphate dikinase genes. *The Plant Cell* 3, 225–245.

482 Stamatakis, A. (2014). Raxml version 8: a tool for phylogenetic analysis and post-analysis of 483 large  
phylogenies. *Bioinformatics* 30, 1312–1313.

484 Studer, A. J., Schnable, J. C., Weissmann, S., Kolbe, A. R., McKain, M. R., Shao, 485 Y., Cousins, A. B., Kellogg, E.  
A. and Brutnell, T. P. (2016). The draft genome of the 486 c 3 panicoid grass species *dichanthelium*  
*oligosanthes*. *Genome biology* 17, 223.

487 Swaminathan, K., Alabady, M. S., Varala, K., De Paoli, E., Ho, I., Rokhsar, D. S.,  
488 Arumuganathan, A. K., Ming, R., Green, P. J., Meyers, B. C. et al. (2010). Genomic 489 and small rna  
sequencing of *miscanthus* × *giganteus* shows the utility of sorghum as a reference 490 genome sequence  
for andropogoneae grasses. *Genome biology* 11, R12.

491 Talavera, G. and Castresana, J. (2007). Improvement of phylogenies after removing divergent 492 and  
ambiguously aligned blocks from protein sequence alignments. *Systematic biology* 56, 564– 493 577. 494 Tang,  
H., Lyons, E., Pedersen, B., Schnable, J. C., Paterson, A. H. and Freeling, M. 495 (2011). Screening synteny blocks  
in pairwise genome comparisons through integer programming. 496 *BMC bioinformatics* 12, 102.

497 VanBuren, R., Bryant, D., Edger, P. P., Tang, H., Burgess, D., Challabathula, D., 498 Spittle, K., Hall, R., Gu, J.,  
Lyons, E. et al. (2015). Single-molecule sequencing of the 499 desiccation-tolerant grass *oropetium thomaeum*.  
*Nature* 527, 508.

500 VanBuren, R., Wai, C. M., Keilwagen, J. and Pardo, J. (2018). A chromosome-scale 501 assembly of the model  
desiccation tolerant grass *oropetium thomaeum*. *Plant Direct* 2, e00096.

502 VanBuren, R., Wai, C. M., Pardo, J., Yocca, A. E., Wang, X., Wang, H., Chaluvadi,  
503 S. R., Bryant, D., Edger, P. P., Bennetzen, J. L. et al. (2019). Exceptional subgenome 504 stability and  
functional divergence in allotetraploid teff, the primary cereal crop in ethiopia.  
505 *bioRxiv* p. 580720.

506 Varshney, R., Liu, X., Shi, C., Vigouroux, Y. and Xu, X. (2017a). Genomic data of pearl 507 millet  
(*Pennisetum glaucum*). [dx.doi.org/10.5524/100192](https://doi.org/10.5524/100192).

508 Varshney, R. K., Shi, C., Thudi, M., Mariac, C., Wallace, J., Qi, P., Zhang, H., Zhao, 509 Y., Wang, X., Rathore, A.  
et al. (2017b). Pearl millet genome sequence provides a resource 510 to improve agronomic traits in arid  
environments. *Nature biotechnology* 35, 969.

<sup>511</sup> Wang, X., Gowik, U., Tang, H., Bowers, J. E., Westhoff, P. and Paterson, A. H. (2009). <sup>512</sup> Comparative genomic analysis of c4 photosynthetic pathway evolution in grasses. *Genome biology* <sup>513</sup> 10, R68.

<sup>514</sup> Washburn, J. D., Schnable, J. C., Conant, G. C., Brutnell, T. P., Shao, Y., Zhang, Y., <sup>515</sup> Ludwig, M., Davidse, G. and Pires, J. C. (2017). Genome-guided phylo-transcriptomic <sup>516</sup> methods and the nuclear phylogenetic tree of the paniceae grasses. *Scientific reports* 7, 13528. <sup>517</sup> Washburn, J. D., Schnable, J. C., Davidse, G. and Pires, J. C. (2015). Phylogeny and <sup>518</sup> photosynthesis of the grass tribe paniceae. *American Journal of Botany* 102, 1493–1505.

<sup>519</sup> Yan, L., Kenchanmane Raju, S. K., Lai, X., Zhang, Y., Dai, X., Rodriguez, O.,  
<sup>520</sup> Mahboub, S., Roston, R. L. and Schnable, J. C. (2019). Parallels between natural  
<sup>521</sup> selection in the cold-adapted crop-wild relative *tripsacum dactyloides* and artificial selection in <sup>522</sup> temperate adapted maize. *The Plant Journal* .

<sup>523</sup> Yates, H. O. (1966). Revision of grasses traditionally referred to *uniola*, ii. *chasmanthium*. *The* <sup>524</sup> *Southwestern Naturalist* pp. 415–455.

<sup>525</sup> Yu, J., Wang, J., Lin, W., Li, S., Li, H., Zhou, J., Ni, P., Dong, W., Hu, S., Zeng, C. <sup>526</sup> et al. (2005). The genomes of *oryza sativa*: a history of duplications. *PLoS biology* 3, e38.

<sup>527</sup> Zhang, G., Sun, M., Wang, J., Lei, M., Li, C., Zhao, D., Huang, J., Li, W., Li, S., Li, <sup>528</sup> J. et al. (2019). Pacbio full-length cDNA sequencing integrated with rna-seq reads drastically <sup>529</sup> improves the discovery of splicing transcripts in rice. *The Plant Journal* 97, 296–305.

<sup>530</sup> Zhang, Y., Ngu, D. W., Carvalho, D., Liang, Z., Qiu, Y., Roston, R. L. and Schnable, <sup>531</sup> J. C. (2017). Differentially regulated orthologs in sorghum and the subgenomes of maize. *The* <sup>532</sup> *Plant Cell* 29, 1938–1951.

<sup>533</sup> Zou, C., Li, L., Miki, D., Li, D., Tang, Q., Xiao, L., Rajput, S., Deng, P., Peng, L., <sup>534</sup> Jia, W. et al. (2019). The genome of broomcorn millet. *Nature communications* 10, 436.