

# GigaScience

## Genome sequence and genetic diversity analysis of an under-domesticated orphan crop, white fonio (*Digitaria exilis*) --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-20-00197	
<b>Full Title:</b>	Genome sequence and genetic diversity analysis of an under-domesticated orphan crop, white fonio ( <i>Digitaria exilis</i> )	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	Giles Fellowship from the University of Georgia	Dr Jeffrey L. Bennetzen
	H2020 European Research Council (833522)	Dr Yves Van de Peer
	Seed Biotechnology Center, UC Davis	DR. Allen Van Deynze
	ICRISAT	Dr Jason Wallace
	McKnight Foundation	Dr Moussa D. Sanogo
<b>Abstract:</b>	<p><b>Background</b></p> <p><i>Digitaria exilis</i>, white fonio, is a minor but vital crop of West Africa that is valued for its resilience in hot, dry and low fertility environments and for the exceptional quality of its grain for human nutrition. The crop is plagued, however, by a low degree of improvement.</p> <p><b>Findings</b></p> <p>We sequenced the fonio genome with long-read SMRT-cell technology, yielding an ~761 Mb assembly in 3333 contigs (N50 1.73 Mb, L50 126). The assembly approaches a high level of completion, with a BUSCO score of &gt;98%. The fonio genome was found to be a tetraploid, with most of the genome retained as homoeologous duplications that differ overall by ~4.3%, neglecting indels. The two genomes within fonio were found to have begun their independent divergence ~3.1 million years ago. The repeat content (&gt;49%) is fairly standard for a grass genome of this size, but the ratio of Gypsy to Copia LTR-retrotransposons (~6.7) was found to be exceptionally high. Several genes related to future improvement of the crop were identified. Analysis of fonio population genetics, primarily in Mali, indicated that the crop has extensive genetic diversity that is largely partitioned across a north-south gradient coinciding with the Sahel and Sudan grassland domains.</p> <p><b>Conclusions</b></p> <p>We provide a high-quality assembly, annotation and diversity analysis for a vital African crop. The availability of this information should empower future research into further domestication and improvement of fonio.</p>	
<b>Corresponding Author:</b>	Allen Van Deynze, Ph.D University of California Davis Davis, CA UNITED STATES	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	University of California Davis	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Jeffrey L. Bennetzen	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Jeffrey L. Bennetzen	

	Shiyu Chen
	Xiao Ma
	Xuewen Wang
	Anna E. J. Yssel
	Srinivasa R. Chaluvadi
	Matthew Johnson
	Prakash Gangashetty
	Falalou Hamidou
	Moussa D. Sanogo
	Arthur Zwaenepoel
	Jason Wallace
	Yves Van de Peer
	Allen Van Deynze, Ph.D
<b>Order of Authors Secondary Information:</b>	
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>  Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . Information essential to interpreting the data presented should be made available in the figure legends.  Have you included all the information requested in your manuscript?	Yes
<b>Resources</b>  A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.	Yes

<p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

1 **Genome sequence and genetic diversity analysis of an under-domesticated orphan crop,**  
2 **white fonio (*Digitaria exilis*)**

3  
4 Jeffrey L. Bennetzen<sup>1,\*#</sup>, Shiyu Chen<sup>2,\*</sup>, Xiao Ma<sup>3,\*</sup>, Xuewen Wang<sup>1,\*</sup>, Anna E. J. Yssel<sup>4,5,\*</sup>,  
5 Srinivasa R. Chaluvadi<sup>1</sup>, Matthew S. Johnson<sup>6</sup>, Prakash Gangashetty<sup>7</sup>, Falalou Hamidou<sup>7</sup>,  
6 Moussa D. Sanogo<sup>8</sup>, Arthur Zwaenepoel<sup>3</sup>, Jason Wallace<sup>9</sup>, Yves Van de Peer<sup>3,4,10</sup> and Allen Van  
7 Deynze<sup>2</sup>

8  
9 <sup>1</sup>Department of Genetics, University of Georgia, Athens, GA 30602 USA,

10 <sup>2</sup>Department of Plant Sciences, Seed Biotechnology Center, University of California, Davis, CA  
11 95616 USA,

12 <sup>3</sup>Department of Plant Biotechnology and Bioinformatics, Ghent University, Ghent, Belgium VIB  
13 - UGent Center for Plant Systems Biology, Technologiepark 71, Ghent, Belgium,

14 <sup>4</sup>Centre for Microbial Ecology and Genomics, Department of Biochemistry, Genetics and  
15 Microbiology, University of Pretoria, Pretoria 0028, South Africa,

16 <sup>5</sup>Centre for Bioinformatics and Computational Biology, Department of Biochemistry, Genetics  
17 and Microbiology, University of Pretoria, Pretoria 0028, South Africa,

18 <sup>6</sup>Institute of Plant Breeding, Genetics, and Genomics, University of Georgia, Athens, GA 30602  
19 USA,

20 <sup>7</sup>International Crops Research Institute for the Semi-Arid Tropics (ICRISAT), BP 12404,  
21 Niamey, Niger,

22 <sup>8</sup>Institut d'Economie Rurale, Ministere de l'Agriculture, Cinzana, BP 214, Ségou, Mali

23 <sup>9</sup>Department of Crop and Soil Sciences, University of Georgia, Athens, GA 30602 USA,

24 <sup>10</sup>College of Horticulture, Nanjing Agricultural University, Nanjing, China.

25  
26  
27 \*These authors contributed equally to this work.

28 #corresponding author  
29  
30  
31



32 **Abstract**

33

34 **Background:** *Digitaria exilis*, white fonio, is a minor but vital crop of West Africa that is valued  
35 for its resilience in hot, dry and low fertility environments and for the exceptional quality of its  
36 grain for human nutrition. The crop is plagued, however, by a low degree of improvement.

37 **Findings:** We sequenced the fonio genome with long-read SMRT-cell technology, yielding a  
38 ~761 Mb assembly in 3333 contigs (N50 1.73 Mb, L50 126). The assembly approaches a high  
39 level of completion, with a BUSCO score of >98%. The fonio genome was found to be a  
40 tetraploid, with most of the genome retained as homoeologous duplications that differ overall by  
41 ~4.3%, neglecting indels. The two genomes within fonio were found to have begun their  
42 independent divergence ~3.1 million years ago. The repeat content (>49%) is fairly standard for  
43 a grass genome of this size, but the ratio of *Gypsy* to *Copia* LTR-retrotransposons (~6.7) was  
44 found to be exceptionally high. Several genes related to future improvement of the crop were  
45 identified. Analysis of fonio population genetics, primarily in Mali, indicated that the crop has  
46 extensive genetic diversity that is largely partitioned across a north-south gradient coinciding  
47 with the Sahel and Sudan grassland domains.

48 **Conclusions:** We provide a high-quality assembly, annotation and diversity analysis for a vital  
49 African crop. The availability of this information should empower future research into further  
50 domestication and improvement of fonio.

51

52 **Key Words:** *Digitaria exilis*, fonio, millet, polyploid, domestication, orphan crop

53

## 54 **Data Description**

55

## 56 **Background information**

57 White fonio (*Digitaria exilis*) is a vital cereal crop of West Africa, where it is commonly known  
58 as fonio or acha. A related *Digitaria* species, black fonio (*D. iburura*), is a very minor crop,  
59 mostly of Nigeria, Benin and Togo. Fonio (*D. exilis*) has exceptionally small but very nutritious  
60 grain, with both high protein and high dietary fiber content [1,2,3]. Fonio can mature in as little  
61 as eight weeks after planting, and is commonly grown without fertilizer or irrigation on poor  
62 quality soils in dry regions of the Sudan grasslands and Sahel. Although yields are low, the West  
63 African crop is harvested in early summer, where it fills a vital dietary gap before the maturation  
64 of sorghum or pearl millet crops in the same region. Perhaps no other crop deserves the title  
65 “orphan” more, because research attention on fonio has been minimal [4].

66 Wild *D. exilis* (sometimes called “hungry rice”) and other West African *Digitaria* have  
67 been harvested by farmers in times of famine throughout recorded history [5], but very little  
68 improvement has been made to the domesticated crop, at least partly evidenced by the fact that  
69 no controlled cross between fonio varieties has been substantiated. Fonio was probably  
70 domesticated in West Africa, presumably before the arrival of pearl millet or sorghum from  
71 Central and East Africa [6], as is suggested by the importance of fonio in Dogon and other  
72 creation myths [4]. Applying the term “domesticated” to fonio cultivars is, however, something  
73 of a stretch. Fonio cultivars do not exhibit the full set of domestication traits, in that they exhibit  
74 the shattering (grain release at maturity) and day-length dependence traits that have been selected  
75 against by early farmers across virtually all cereal crops [7,8]. The selected mutations to non-  
76 shattering and daylength independence are routinely recessive, so the absence of these  
77 agricultural improvements may be an outcome of the polyploid nature of the fonio genome [9].

78 As an orphan crop, fonio has received very little research attention. Over the last 20  
79 years, for instance, only nine refereed publications report any new investigation of any aspect of  
80 fonio biology, although an additional 30 plus publications in that time period investigated fonio  
81 agronomy, cultural significance or nutritional properties [10,11]. In 2007, Adoukonou-Sagbadja  
82 and colleagues [12] published a DNA marker-based analysis of fonio genetic diversity, and there  
83 is some transcript sequence data [13] at NCBI. Beyond this, most fonio investigations have been

84 conducted in West Africa to determine appropriate conditions for subsistence farmers to grow  
85 and/or process the grain from local landraces. In contrast, several other orphan cereal crops of  
86 Africa and Asia have begun to receive extensive attention, including comprehensive analyses of  
87 germplasm resources, even to the extent of full genome sequence analysis. Three of these cereals  
88 with relatively deep recent analyses are, like fonio, panicoid grasses: foxtail millet (*Setaria*  
89 *italica*), pearl millet (*Cenchrus americanus*) and proso millet (*Panicum miliaceum*) [14,15,16].  
90 With these panicoid grass resources, and a comparative genomics strategy [17], it should be  
91 possible to rapidly elevate fonio research to benefit fonio consumers and producers. This  
92 manuscript describes our genomic sequence analysis of the fonio landrace Niatia, and a genetic  
93 comparison of fonio germplasms from across West Africa.

94

#### 95 **Plant material and nucleic acid preparation**

96 Fonio millet (cv. Niatia) seed were obtained from Dr. Sara Patterson (University of Wisconsin,  
97 USA). Niatia is a popular local variety in Mali [18] (see Genetic Diversity for Nagoya protocol).  
98 The seed were multiplied in a University of Georgia greenhouse. Seeds collected from a single  
99 plant were used for all DNA isolation. The seeds were surface sterilized with 8% sodium  
100 hypochlorite (Bioworld, United States) for 10 min, followed by three rinses with sterile distilled  
101 water. The plants were grown in standard potting soil in a greenhouse (with 14 h daylight and  
102 day/night temperatures of 26/20°C). They were watered daily to ~70% soil water holding  
103 capacity. The leaves of four-week-old plants were used for DNA isolation, using a previously  
104 described protocol [19].

105

#### 106 **PacBio SMRT sequencing, sequence polishing and genome assembly**

107 DNA samples were used to construct a PacBio (Pacific Biosciences, Menlo Park, USA) SMRT  
108 sequencing library according to manufacturer recommendations at the University of California at  
109 Davis Genome Center. Fragments >10 kb were selected for sequencing via BluePippen (Sage  
110 Science, LLC, Beverly, USA). A total of 88 Gb of raw PacBio reads from 76 SMRT cells were  
111 passed through the secondary analysis pipeline in SMRT Link (v6.0) and filtered for read quality  
112 higher than 0.75 and length longer than 1 kb. The resultant 75 Gb of filtered reads were  
113 assembled in Canu (v1.8) with the default settings for raw PacBio reads.

114 Racon was used to polish the original assembly for two rounds with the Canu-corrected  
115 PacBio reads. Sequentially, Arrow (VariantCaller v2.3.3) and Pilon (v 1.23) were used to further  
116 polish the assembly with 36 Gb of Illumina paired-end reads obtained on the HiSeq4000 at the  
117 Georgia Genomics and Bioinformatics Core at the University of Georgia.

118 The final assembly (v1.0) has a total length of 760.66 Mb and 3333 contigs, with N50 of  
119 1.73 Mb (L50 of 126) and N90 of 75.85 kb (L90 of 889). The longest contig is 10.17 Mb and the  
120 shortest contig is 1013 bp.

121

### 122 **Estimation of genome size and heterozygosity**

123 Kmer Analysis Toolkits [20] was used to count kmers in Illumina raw reads and to compare the  
124 results with the kmers counted from the genome assembly at several different kmer sizes, from  
125 17-30. These all yielded similar results, but with a somewhat larger genome predicted at  
126 smaller kmer lengths. The distribution of kmer counts was modeled and the heterozygosity level  
127 was estimated using GenomeScope2.0 (<http://qb.cshl.edu/genomescope/genomescope2.0/>).

128 Two distinct peaks were observed in the raw read kmer distribution. We interpret the  
129 peaks at ~50 and ~100 counts/coverage as the two subgenomes in fonio (**Suppl. Fig. S1**).  
130 Genome size estimated from the peaks was 668-707 Mb, depending on the kmer size employed.  
131 This range of values is low compared to previous results from flow cytometry that indicated a  
132 genome size range of 830-1000 Mb for a broad selection of *D. exilis* germplasm [9]. Single  
133 nucleotide variation was estimated to be 4.3% when comparing the A and B genomes in this  
134 tetraploid, but slightly less than 0.01% heterozygosity was observed within either the A or B  
135 genomes, as assayed by kmer allelic ratios. Kmer counts in the assembled genome suggests that  
136 the peak at 100 counts represents common sequences between the two subgenomes, and the  
137 kmers under the peak at 50 counts represent the divergent regions between the two subgenomes.

138

### 139 **Repeat annotation**

140 Repeated sequences were mined and annotated with a combination of *de novo* and homology-  
141 based methods. First, simple sequence repeats were identified and masked with GMATA [21].  
142 Long-Terminal-Repeat-Retrotransposons (LTR-RTs) were identified *de novo* using the  
143 bioinformatic tools LTR\_FINDER [22] and LTRharvest [23] that employ structural criteria to find  
144 intact LTR-RTs, followed by LTR\_retriever analysis [24] to minimize false positives. SINE\_scan

145 (version 1.1.1) [25] was used to find small interspersed nuclear elements (SINEs), a class of  
 146 retroelements, and these were confirmed by manual investigation. Long interspersed nuclear  
 147 elements (LINEs), another class of retroelements, were found with MGEscan-nonLTR (version 2)  
 148 [26]. Small DNA transposable elements (TEs) were found with MITE Tracker [27] and  
 149 HelitronScanner [28] was used to identify the DNA transposons called *Helitrons*. All of the TEs  
 150 from the genome assembly were used to generate a fonio-specific TE library, with individual TE  
 151 families named according to the prevalent current nomenclature system [29]. The fonio TE library  
 152 was compared to the multispecies repeat repository called Repbase [13] to validate annotations,  
 153 and to discover any additional candidate repeats represented in Repbase. Then, the fonio TE library  
 154 was used to identify both full-length and truncated TE elements by homologous search with  
 155 RepeatMasker version 4.0.7 [30] in the genome assembly. Parameter settings were adopted from  
 156 the analysis described in a previous publication [31]. The predicted insertion dates of intact LTR-  
 157 RTs were calculated with LTR\_retriever [24]. The SSRs and TEs were masked by Ns and a TE  
 158 annotation file in GFF3 format was generated for subsequent gene annotation. Types and  
 159 abundances of TEs and other repeats discovered in the fonio genome are presented in **Table 1**.

160

161 **Table 1. Summary of repeat sequence properties in the genome assembly.**

Class	Subclass	Type	Number of families	Number of repeats	Length (Mb)	Percent of genome
Class I TEs:						
Retroelements	LTR-RT	<i>Copia</i>	353	45,194	22.8	3.0
		<i>Gypsy</i>	1223	125,773	153.9	20.2
		Other	824	90,110	57.8	7.6
	LINE	I	17	3040	1.5	0.2
	SINE		3790	181,505	30.6	4.0
Class II TEs:						
DNA Transposons	TIR	CACTA	348	42,737	7.4	1
		<i>Mutator</i>	34	8493	1.8	0.2
		PIF	120	13,973	2.4	0.3
		Tc1	896	124,252	21.5	2.8
		hAT	93	13097	2.5	0.3
	<i>Helitron</i>	<i>Helitron</i>	313	104,271	21.6	2.8
Tandem repeats	SSRs			133,570	5.9	0.8
Unclass. repeats	(Repbase)				48	6.3
Total					329.8	49.7

162

163 **Transcriptome assembly, candidate gene annotation and BUSCO quality assessment**

164 Illumina RNA sequencing data (paired-end 100 bp) of *Digitaria exilis* [13] were downloaded  
165 from the NCBI Sequence Read Archive (accession number SRX1967865) from RNA consisting  
166 of ~80% inflorescence and ~20% leaf tissue. FastQC [32] was used to evaluate data quality, and  
167 low-quality reads and adapter sequences were removed using Trimmomatic [33]. The remaining  
168 reads were aligned to the genome assembly using HISAT2 [34]. The spliced alignments were  
169 used as input for StringTie [35] and assembled into transcripts. TransDecoder, a companion  
170 software of the Trinity platform [36], was used to predict open reading frames.

171 For gene prediction and genome annotation, we used the Maker-P pipeline [37], in  
172 combination with Augustus [38], SNAP [39] and GeneMark [40]. Augustus gene models came  
173 from the BUSCO [41] data set identified during the assembly (see below). GeneMark\_ES was  
174 used to produce *ab initio* gene predictions. Detailed settings for each round of Maker can be  
175 found in the Supplemental Methods. The first round of gene prediction with Maker used the  
176 following inputs: the RNAseq assembly described in the previous section, protein fasta  
177 sequences from *S. bicolor* and *S. italica* [42] as well as the repeat models for *Digitaria*  
178 (described above), and the soft-masked genome assembly. A second round of Maker used as  
179 input the genome file, the annotation produced by the previous round and a SNAP species  
180 parameter/hmm file based on the prior annotation. Finally, the third round of Maker was run  
181 using the following input: the genome assembly, the annotation produced by round two and the  
182 GeneMark models. Functional annotation was done using the accessory scripts of Maker as  
183 described by Campbell and coworkers [43]. Briefly, a BLAST [44] search against the Swissprot  
184 database was used to assign putative functions to the newly annotated gene models, while  
185 InterProScan 5 [45] was used to obtain domain information.

186 Following mapping of RNAseq data with HISAT2, 88% of the RNAseq reads could be  
187 well-aligned to the genome. Transcripts were assembled with Stringtie and ORFs were predicted  
188 with TransDecoder. A total of 58,305 candidate transcripts were obtained, of which 50,389 had  
189 predicted open reading frames.

190 Our first round of Maker predicted 60,300 protein-coding genes (based only on RNA  
191 evidence and protein evidence from sorghum and Setaria). After the 2nd and 3rd round, where  
192 Augustus, SNAP and Genemark-ES models were included, the number of predicted protein

193 coding genes increased to 67,921 and finally to 68,302. We removed 447 candidate genes that  
194 were judged as spurious because they were fragments of otherwise fully assembled genes in the  
195 annotation, so the final number of genes annotated as protein coding genes is 67,855. The  
196 statistics for the gene annotation can be found in the Supplemental Materials (**Table S1**). In total,  
197 88.3% of the gene models were supported by RNAseq data. The Annotation Edit Distance  
198 (AED) measurements indicate how well an annotation agrees with overlapping evidence  
199 (protein, mRNA or EST data). In the fonio assembly, >90% of the gene models have an AED  
200 score less than 0.4%, indicating that gene models are well supported by evidence.

201 BUSCO [41,46] analysis of the filtered predicted protein sequences against the reference  
202 set for plants, on the gVolante platform [47], showed that 97.99% of the BUSCO genes were  
203 found as complete genes, while this representation number increased to 99.31% if partially  
204 covered BUSCO genes were added. A total of 12.4% of the BUSCO genes were single copy,  
205 while 85.6% of the BUSCO genes were found in duplicate. Approximately 1.3% of the BUSCO  
206 genes were fragmented and ~0.7% were missing.

207 A total of 4741 non-coding RNAs (see **Suppl. Table S2**) were predicted with Infernal  
208 [48] by comparing the genome fasta file with the RFAM CM database, version 14.2 [49] using  
209 the protocol described in [50]. Most of these non-coding RNAs were found to be tRNAs, rRNAs  
210 and snoRNAs, as seen in other plant genomes.

211

## 212 **Phylogenetic divergence and dating the most recent whole genome duplication**

213 The coding DNA sequences (CDS) and annotations for *S. bicolor* and *S. italica* were  
214 downloaded from the PLAZA database [42].  $K_s$  distribution analyses were performed using the  
215 wgd package (v1.1) [43]. For each species, the paranome (entire collection of duplicated genes)  
216 was obtained with ‘wgd mcl’ using all-against-all BlastP [43] and MCL clustering [51].  $K_s$   
217 distributions were then constructed using ‘wgd ksd’ with default settings (using MAFFT for  
218 multiple sequence alignment [52], codeml for maximum likelihood estimation of pairwise  
219 synonymous distances [53], and FastTree [54] for inferring phylogenetic trees used in the node  
220 weighting procedure. Anchors or anchor pairs (duplicates lying in collinear or syntenic regions  
221 of the genome) were obtained using i-ADHoRe [55], employing the default settings in ‘wgd

222 syn'.

223 We obtained gene families for a set of nine species in the Poaceae family using  
224 OrthoFinder with default settings [56]. All sequence data were obtained from PLAZA [42]. From  
225 this set of gene families, we identified all gene families that were single-copy in all species but  
226 duplicated in *D. exilis*, and where the *D. exilis* duplicates were anchor pairs (1967 gene families).  
227 For these gene families, we performed pre-alignment homology filtering using PREQUAL [56]  
228 and multiple sequence alignment of the masked amino acid sequences using MAFFT [52]. For  
229 each multiple sequence alignment, we obtained the corresponding codon-level nucleotide  
230 alignment. For each thus-obtained nucleotide alignment, we sampled tree topologies from the  
231 posterior using MrBayes v3.2 [58] under the GTR model with a discrete Gamma mixture for  
232 relative substitution rates across sites (using four classes), sampling every 10 iterations, for a  
233 total of 250,000 iterations. We then identified all gene families for which the expected species  
234 tree topology had posterior probability above 0.9, resulting in a set of 1242 gene families. A  
235 concatenated codon alignment was obtained for these families, which was in three partitions  
236 corresponding to each codon position. We then performed posterior inference of substitution  
237 rates and divergence times for the partitioned alignment using MCMCTree [51, 59] using the  
238 multivariate Normal (MVN) approximation of the likelihood (where the MVN approximation  
239 was based on the maximum likelihood estimates under the GTR model with Gamma distributed  
240 relative rates across sites (5 categories)). We used a Gamma (2, 11) prior for the mean  
241 substitution rate per site per 100 My (million years), based on a rough estimate of the  
242 substitution rate under the molecular clock with a root age of 50 My obtained using baseml from  
243 the PAML package [53]. We use an independent log-normal rates relaxed molecular clock prior  
244 on branch-specific substitution rates, using a Gamma (2, 10) prior for the variance parameter of  
245 the clock. We set the birth-death-sampling prior such that a uniform prior over node ages is  
246 obtained. We include two fossil calibrations. First, we used a minimum age for the *Oryza* -  
247 *Hordeum* divergence of 34 My based on the review of [60]. Next, a secondary calibration for the  
248 root based on previous dating studies included in the Time Tree [61] database was used, where  
249 we excluded all time estimates younger than the 34 My constraint and older than 80 My. We  
250 then fitted a log-normal distribution to the age estimates in the time tree data, which we  
251 approximated by a Gamma (47,100) distribution. We used MCMCTree to obtain 5000 from the  
252 posterior sampling every 200 iterations after a burn-in of 50,000 iterations. We compared two



253 independent runs with each other to verify convergence and with a run of the MCMC algorithm  
254 under the prior alone to compare the posterior distribution for the node ages to the effective prior  
255 implied by the fossil calibrations (**Suppl. Fig. S2**). The results of this analysis provide the  
256 phylogenetic tree shown in **Figure 1D**.

257

### 258 **Transposable element properties**

259 The ~42.6% TE content of the fonio genome is a minimal estimate, given that degraded TE  
260 fragments are often missed by the *de novo* discovery analysis that was employed. This  
261 underestimation is routine in other plant genome annotations as well [62], so it is reasonable to  
262 compare TE descriptions across plant genomes. In fonio, the very high level of *Gypsy* LTR-RTs  
263 compared to *Copia* LTR-RTs is exceptional. Although most grass genomes have more *Gypsy*  
264 TEs than *Copia* (for instance, ~50% *Gypsy* and ~25% *Copia* in the ~2.4 Gb maize genome [63]  
265 or ~36% *Gypsy* and ~33% *Copia* in the ~2.8 Gb pearl millet genome [15], the ~6.7:1 *Gypsy* to  
266 *Copia* ratio in the ~900 Mb fonio genome is unprecedented. One should remember, however,  
267 that the diploid constituent genomes of fonio are ~450 Mb, so somewhat similar results are  
268 observed in other small panicoid genomes like sorghum (~750 Mb) and rice (~430 Mb), with  
269 *Gypsy/Copia* of ~3.7 and ~4.9, respectively [64]. This fonio observation is surprising because the  
270 quantity of *Gypsy* LTR-retrotransposons is the major determinant of genome size in grasses [65],  
271 so one would expect higher *Gypsy* to *Copia* ratios as genome size increases, rather than the  
272 opposite that we observe. These results suggest that either different factors initiate *Gypsy*  
273 amplification bursts than *Copia* amplifications, or that *Copia* elements are particularly sensitive  
274 to shared activation factors. It would be useful to investigate additional *Digitaria* species to see if  
275 this *Gypsy/Copia* ratio trait is shared by other close relatives, and thus a possible outcome of  
276 common ancestral properties.

277         Analysis of LTR-RT insertion dates demonstrated that most of the elements inserted  
278 within the last 2 My. This high level of recent activity is a standard observation in the grasses, at  
279 least partly caused by the fact that the rapid DNA removal by accumulated small deletions  
280 quickly excise and otherwise obscures any DNA that is not under positive selection [66, 67].

281

### 282 **Whole genome duplication and subsequent stability**

283 We inferred whole-paranome and one-vs.-one ortholog  $K_s$  distributions and performed syntenic

284 analyses to further assess the clear signature of a relatively recent whole-genome duplication  
285 (WGD) in *Digitaria exilis*.  $K_S$  distributions present a very clear signature of WGD in the recent  
286 evolutionary past of *D. exilis*, with this event not shared with the closest relative in our analyses  
287 (*S. italica*) (**Figure 1A**). We note that a trace of an older, likely Poaceae-shared WGD [68] event  
288 was also clearly observed in both the whole-paranome and anchor pair  $K_S$  distributions of *D.*  
289 *exilis*, coinciding with similar signatures in sorghum and Setaria (**Figure 1B**). Analysis of co-  
290 linearity and synteny show that the genome of *D. exilis* is still largely conserved in duplicate  
291 (**Figure 1C**). Phylogenetic divergence time estimation (**Figure 1D**) estimated the timing of the  
292 WGD event (or divergence of parental genomes in the case of an allopolyploidy event) at ~3.1  
293 million years ago (mya) with a 95% posterior uncertainty interval of (2.2, 4.2 My) and the  
294 divergence of *Digitaria* from Setaria at 17.8 (12.5, 23.1) mya; with these estimates associated  
295 with a posterior mean substitution rate across the three codon positions of  $2.5 \times 10^{-9}$  ( $1.1 \times 10^{-9}$ ,  
296  $5.0 \times 10^{-9}$ ) substitutions per year per site.

297 It is interesting that **Figure 1C** shows extreme conservation of gene content and order  
298 across long scaffolds, but also the presence of large rearrangements that differentiate  
299 chromosome-size blocks. This suggests a possible selection for major rearrangements after the  
300 polyploids were formed, perhaps to minimize tetrasomic inheritance [69, 70].

301 In the ~3.1 My since the latest WGD, most of the duplicated genes have had both copies  
302 retained. For instance, the BUSCO gene set yielded 81% of the genes still in a duplicated state.  
303 Our genome assemblies did not yield complete chromosomes, so we could not investigate the  
304 details of major chromosomal rearrangements, preferential gene loss (also known as  
305 fractionation), or parent-specific gene expression differences that might differentiate the two  
306 ancestral genomes in this tetraploid [71]. The large stretches of gene content and gene  
307 collinearity retention observed between our largest contiguous assemblies (**Figure 1C**) do  
308 demonstrate, however, that there has been no large number of small rearrangements of these  
309 genomes over the last 3.1 My.

310

311 **Candidate domestication genes**

312 Improvement of fonio will require further domestication, particularly to solve the issues of  
313 shattering and lodging. This process should be greatly assisted by the provision of a  
314 comprehensive genome sequence.

315 In rice, sorghum and maize, mutations in the gene SSH1 (SUPPRESSION of SEED  
316 SHATTERING-1) are associated with panicle retention of the grain after seed maturation (the  
317 “non-shattering trait) in domesticated accessions [72]. Nine sequenced grass genomes were  
318 scanned with OrthoFinder (as described in the section “Phylogenetic divergence and dating the  
319 most recent whole genome duplication”) to find the orthologues of this gene. The gene family  
320 fasta files were used to construct trees using Mafft and Iqtree, trees were visualized in FigTree.  
321 Interproscan was used to annotate the proteins with their pFam domains, and alignments were  
322 visualized in Geneious Prime [73].

323 Fonio has 4 genes related to SSH1, but the phylogenetic tree indicated that two are more  
324 closely related to the rice SSH1 gene associated with shattering than to the other SSH1-like gene  
325 in rice (**Suppl. Fig. S3**). Other species included in our dataset have between 1 and 3 SSH1-like  
326 genes (**Suppl. Table S3**). The extra copies in *D. exilis* are expected because of its polyploid  
327 nature, and thus can explain why no ancient or modern farmers have detected recessive single  
328 gene mutations at each of these loci in a single fonio plant. By modern forward or reverse  
329 mutational techniques, inactivations of both of these genes should be targeted in order to solve  
330 the shattering problem in fonio.

331 Inactivation of the *dw3* (Dwarfing-3) genes of sorghum is responsible for the semi-dwarf  
332 trait that diminishes lodging and thereby greatly improves yield and input response in this  
333 important crop of arid and semi-arid agriculture [74]. Inactivational mutations of orthologues of  
334 the same gene are also responsible for the pearl millet cultivars with highest lodging resistance  
335 and the highest grain yield [75]. Hence, orthologues of *dw3* also should be targets for  
336 inactivational mutation and molecular breeding in fonio. Once again, fonio has more copies of  
337 this gene than do any of the other grasses screened, all of which are diploids (**Suppl. Fig. S4 and**  
338 **Table S4**).

339 The GW2 (GRAIN WEIGHT-2) gene controls seed weight in wheat and rice, with  
340 inactivation of the gene leading to larger grain [76,77]. Orthofinder results indicated that  
341 members of this gene family are present in single copy in all of the examined grass species,  
342 except fonio and maize (**Suppl. Fig. S5 and Table S5**). The two copies in *D. exilis* only differ

343 from each other by 3 amino acid residue substitutions. The fonio genes were found to be nearly  
344 identical to the unmutated GW2 version that yields smaller grain in rice and wheat (data not  
345 shown). Although increased seed weight does not always increase yield (due to correlated traits,  
346 like seed number), it is particularly important trait in fonio to enable sowing for uniform stands  
347 and mechanical threshing.

348

### 349 **Genetic diversity**

350 Fonio genetic diversity was assessed using 184 samples from ~130 accessions collected  
351 from Mali and Niger, signatories to the Cartagena Protocol on Biosafety (Suppl. Table S6).  
352 Consistent with the Nagoya Protocol, fonio materials from Mali were collected in Mali by  
353 Institut d’Economie Rural (IER) while those from Niger were collected in Niger by Institute  
354 National de Recherche Agronomique du Niger (INRAN) and conserved at the ICRISAT Niamey  
355 genebank. Authors Sanogo, Hamidou and Gangashetty were involved in the germplasm  
356 collection, seed conservation at the genebank and/or DNA extraction from young seedlings. All  
357 DNA samples were sent to the USA for analysis for research purposes only. This research has no  
358 commercial application.

359 Seedlings of each sample were grown at the respective institutions in West Africa, and  
360 DNA was extracted from young leaves with a QIAGEN DNeasy Plant Mini Kit (Germantown,  
361 USA). Lyophilized DNA was then sent to Data2Bio (Ames, USA) for tunable genotyping-by-  
362 sequencing (tGBS) using 2-bp selection and 5 runs on an Ion Torrent Ion Proton Instrument  
363 (Thermo Fisher Scientific, Waltham, USA). The resulting raw sequences were quality-trimmed  
364 by Data2Bio, which removed bases with PHRED quality scores <15. These trimmed sequences  
365 were then aligned to the genome assembly with GSNAP v2020-04-08 [78] using default  
366 parameters. SNPs were called using the bcftools mpileup command v1.9 [79] with max-depth set  
367 to 1000 and minimum base quality set to 20. Raw SNPs were then filtered using TASSEL  
368 v5.2.40 [80], custom R scripts with R v3.5.1 [81], and bcftools to include only sites with  $\leq 25\%$   
369 heterozygosity,  $\leq 500$  total read depth,  $\leq 60\%$  missing data, and  $\geq 2.5\%$  minor allele frequency  
370 (Suppl. Table S7). Population substructure was determined with fastStructure v1.0 [82], testing  
371 from 1 to 10 population clusters and identifying the optimal number with the included  
372 chooseK.py program. This identified 3 clear clusters of material, with genetic separation strongly  
373 correlated with geography (**Figure 2A**). The genetic distinctions among these clusters are clear

374 when plotting the genetic principal coordinates and relationship dendrogram (**Figure 2B**). A  
375 small number of accessions (<5) appear “misplaced” on the geographic map, which could be due  
376 to recent transfer of germplasm or human error during collection, storage, or processing.

377 Principal coordinates were calculated by using classical multidimensional scaling (R  
378 function `cmdscale()`) on a genetic distance matrix calculated in TASSEL (option –  
379 `distanceMatrix`). The same distance matrix was used to create the dendrogram by neighbor-  
380 joining (function `nj()`) with the R package `app` v5.3 [83]. Accessions were plotted geographically  
381 using the R package `ggmap` v3.0.0 [81]. Additional software used in this analysis included  
382 `samtools` v0.1.19-96b5f2294a [84], `conda` 4.8.3 [85], `PLINK` v1.90b5.2 [86] and the R packages  
383 `argparse` v2.0.1 [87], `ggplot2` v3.2.1 [88], `gridExtra` v2.3 [89], and `RColorBrewer` v1.1.2 [90].

384

## 385 **Conclusions**

386

387 Genome analysis of any polyploid is challenging, especially when no diploid ancestors are  
388 known. Our sequence of the white fonio (*D. exilis*) genome indicates its recent tetraploid origin  
389 and the retention of most of the genes duplicated in this process. This retention of duplicated  
390 genes likely explains why recessive mutations for important agronomic traits like shattering, day  
391 length dependence and semi-dwarfism have not yet been detected in fonio. However, it is now  
392 possible to identify such mutations by using modern mutation detection schemes, like those used  
393 for the tetraploid cereal *Eragrostis tef* [91]. One purpose for generating a fonio genome sequence  
394 was to attract molecular genetics researchers into the study of this crop, and thereby enable  
395 hypothesis-driven breeding through genomics-assisted selection. If future researchers develop a  
396 transformation technology for fonio [92] or develop other genome editing strategies [93], then  
397 directed mutagenesis could be used to knock out pairs of these domestication genes in a single  
398 step [94].

399 The importance of correcting such problems as shattering, seed size, lodging in fonio  
400 cannot be over-estimated. Until shattering is solved, farmers will continue to be required to  
401 harvest before grains fully mature, thus dramatically decreasing overall yield. Without semi-  
402 dwarf varieties, already serious lodging problems in fonio will continue to prohibit the use of  
403 more inputs (because fertilizer increases plant height and thus lodging) or even the selection of  
404 larger grain yield from the panicles, because greater weight on the top of the plant can cause

405 more lodging. The same will almost certainly be true for fonio, hence providing a partial  
406 explanation for its tiny seed size in cultivated landraces. With domestication traits fully penetrant  
407 into fonio cultivars, one can expect dramatic increases in fonio performance, with expectations  
408 of a two-fold or greater yield enhancement easily within the short-term range of possibilities.

409 The absence of an outcrossing protocol for fonio is another technical deficiency that  
410 severely limits this crop's potential for improvement. Our diversity analysis on cultivar Niatia  
411 indicates <0.01% heterozygosity, showing that crosses occur very rarely by natural processes.  
412 Hence, generating controlled crosses will probably require a serious dedication to this pursuit.  
413 Our results indicate a great deal of genetic variability within fonio landraces, so we have no  
414 doubt that hybridization could be used in breeding projects to optimize fonio germplasm quality  
415 for future W. African and other farmers.

416

#### 417 **Availability of Supporting Data**

418

419 The genome and annotation can be accessed on the AOCC-specific branch of the  
420 ORCAE platform [95,96] at: <https://bioinformatics.psb.ugent.be/orcae/aocc/overview/Digex>).  
421 The GenBank project number for the assembly is PRJNA640067. All scripts for diversity  
422 analysis and data tables are available at <https://github.com/wallacelab/paper-fonio-diversity-2020>  
423 including full genotyping table. Genotyping table also available at GenBank Project number  
424 PRJNA644458.

425

426

427

#### 428 **Abbreviations**

429

430 Dw3: dwarf3; Gb: gigabase; GW2: grain weight2; LINE: long interspersed nuclear element;  
431 LTR: long terminal repeat; LTR-RT: long terminal repeat retrotransposon; Mb: megabase;  
432 MITE: miniature inverted repeat transposable element; My: million years; mya: million years  
433 ago; NCBI: National Center for Biotechnology Information; ORF: open reading frame; SINE:  
434 small interspersed nuclear element; SMRT: single molecule, real time sequencing; SSH1:

435 suppression of shattering1; SSR: simple sequence repeat; TE: transposable element; TIR:  
436 terminal inverted repeat transposable element; Unclass: unclassified repeat.

437

#### 438 **Conflict of Interest**

439

440 The authors declare that they have no competing interests.

441

#### 442 **Consent for publication**

443 Not Applicable

444

#### 445 **Funding**

446

447 JLB acknowledges the Giles Fellowship from the University of Georgia as a source of funding for  
448 this project. YVdP acknowledges funding from the European Research Council (ERC) under the  
449 European Union's Horizon 2020 research and innovation program (grant agreement No 833522).  
450 AV acknowledges funding from the Seed Biotechnology Center, University of California, USA.  
451 JGW acknowledges funding from the International Crops Research Institute for the Semi-Arid  
452 Tropics (ICRISAT) and the University of Georgia. MDS acknowledges funding from the  
453 McKnight foundation.

454

#### 455 **Author Contributions**

456

457 J.L.B., J.W., Y.vdP., and A.V.D. conceived, designed and interpreted the study; S.C., X.M., X.  
458 W., A.E.J.Y., S.R.C., M.S.J., P.G., F.H., M.D.S., and A.Z. prepared the materials, conducted the  
459 experiments, and analyzed all data; J.L.B. led on manuscript preparation, while all other authors  
460 revised the manuscript and approved the final version.

461

#### 462 **Acknowledgements:**

463

464 Shu-Min Kao for providing scripts and Sara Patterson for providing *Niatia* seed and for helpful  
465 discussions. Oanh Nguyen from the UC Davis Genome and Biomedical Sciences Facility for technical  
466 expertise for sequencing with Pacific Biosciences. Armando Garcia-Llanos for DNA and library quality  
467 control for sequencing.

468

## 469 **Figure Legends**

470

471 **Figure 1:** Whole genome duplication and polyploidy analysis. (A)  $K_S$  estimation of age  
472 distribution for paralogs and orthologs of white fonio (*Digitaria*) and some close relatives. The  
473 distribution in light pink represents the entire white fonio paranome, while the distribution in  
474 darker pink represents the anchor points (duplicated genes lying in syntenic or colinear regions  
475 (see C)). Distributions in black, dark green and light green represent the one-vs.-one ortholog  
476 comparisons between *Digitaria-Setaria*, *Digitaria-Sorghum* and *Sorghum-Setaria*, respectively.  
477 (B)  $K_S$  distributions for paralogs of white fonio, sorghum and *Setaria* (zoom in), showing an  
478 older, likely Poaceae-shared, WGD. (C) Syntenic relationships between putative homoeologous  
479 contigs, with colored lines connecting homoeologous gene pairs in the white fonio genome  
480 assembly. (D) Time-calibrated phylogenetic tree of several major Poaceae lineages, including  
481 white fonio, based on 1242 gene families consisting of a single gene copy in each lineage and an  
482 anchor pair (A and B) in *Digitaria*. The time scale is shown in million years (My). See text for  
483 details.

484

485 **Figure 2 – Fonio Genetic Diversity.** The genetic diversity of fonio samples was surveyed by  
486 genotyping-by-sequencing. (A) Fonio samples originated from Mali and Niger. They separate  
487 into 3 primary subpopulations based on population structure analysis. Both principal coordinate  
488 analysis of the genetic diversity (B) and a neighbor-joining tree of the population (C) confirm  
489 these groupings. A few discrepancies between population assignment and geography may be due  
490 to recent long-distance germplasm exchanges or labelling errors during collection and storage.

491

## 492 **References**

- 493 1. Temple VJ and Bassa JD. Proximate chemical composition of acha (*Digitaria exilis*) grain.  
494 Journal of the Science of Food and Agriculture. 1991;**56**(4):561-3.
- 495 2. Fanou N, Hulshof P, Koreissi Y and Brouwer I. Nutritive values of fonio and fonio



- 496 products. *Annals of Nutrition and Metabolism*. 2009;**55**:110-18.
- 497 3. Ballogou V, Soumanou M, Toukourou F and Hounhouigan J. Structure and nutritional  
498 composition of fonio (*Digitaria exilis*) grains: a review. *International Research Journal of*  
499 *Biological Sciences*. 2013;**2**(1):73-9.
- 500 4. Vietmeyer N, Borlaugh N, Axtell J, Burton G, Harlan J and Rachie K. Fonio (Acha).  
501 NRC/BSTID ed *Lost Crops of Africa*. 1996;**1**:59-76.
- 502 5. Council NR. *Lost Crops of Africa: volume I: grains*. National Academies Press; 1996.
- 503 6. De Wet J. The three phases of cereal domestication. In: Chapman GP, editor. *Grass*  
504 *Evolution and Domestication*. London: Cambridge University Press; 1992. p. 176-98.
- 505 7. Aliero A and Morakinyo J. Photoperiodism in *Digitaria exilis* (Kipp) Stapf accessions.  
506 *African Journal of Biotechnology*. 2005; **4**(3):241-3.
- 507 8. Patterson SE, Bolivar-Medina JL, Falbel TG, Hedtcke JL, Nevarez-McBride D, Maule AF,  
508 et al. Are we on the right track: can our understanding of abscission in model systems  
509 promote or derail making improvements in less studied crops? *Frontiers in Plant Science*.  
510 2016;**6**:1268.
- 511 9. Adoukonou-Sagbadja H, Schubert V, Dansi A, Jovtchev G, Meister A, Pistrick K, et al.  
512 Flow cytometric analysis reveals different nuclear DNA contents in cultivated Fonio  
513 (*Digitaria* spp.) and some wild relatives from West-Africa. *Plant Systematics and*  
514 *Evolution*. 2007;**267**:163-76.
- 515 10. Gigou J, Stilmant D, Diallo TA, Cisse N, Sanogo MD, Vaksman M, et al. Fonio millet  
516 (*Digitaria exilis*) response to N, P and K fertilizers under varying climatic conditions in  
517 West Africa. *Experimental Agriculture*. 2009;**45**(4):401-15.
- 518 11. Chukwurah PN, Uyoh EA, Usen IN, Ekerette EE and Ogbonna NC. Assessment of intra  
519 and inter species variation in antioxidant composition and activity in marginalized Fonio  
520 millet (*Digitaria* spp.). *Journal of Cereals and Oilseeds*. 2016;**7**(1):1-6.
- 521 12. Adoukonou-Sagbadja H, Wagner C, Dansi A, Ahlemeyer J, Daïnou O, Akpagana K, et al.  
522 Genetic diversity and population differentiation of traditional fonio millet (*Digitaria* spp.)  
523 landraces from different agro-ecological zones of West Africa. *Theoretical and Applied*  
524 *Genetics*. 2007;**115**(7):917-31.
- 525 13. Jurka, J., V.V. Kapitonov, A. Pavlicek, P. Klonowski, O. Kohany, and J. Walichiewicz.  
526 2005. Repbase Update, a database of eukaryotic repetitive elements. *Cytogenetic and*  
527 *genome research* **110**:462-467.
- 528 14. Bennetzen JL, Schmutz J, Wang H, Percifield R, Hawkins J, Pontaroli AC, et al. Reference  
529 genome sequence of the model plant *Setaria*. *Nature Biotechnology*. 2012;**30**(6):555-61.
- 530 15. Varshney RK, Shi C, Thudi M, Mariac C, Wallace J, Qi P, et al. Pearl millet genome  
531 sequence provides a resource to improve agronomic traits in arid environments. *Nature*  
532 *Biotechnology*. 2017;**35**(10):969-76.
- 533 16. Zou C, Li L, Miki D, Li D, Tang Q, Xiao L, et al. The genome of broomcorn millet. *Nature*  
534 *Communications*. 2019;**10**(1):1-11.
- 535 17. Bennetzen JL and Freeling M. The unified grass genome: synergy in synteny. *Genome*  
536 *Research*. 1997;**7**(4):301-6.
- 537 18. Cruz J-F. Fonio. Upgrading quality and competitiveness of fonio for improved livelihoods  
538 in West Africa: Second activity report. Montpellier: CIRAD; 2008.
- 539 19. Murray MG and Thompson WF (1980) Rapid isolation of high molecular weight plant  
540 DNA. *Nucleic Acids Res* **8**: 4321-4325.
- 541 20. Mapleson D, Garcia Accinelli G, Kettleborough G, Wright J and Clavijo BJ. KAT: a K-

- mer analysis toolkit to quality control NGS datasets and genome assemblies. *Bioinformatics*. 2017;**33**(4):574-6.
- 544 21. Wang X and Wang L. GMATA: an integrated software package for genome-scale SSR  
545 mining, marker development and viewing. *Frontiers in Plant Science*. 2016;**7**:1350.
- 546 22. Xu Z and Wang H. LTR\_FINDER: an efficient tool for the prediction of full-length LTR  
547 retrotransposons. *Nucleic Acids Research*. 2007;**35**(Web Server issue):W265-8.
- 548 23. Ellinghaus D, Kurtz S and Willhoeft U. LTRharvest, an efficient and flexible software for  
549 *de novo* detection of LTR retrotransposons. *BMC Bioinformatics*. 2008;**9**(1):18.
- 550 24. Ou S and Jiang N. LTR\_retriever: a highly accurate and sensitive program for identification  
551 of long terminal repeat retrotransposons. *Plant Physiology*. 2018;**176**(2):1410-22.
- 552 25. Mao H and Wang H. SINE\_scan: an efficient tool to discover short interspersed nuclear  
553 elements (SINEs) in large-scale genomic datasets. *Bioinformatics*. 2017;**33**(5):743-5.
- 554 26. Rho M and Tang H. MGEScan-non-LTR: computational identification and classification  
555 of autonomous non-LTR retrotransposons in eukaryotic genomes. *Nucleic Acids Research*.  
556 2009;**27**(21):e143.
- 557 27. Crescente JM, Zavallo D, Helguera M and Vanzetti LS. MITE Tracker: an accurate  
558 approach to identify miniature inverted-repeat transposable elements in large genomes.  
559 *BMC Bioinformatics*. 2018;**19**(1):348.
- 560 28. Xiong W, He L, Lai J, Dooner HK and Du C. HelitronScanner uncovers a large overlooked  
561 cache of Helitron transposons in many plant genomes. *Proceedings of the National  
562 Academy of Sciences*. 2014;**111**(28):10263-8.
- 563 29. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, et al. A unified  
564 classification system for eukaryotic transposable elements. *Nature Reviews Genetics*.  
565 2007;**8**(12):973-82.
- 566 30. RepeatMasker. <http://www.repeatmasker.org>, version 4.0.7. Accessed March 15, 2020  
567
- 568 31. Luo M-C, Gu YQ, Puiu D, Wang H, Twardziok SO, Deal KR, et al. Genome sequence of  
569 the progenitor of the wheat D genome *Aegilops tauschii*. *Nature*. 2017;**551** (7681):498-  
570 502.
- 571 32. Andrews S. Babraham bioinformatics-FastQC a quality control tool for high throughput  
572 sequence data. URL: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.  
573 Accessed 19 June 2020.
- 574 33. Bolger AM, Lohse M and Usadel B. Trimmomatic: a flexible trimmer for Illumina  
575 sequence data. *Bioinformatics*. 2014;**30**(15):2114-20.
- 576 34. Kim D, Langmead B and Salzberg SL. HISAT: a fast spliced aligner with low memory  
577 requirements. *Nature Methods*. 2015;**12**(4):357-60.
- 578 35. Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT and Salzberg SL. StringTie  
579 enables improved reconstruction of a transcriptome from RNA-seq reads. *Nature  
580 Biotechnology*. 2015;**33**(3):290.
- 581 36. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. *De novo*  
582 transcript sequence reconstruction from RNA-seq using the Trinity platform for reference  
583 generation and analysis. *Nature Protocols*. 2013;**8**(8):1494.
- 584 37. Campbell MS, Law M, Holt C, Stein JC, Moghe GD, Hufnagel DE, et al. MAKER-P: a  
585 tool kit for the rapid creation, management, and quality control of plant genome  
586 annotations. *Plant Physiology*. 2014;**164**(2):513-24.
- 587 38. Stanke M, Steinkamp R, Waack S and Morgenstern B. AUGUSTUS: a web server for gene

- 588 finding in eukaryotes. *Nucleic Acids Research*. 2004;**32** (Web Server issue):W309-W12.
- 589 39. Korf I. Gene finding in novel genomes. *BMC Bioinformatics*. 2004;**5**(1):59.
- 590 40. Borodovsky M, Lomsadze A. Eukaryotic gene prediction using GeneMark.hmm-E and  
591 GeneMark-ES. *Current Protocols in Bioinformatics*. 2011;**4**:4.6.1-10.
- 592 41. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO:  
593 assessing genome assembly and annotation completeness with single-copy orthologs.  
594 *Bioinformatics*. 2015;**31**(19):3210-2.
- 595 42. Van Bel M, Diels T, Vancaester E, Kreft L, Botzki A, Van de Peer Y, et al. PLAZA 4.0:  
596 an integrative resource for functional, evolutionary and comparative plant genomics.  
597 *Nucleic Acids Research*. 2018;**46** (D1):D1190-6.
- 598 43. Campbell MS, Holt C, Moore B and Yandell M. Genome annotation and curation using  
599 MAKER and MAKER-P. *Current Protocols in Bioinformatics*. 2014;**48** (1):4.11.1-39.
- 600 44. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST  
601 and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids*  
602 *Research*. 1997;**25**(17):3389-402.
- 603 45. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-  
604 scale protein function classification. *Bioinformatics*. 2014;**30** (9):1236-40.
- 605 46. Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G, et al.  
606 BUSCO applications from quality assessments to gene prediction and phylogenomics.  
607 *Molecular Biology and Evolution*. 2018;**35**(3):543-8.
- 608 47. Nishimura O, Hara Y and Kuraku S. gVolante for standardizing completeness assessment  
609 of genome and transcriptome assemblies. *Bioinformatics*. 2017;**33**(22):3635-7.
- 610 48. Nawrocki EP and Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches.  
611 *Bioinformatics*. 2013;**29** (22):2933-5.
- 612 49. Kalvari I, Argasinska J, Quinones-Olvera N, Nawrocki EP, Rivas E, Eddy SR, et al. Rfam  
613 13.0: shifting to a genome-centric resource for non-coding RNA families. *Nucleic Acids*  
614 *Research*. 2018;**46**(D1):D335-42.
- 615 50. Kalvari I, Nawrocki EP, Argasinska J, Quinones-Olvera N, Finn RD, Bateman A, et al.  
616 Non-coding RNA analysis using the Rfam database. *Current Protocols in Bioinformatics*.  
617 2018;**62**(1):e51. doi: 10.1002/cpbi.51.
- 618 51. Dongen SV. *Graph clustering by flow simulation*. PhD thesis, University of Utrecht  
619 Amsterdam, Netherlands, 2000.
- 620 52. Katoh K and Standley DM. MAFFT multiple sequence alignment software version 7:  
621 improvements in performance and usability. *Molecular Biology and Evolution*. 2013;**30**(4)  
622 772-80.
- 623 53. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular Biology and*  
624 *Evolution*. 2007;**24**(8):1586-91.
- 625 54. Price MN, Dehal PS and Arkin AP. FastTree 2—approximately maximum-likelihood trees  
626 for large alignments. *PLoS One*. 2010;**5**(3):e9490.
- 627 55. Proost S, Fostier J, De Witte D, Dhoedt B, Demeester P, Van de Peer Y, et al. i-ADHoRe  
628 3.0—fast and sensitive detection of genomic homology in extremely large data sets.  
629 *Nucleic Acids Research*. 2012;**40** (2):e11 doi:10.1093/nar/gkr955.
- 630 56. Emms DM and Kelly S. OrthoFinder: phylogenetic orthology inference for comparative  
631 genomics. *Genome Biology*. 2019;**20**(1):238.
- 632 57. Whelan S, Irisarri I and Burki F. PREQUAL: detecting non-homologous characters in sets  
633 of unaligned homologous sequences. *Bioinformatics*. 2018;**34**(22):3929-30.

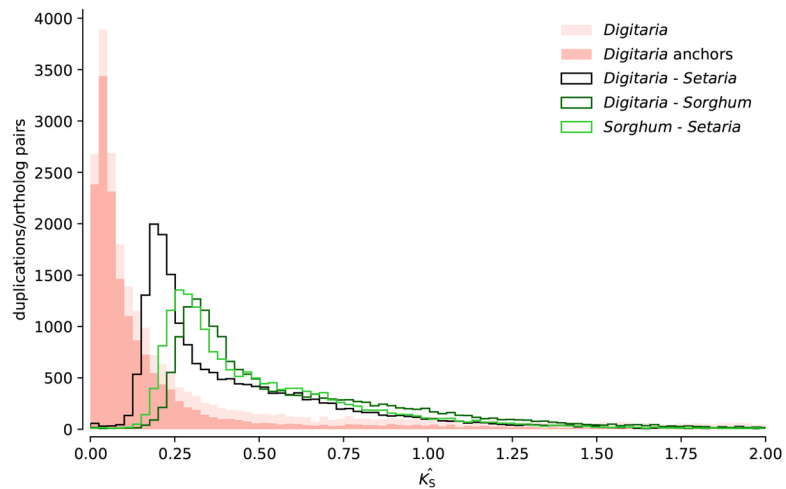
- 634 58. Ronquist F, Teslenko M, Van Der Mark P, Ayres DL, Darling A, Höhna S, et al. MrBayes  
635 3.2: efficient Bayesian phylogenetic inference and model choice across a large model  
636 space. *Systematic Biology*. 2012;**61**(3):539-42.
- 637 59. Rannala B and Yang Z. Inferring speciation times under an episodic molecular clock.  
638 *Systematic Biology*. 2007;**56** (3):453-66.
- 639 60. Iles WJD, Smith SY, Gandolfo MA and Graham SW. Monocot fossils suitable for  
640 molecular dating analyses. *Bot J Linn Soc*. 2015;**178** 3:346-74. doi:10.1111/boj.12233.
- 641 61. Kumar S, Stecher G, Suleski M and Hedges SB. TimeTree: a resource for timelines,  
642 timetrees, and divergence times. *Molecular Biology and Evolution*. 2017;**34**(7):1812-9.
- 643 62. Bennetzen JL and Park M. Distinguishing friends, foes, and freeloaders in giant genomes.  
644 *Current Opinion in Genetics & Development*. 2018;**49**:49-55.
- 645 63. Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, et al. The B73 maize  
646 genome: complexity, diversity, and dynamics. *Science*. 2009;**326**(5956):1112-5.
- 647 64. Paterson AH, Bowers JE, Bruggmann R, Dubchak I, Grimwood J, Gundlach H, et al. The  
648 *Sorghum bicolor* genome and the diversification of grasses. *Nature*. 2009;**457**(7229):551-  
649 6.
- 650 65. Bennetzen JL and Wang H. The contributions of transposable elements to the structure,  
651 function, and evolution of plant genomes. *Annual Review of Plant Biology*. 2014;**65**:505-  
652 30.
- 653 66. Devos KM, Brown JKM and Bennetzen JL. Genome size reduction through illegitimate  
654 recombination counteracts genome expansion in *Arabidopsis*. *Genome Research*.  
655 2002;**12**:1075-9.
- 656 67. Ma J, Devos KM and Bennetzen JL. Analyses of LTR-retrotransposon structures reveal  
657 recent and rapid genomic DNA loss in rice. *Genome Research*. 2004;**14**(5):860-9.
- 658 68. Jiao Y, Li J, Tang H and Paterson AH. Integrated syntenic and phylogenomic analyses  
659 reveal an ancient genome duplication in monocots. *Plant Cell*. 2014;**26**(7):2792-802.
- 660 69. Sybenga J. Allopolyploidization of autopolyploids I. Possibilities and limitations.  
661 *Euphytica*. 1969;**18**:355-371.
- 662 70. Soltis DE and Soltis PS. Molecular data and the dynamic nature of polyploidy. *Critical*  
663 *Reviews in Plant Sciences*. 1993;**12**:243-273.
- 664 71. Bird KA, VanBuren R, Puzey JR and Edger PP. The causes and consequences of  
665 subgenome dominance in hybrids and recent polyploids. *New Phytologist*.  
666 2018;**220**(1):87-93.
- 667 72. Lin Z, Li X, Shannon LM, Yeh CT, Wang ML, Bai G, et al. Parallel domestication of the  
668 *Shattering1* genes in cereals. *Nature Genetics*. 2012;**44**(6):720-4.
- 669 73. Geneious Prime 2020. <https://www.geneious.com>.
- 670 74. Multani DS, Briggs SP, Chamberlin MA, Blakeslee JJ, Murphy AS and Johal GS. Loss  
671 of an MDR transporter in compact stalks of maize br2 and sorghum dw3 mutants.  
672 *Science*. 2003;**302**(5642):81-4.
- 673 75. Parvathaneni RK, Jakkula V, Padi FK, Faure S, Nagarajappa N, Pontaroli AC, et al. Fine-  
674 mapping and identification of a candidate gene underlying the *d2* dwarfing phenotype in  
675 pearl millet, *Cenchrus americanus* (L.) Morrone. G3. 2013;**3**(3):563-72.
- 676 76. Song S-L, Huang W, Shi M, Zhu M-Z and Lin H-X. A QTL for rice grain width and weight  
677 encodes a previously unknown RING-type ubiquitin ligase. *Nature Genetics*. 2007;**39**:623-  
678 30.
- 679 77. Simmonds J, Scott P, Brinton J, Mestre TC, Bush M, del Blanco A, Dubcovsky J and Uauy

- 680 C. A splice acceptor site mutation in TaGW2-A1 increases thousand grain weight in  
681 tetraploid and hexaploidy wheat through wider and longer grains. *Theoretical and Applied*  
682 *Genetics*. 2016;**129**:1099-112.
- 683 78. Wu TD, Reeder J, Lawrence M, Becker G, Brauer MJ. GMAP and GSNAP for Genomic  
684 Sequence Alignment: Enhancements to Speed, Accuracy, and Functionality. *Methods in*  
685 *Molecular Biology*. 2016;**1418**:283-334.
- 686 79. Li H. A statistical framework for SNP calling, mutation discovery, association mapping  
687 and population genetical parameter estimation from sequencing data. *Bioinformatics*.  
688 2011;**27** (21):2987-93.
- 689 80. Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y and Buckler ES. TASSEL:  
690 software for association mapping of complex traits in diverse samples. *Bioinformatics*.  
691 2007;**23**(19):2633-5.
- 692 81. Kahle D and Wickham H. ggmap: Spatial Visualization with ggplot2. *The R Journal*.  
693 2013;**5**(1):144-61.
- 694 82. Raj A, Stephens M and Pritchard JK. fastSTRUCTURE: variational inference of  
695 population structure in large SNP data sets. *Genetics*. 2014;**197**(2):573-89.
- 696 83. Paradis E and Schliep K. ape 5.0: an environment for modern phylogenetics and  
697 evolutionary analyses in R. *Bioinformatics*. 2019;**35**(3):526-8.
- 698 84. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence  
699 alignment/map format and SAMtools. *Bioinformatics*. 2009;**25** (16):2078-9.
- 700 85. Analytics C. Anaconda Software Distribution. Computer software. Vers. 2-2.4.0.  
701 <https://anaconda.com>. Accessed 19 June 2020.
- 702 86. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a  
703 tool set for whole-genome association and population-based linkage analyses. *The*  
704 *American Journal of Human Genetics*. 2007;**81**(3):559-75.
- 705 87. Davis TL. argparse: Command Line Optional and Positional Argument Parser. R package  
706 version 2.0.1. <https://CRAN.R-project.org/package=argparse>. Accessed 19 June 2020.
- 707 88. Wickham H. Ggplot2: elegant graphics for data analysis. 3rd ed. Springer; 2016.
- 708 89. Auguie B. gridExtra: miscellaneous functions for “grid” graphics. R package version. Vers,  
709 2.3. <https://CRAN.R-project.org/package=gridExtra>. Accessed 19 June 2020.
- 710 90. Neuwirth E and Brewer RC. ColorBrewer palettes. R package version. Vers, 1.1-2.  
711 <https://CRAN.R-project.org/package=RColorBrewer>. Accessed 19 June 2020.
- 712 91. Zhu Q, Smith SM, Ayele M, Yang L, Jogi A, Chaluvadi SR, et al. High-throughput  
713 discovery of mutations in tef semi-dwarfing genes by next-generation sequencing analysis.  
714 *Genetics*. 2012;**192**(3):819-29.
- 715 92. Ntui VO, Azadi P, Supaporn H and Mii M. Plant regeneration from stem segment-derived  
716 friable callus of “Fonio” (*Digitaria exilis* (L.) Stapf.). *Scientia Horticulturae*.  
717 2010;**125**(3):494-9.
- 718 93. Ji X, Yang B, and Wang D. Achieving plant genome editing while bypassing tissue culture.  
719 *Trends in Plant Science*. 2020;**25**(5):427-9.
- 720 94. Hu N, Xian Z, Li N, Liu Y, Huang W, Yan F, et al. Rapid and user-friendly open-source  
721 CRISPR/Cas9 system for single- or multi-site editing of tomato genome. *Horticulture*  
722 *Research*. 2019;**6**:7.
- 723 95. Sterck L, Billiau K, Abeel T, Rouzé P, Van de Peer Y. ORCAE: online resource for  
724 community annotation of eukaryotes. *Nature Methods*. 2012;**9**(11):1041.
- 725 96. Yssel, A.E.J.; Kao, S.-M.; Van de Peer, Y.; Sterck, L. ORCAE-AOCC: A centralized portal

726  
727

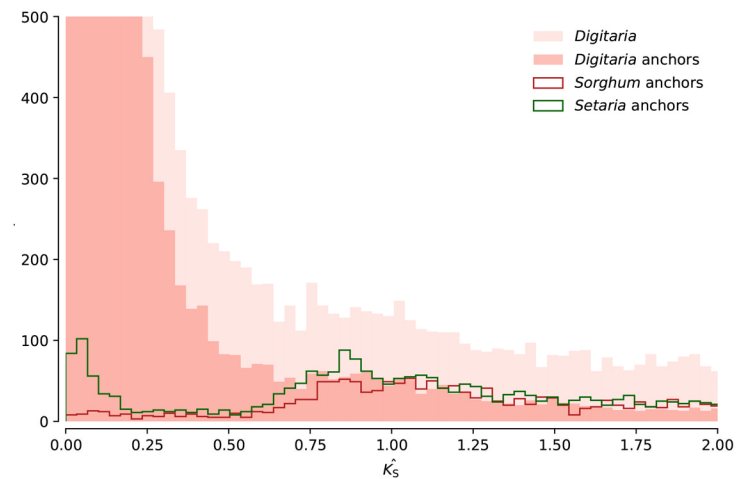
for the annotation of african orphan crop genomes. *Genes* 2019; **10**(12): 950.

Figure 1

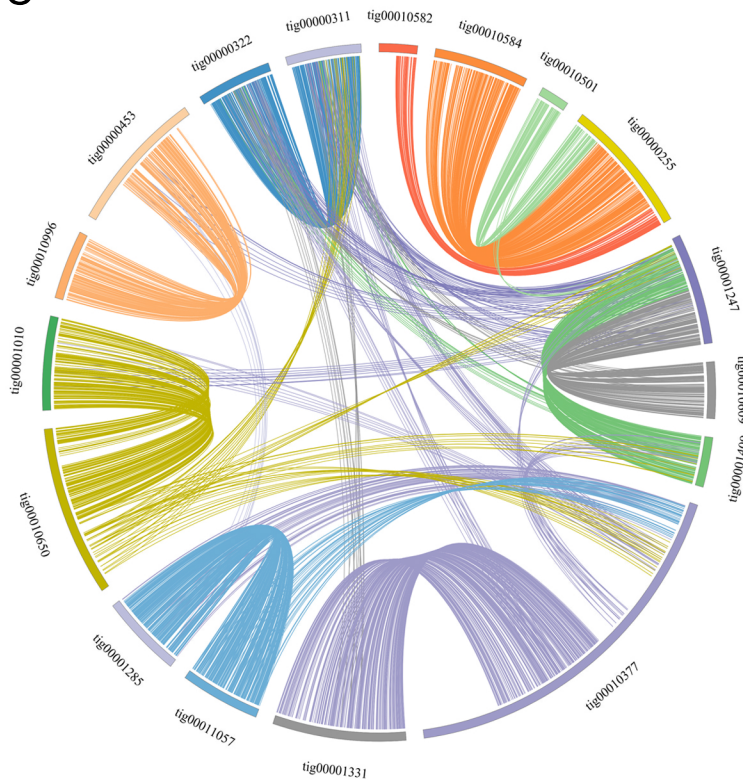


B

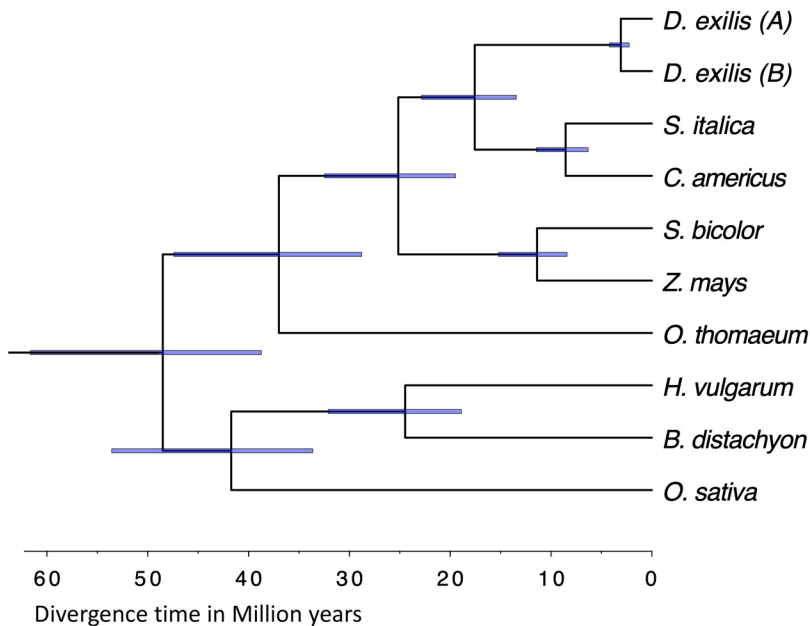
[Click here to access/download;Figure;Figure 1.pdf](#)



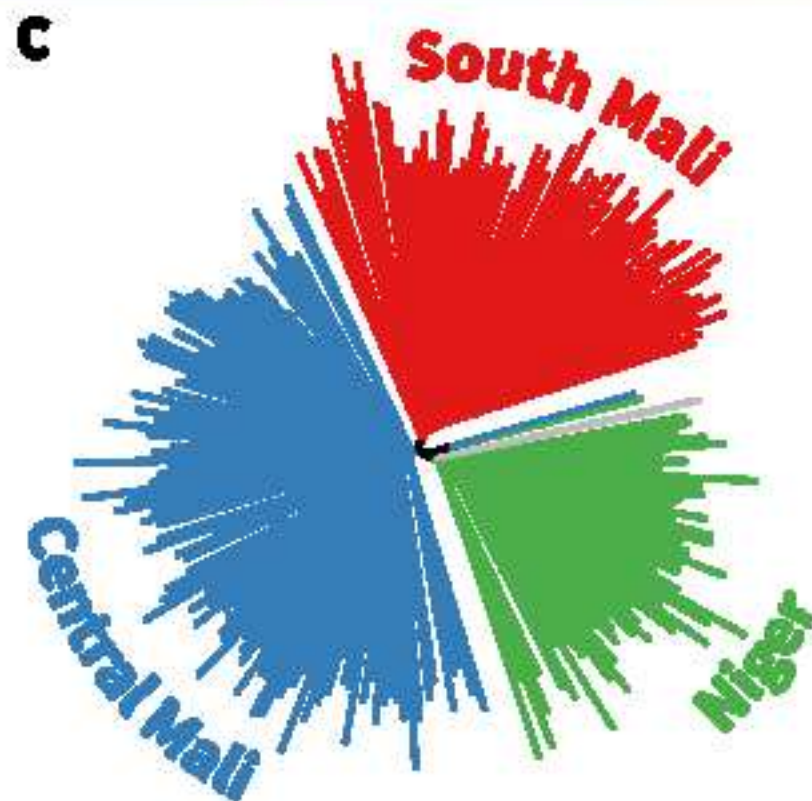
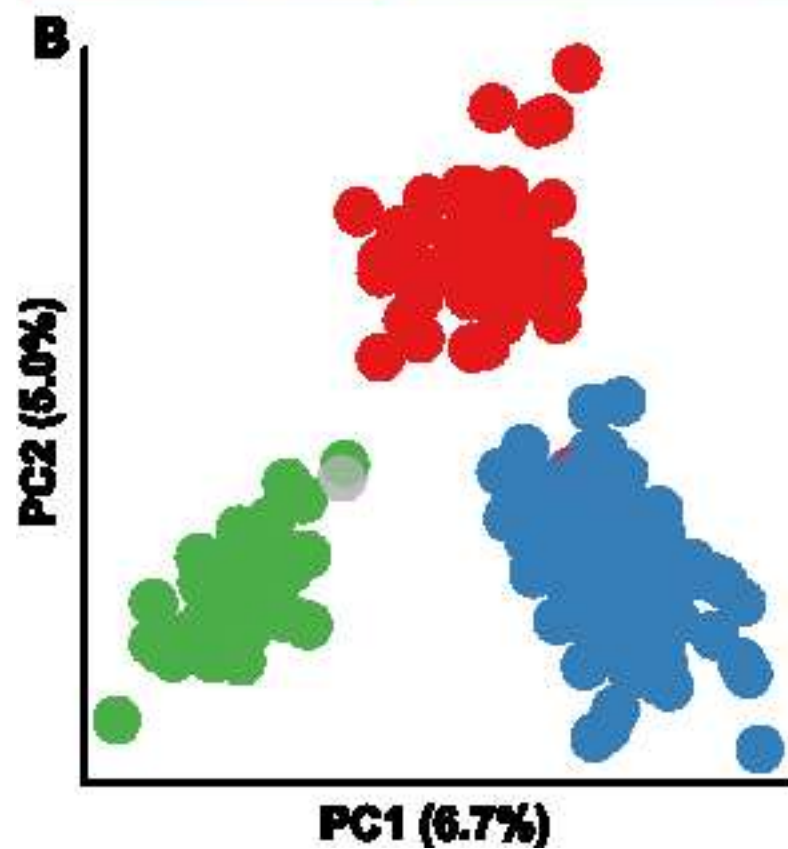
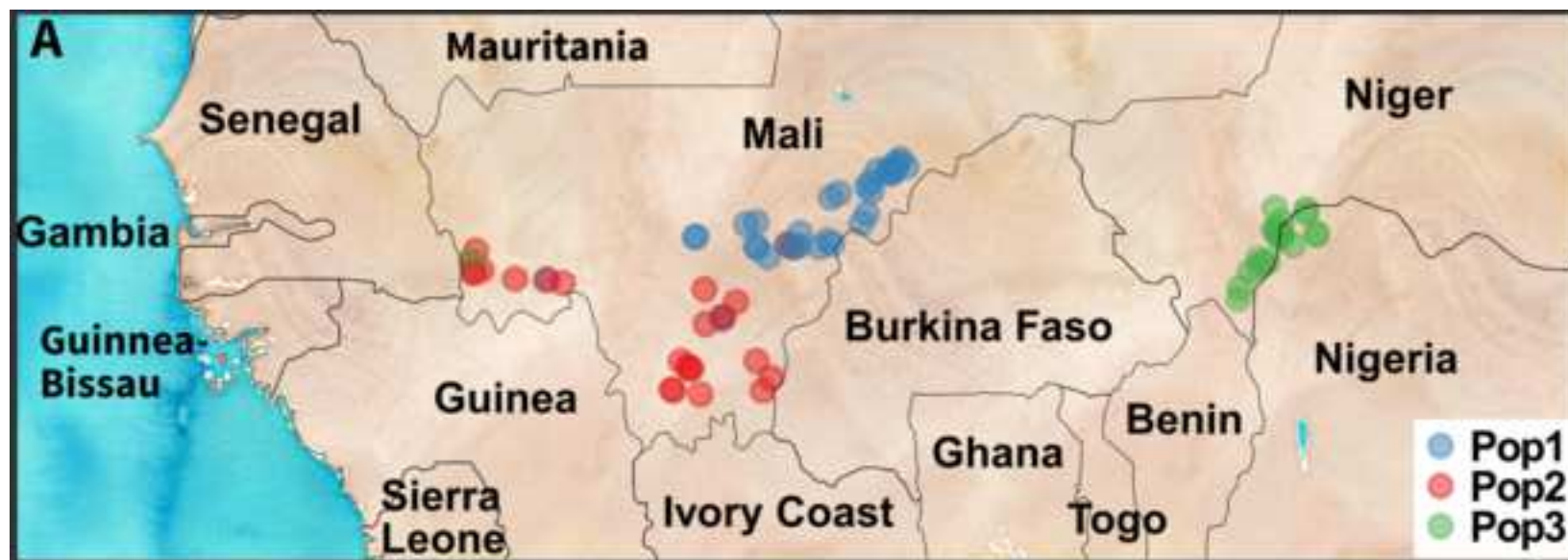
C



D











[Click here to access/download](#)

**Supplementary Material**

Supplementary Methods 072020.docx





Click here to access/download  
**Supplementary Material**  
Suppl. Tables and Figures.docx





Click here to access/download  
**Supplementary Material**  
Fono Suppl Table S6 and S7.xlsx



## UNIVERSITY OF CALIFORNIA, DAVIS

BERKELEY • DAVIS • IRVINE • LOS ANGELES • MERCED • RIVERSIDE • SAN DIEGO • SAN FRANCISCO

SANTA BARBARA • SANTA CRUZ

COLLEGE OF AGRICULTURAL AND ENVIRONMENTAL SCIENCES

SEED BIOTECHNOLOGY CENTER  
PLANT REPRODUCTIVE BIOLOGY  
EXTENSION CENTER DRIVE  
DAVIS, CA 95616PHONE: (530) 304-9329  
E-MAIL: [avandeynze@ucdavis.edu](mailto:avandeynze@ucdavis.edu)  
<http://sbc.ucdavis.edu>Allen Van Deynze  
Director of Seed Biotechnology Center  
Associate Director of Plant Breeding Center

July 10, 2020

Editors  
*GigaScience*

Dear Editors,

This letter accompanies our submission of the manuscript “Genome sequence and genetic diversity analysis of an under-domesticated orphan crop, white fonio (*Digitaria exilis*)” by Bennetzen and colleagues. This manuscript provides a very high-quality genome sequence for a vital African cereal, fonio (*Digitaria exilis*), based primarily on long read sequence assembly. We also provide information regarding the genetic diversity and population structure of fonio in its main regions of production, in West Africa. We discover that the fonio genome is the product of a fairly recent tetraploidy, but that internal heterozygosity is so low that outcrossing seems to be a very rare event. Because fonio needs further domestication and improvement, we identify specific genes that should be targets for mutational breeding/editing within the fonio genome. The data, analysis and discussion in this manuscript will greatly empower fonio researchers worldwide, leading to a newly accelerated potential for improvement of this orphan crop.

Appropriate reviewers for this manuscript would include Robert VanBuren of Michigan State Univ. ([bob.vanburen@gmail.com](mailto:bob.vanburen@gmail.com)), or James Schnable at the University of Nebraska ([schnable@unl.edu](mailto:schnable@unl.edu)). We prefer that the manuscript not be reviewed by Dr. Simon Krattinger or his colleagues at KAUST, because they are (friendly) competitors in the genomic study of fonio. The data, analyses and discussions in this manuscript have not been published or submitted for publication in any other journal.

We have addressed questions about data accessibility and the Nagoya Protocol in the Diversity analysis section and Data Availability Sections. We have updated the supplemental methods with analysis scripts and parameters and added 2 supplementary tables (S6 and S7) supporting data for diversity analysis sections.

Thank you for your time and efforts in the review of this manuscript. If I can be of further assistance, please feel free to contact me at your convenience.

Sincerely,

A. K. J.