

1 **Exceptional subgenome stability and functional divergence in allotetraploid teff, the**
2 **primary cereal crop in Ethiopia**

3
4 Robert VanBuren^{1,2*}, Ching Man Wai^{1,2}, Jeremy Pardo^{1,2,3}, Alan E. Yocca³, Xuewen Wang⁴, Hao
5 Wang⁴, Srinivasa R. Chaluvadi⁴, Doug Bryant⁵, Patrick P. Edger¹, Jeffrey L. Bennetzen⁴, Todd
6 C. Mockler⁵, Todd P. Michael^{6*}

7
8 ¹Department of Horticulture, Michigan State University, East Lansing, MI 48824, USA

9 ²Plant Resilience Institute, Michigan State University, East Lansing, MI 48824, USA

10 ³Department of Plant Biology, Michigan State University, East Lansing, MI 48824, USA

11 ⁴Department of Genetics, University of Georgia, Athens, GA 30602, USA

12 ⁵Donald Danforth Plant Science Center, St. Louis, MO 63132, USA

13 ⁶J. Craig Venter Institute, La Jolla, CA, 92037, USA

14 *Corresponding authors: bobvanburen@gmail.com, tmichael@jcvl.org

15 **Abstract**

16 Teff (*Eragrostis tef*) is a cornerstone of food security in the Horn of Africa, where it is prized for
17 stress resilience, grain nutrition, and market value. Despite its overall importance to small-scale
18 farmers and communities in Africa, teff suffers from low production compared to other cereals
19 because of limited intensive selection and molecular breeding. Here we report a chromosome-
20 scale genome assembly of allotetraploid teff (variety ‘Dabbi’) and patterns of subgenome
21 dynamics. The teff genome contains two complete sets of homoeologous chromosomes, with
22 most genes maintained as syntenic gene pairs. Through analyzing the history of transposable
23 element activity, we estimate the teff polyploidy event occurred ~1.1 million years ago (mya)
24 and the two subgenomes diverged ~5.0 mya. Despite this divergence, we detected no large-scale
25 structural rearrangements, homoeologous exchanges, or bias gene loss, contrasting most other
26 allopolyploid plant systems. The exceptional subgenome stability observed in teff may enable
27 the ubiquitous and recurrent polyploidy within Chloridoideae, possibly contributing to the
28 increased resilience and diversification of these grasses. The two teff subgenomes have
29 partitioned their ancestral functions based on divergent expression patterns among
30 homoeologous gene pairs across a diverse expression atlas. The most striking differences in
31 homoeolog expression bias are observed during seed development and under abiotic stress, and
32 thus may be related to agronomic traits. Together these genomic resources will be useful for
33 accelerating breeding efforts of this underutilized grain crop and for acquiring fundamental
34 insights into polyploid genome evolution.

35 Introduction

36 Thirty crop species supply over 90% of the world's food needs and this narrow diversity reduces
37 global food security. Humans have domesticated several hundred distinct plant species, but most
38 are underutilized, under-improved, and restricted to their regions of origin ¹. Although food
39 systems have become increasingly diverse in the last few decades, many locally adapted species
40 have been replaced by calorically dense staple crops, resulting in global homogeneity ². Many
41 underutilized and “orphan” crop species have desirable nutritional profiles, abiotic and biotic
42 stress resilience, and untapped genetic potential for feeding the growing population under the
43 changing climate.

44 Teff is the staple grain crop in Ethiopia, and it is preferred over other cereals because of
45 its nutritional profile, low input demand, adaptability, and cultural significance. Unlike other
46 major cereals, teff is grown primarily by small-scale, subsistence farmers. An estimated 130,000
47 locally adapted cultivars have been developed. Teff is among the most resilient cereals,
48 tolerating marginal and semi-arid soils that are unsuitable for wheat, maize, sorghum, and rice
49 production. Teff was likely domesticated in the northern Ethiopian Highlands where much of the
50 genetic diversity can be found ³⁻⁵. Consistent yields of small, nutritious seeds were the primary
51 domestication targets of teff, contrasting most cereals where large seed heads and high
52 productivity under tillage were desirable ⁵. Despite its stress tolerance, yield improvements lag
53 behind other cereals because of issues related to lodging, seed shattering, extreme drought, and
54 poor agronomic practices ⁶. Teff and other orphan cereals have undergone limited intensive
55 selection for high productivity under ideal conditions, and rapid gains should be possible with
56 advanced breeding and genome selection. A draft genome is available for the teff cultivar
57 ‘Tsedey’ (DZ-Cr-37) ⁷, but the utility of this reference is limited given its fragmented and
58 incomplete nature.

59 The wild progenitor of teff is likely *Eragrostis pilosa*; a hardy wild grass sharing
60 considerable overlap in morphological, genetic, and karyotype traits with teff ^{8,9}. *E. tef* and *E.*
61 *pilosa* are allotetraploids that arose from a shared polyploidy event of merging two distant,
62 unknown diploid genomes ⁹. Many crop plants are polyploid, and genome doubling can give rise
63 to emergent traits such as spinnable fibers in cotton ¹⁰, morphological diversity in *Brassica* sp. ¹¹,
64 and new aromatic profiles of strawberry fruits ¹². Successful establishment of allopolyploids
65 requires coordination of two distinct sets of homoeologous genes and networks, and often a
66 ‘dominant’ subgenome emerges to resolve genetic and epigenetic conflicts ¹³⁻¹⁵. The effect of
67 polyploidy on desirable traits and interactions between the two subgenomes remains untested in
68 teff. Polyploidy is found in more than 90% of species within the grass subfamily containing teff
69 (Chloridoideae), and this has been hypothesized to contribute to the stress tolerance and
70 diversification of these grasses ¹⁶. Here, we report a chromosome-scale assembly of the teff A
71 and B subgenomes and test for patterns of subgenome interactions and divergence.

72

73 Results

74 *Genome assembly and annotation*

75 We built a chromosome-scale assembly of the allotetraploid teff genome using a combination of
76 long read SMRT sequencing and long-range high-throughput chromatin capture (Hi-C). In total,
77 we generated 5.5 million filtered PacBio reads collectively spanning 52.9 Gb or 85x coverage of
78 the estimated 622 Mb ‘Dabbi’ teff genome. PacBio reads were error corrected and assembled
79 using Canu¹⁷ and the resulting contigs were polished to remove residual errors with Pilon¹⁸ using
80 high coverage Illumina data (45x). The PacBio assembly has a contig N50 of 1.55 Mb across
81 1,344 contigs with a total assembly size of 576 Mb; 92.6% of the estimated genome size. The
82 graph-based structure of the assembly has few bubbles corresponding to heterozygous regions
83 between haplotypes but contains numerous ambiguities related to high copy number long
84 terminal repeat (LTR) retrotransposons (Supplemental Figure 1). This pattern was also observed
85 in the genome assembly graph of the closely related grass, *Oropetium thomaeum*¹⁹. The average
86 nucleotide identity between homoeologous regions in teff is 93.9% in protein coding regions.
87 Thus, high sequence divergence facilitated accurate phasing and assembly. We utilized twenty
88 random fosmids to assess the accuracy of the PacBio-based assembly (Supplemental Table 1).
89 The fosmids collectively span 351kb and have an average identity of 99.9% to the teff genome
90 with individual fosmids ranging from 99.3 to 100%. This suggests that our assembly is mostly
91 complete and accurately polished.

92 Contigs from the Canu based draft genome were anchored into a chromosome-scale
93 assembly using a Hi-C based scaffolding approach. Illumina reads from the Hi-C library were
94 aligned to the PacBio contigs with BWA²⁰ followed by proximity based clustering using the
95 Juicer pipeline²¹. 150bp paired-end reads and aggressive filtering of non-uniquely mapped reads
96 were used to minimize chimeric mapping errors between homoeologous regions. After filtering,
97 twenty high-confidence clusters were identified, consistent with the haploid chromosome
98 number of teff (2n=40; Figure 1). In total, 687 contigs collectively spanning 96% of the
99 assembly (555 Mb) were anchored and oriented across the 20 pseudomolecules (Table 1).
100 Pseudomolecules ranged in size from 19 to 40 Mb, consistent with the teff karyotype²². Seven
101 chimeric contigs corresponding to joined telomeres were identified and split based on Hi-C
102 interactions. As described in the accompanying manuscript (see Wang et al. 2019), this genome
103 assembly was compared to a detailed genetic map of teff to revise and confirm chromosome-
104 scale assemblies for all 20 teff chromosomes, thus providing the opportunity to discover the A
105 and B genomes from the diploid progenitors of this allotetraploid (see below).

106 The teff genome was annotated using the MAKER pipeline. Transcript support from a
107 large-scale expression atlas and protein homology to Arabidopsis and other grass genomes were
108 used as evidence for *ab initio* gene prediction. After filtering transposon-derived sequences, *ab*
109 *initio* gene prediction identified 68,255 gene models. We assessed the annotation quality using
110 the Benchmarking Universal Single-Copy Ortholog (BUSCO) Embryophyta dataset. The
111 annotation contains 98.1% of the 1,440 core Embryophyta genes and the majority (1,210) are
112 found in duplicate in the A and B subgenomes.

113 The teff cultivar ‘Tsedey’ (DZ-Cr-37) was previously sequenced using an Illumina based
114 approach, yielding a highly fragmented draft genome with 14,057 scaffolds and 50,006 gene
115 models⁷. The fragmented nature of this assembly and incomplete annotation hinders

116 downstream functional genomics, genetics, and marker-assisted breeding of teff. We compared
117 the ‘Tsedey’ assembly with our ‘Dabbi’ reference to identify cultivar-specific genes and
118 differences in assembly quality. Only 30,424 (60.8%) of the ‘Tsedey’ gene models had
119 homology (>95% sequence identity) to gene models in our ‘Dabbi’ reference, including 9,866
120 homoeologous gene pairs. Only 20,208 (29.6%) of our ‘Dabbi’ gene models had homology to
121 ‘Tsedey’ gene models. The remaining gene models were unannotated or unassembled in the
122 ‘Tsedey’ assembly. Only one-third of the ‘Tsedey’ genome is assembled into scaffolds large
123 enough to be classified as syntenic blocks to ‘Dabbi’, which is an unavoidable artifact of the poor
124 assembly quality and low contiguity. Because of the fragmented nature of the ‘Tsedey’
125 assembly, we were unable to identify lineage-specific genes. Hence, the genomic resources
126 presented here represent a significant advance over previous efforts.

127

128 *Origins and subgenome dynamics*

129 Teff is an allotetraploid with unknown diploid progenitors, but the polyploidy event is
130 likely shared with other closely related *Eragrostis* species⁹. Because the diploid progenitors are
131 unknown and possibly extinct, we utilized the centromeric array sequences to distinguish the
132 homoeologous chromosomes from the A and B subgenomes of teff. Centromeric (CenT) repeat
133 arrays in teff range from 3.7 kb to 326 kb in size for each chromosome and individual arrays
134 contain 22 to 824 copies (Supplemental Table 2). We identified two distinct CenT arrays in teff
135 (hereon referend to as CenTA and CenTB). CenTA and CenTB are the same length (159 bp) but
136 have different sequence composition (Supplemental Figure 2b). Alignment of the consensus
137 CenT arrays identified several distinguishing polymorphisms and a maximum likelihood
138 phylogenetic tree separated the CenT arrays into two well-supported clades (Supplemental
139 Figure 2a). Each clade contains one member from each of the ten homoeologous chromosome
140 pairs and this classification likely represents differences in Cen array composition between the
141 diploid progenitor species. This approach allowed us to accurately distinguish homoeologous
142 chromosome pairs from the A and B subgenomes and verifies the allopolyploid origin of teff.

143 The Teff subgenomes have 93.9% sequence homology in the coding regions, suggesting
144 that either the polyploidy event was relatively ancient or that the progenitor diploid species were
145 highly divergent²³. To estimate the divergence time of the A and B subgenomes, we calculated
146 Ks (synonymous substitutions per synonymous site) between homoeologous gene pairs. Teff
147 homoeologs have a single Ks peak with a median of 0.15 (Supplemental Figure 3),
148 corresponding to a divergence time of ~5 million years based on a widely used mutation rate for
149 grasses²⁴. The ten pairs of homoeologous chromosomes are highly syntenic with no large-scale
150 structural rearrangements. The A subgenome is 13% (37 Mb) larger in size but contains only 5%
151 more genes than the B subgenome (34,032 vs. 32,255; Table 1). Most genes (54,846) are
152 maintained as homoeologous pairs and 13,409 are found in only one subgenome. We identified
153 6,876 tandemly duplicated genes with array sizes ranging from 2 to 15 copies. Of the 2,748
154 tandem arrays, 998 are found in both subgenomes, while 864 and 1,008 occur in only the A and
155 B subgenomes, respectively (Table 1). Copy number varies extensively in shared arrays between
156 the subgenomes.

157 The monoploid genome size of teff is relatively small (~300 Mb) compared to other
158 polyploid grasses, and repetitive elements constitute a low percentage (25.6%) of the genome.
159 Long terminal repeat retrotransposons (LTR-RTs) are the most abundant repetitive elements,
160 spanning at least 115.9 Mb or ~20.0% of the genome (Supplemental Table 3). This predicted
161 percentage is somewhat lower than that reported for other small grass genomes such as
162 *Oropetium* (250 Mb; 27%)^{19,25} and *Brachypodium* (272 Mb; 21.4%)²⁶. We classified LTRs into
163 families and compared their abundance and insertion times (Figure 2). A particular window of
164 activity was seen for six families of LTR-RTs that were active only in the A genome progenitor
165 or the B genome progenitor (Supplemental Figure 4, Supplemental Table 4). The insertion times
166 for these genome-specific LTR-RTs were all greater than 1.1 mya, indicating the two
167 subgenomes were evolving independently during this period. Hence, this LTR-RT analysis both
168 confirms the A and B genome designations, and provides a novel methodology for determining
169 the date of polyploid formation. In teff, these data indicate that the ancestral polyploidy was
170 established ~1.1 mya.

171 Five of the six subgenome-specific LTR-RT families were found only in the A
172 subgenome, suggesting that LTR-RTs accumulate more rapidly in the A subgenome or are
173 purged more effectively in the B subgenome. This recent bursts of LTR-RT activity contributes
174 to the 13% larger size of the A subgenome. There are 24 families with median insertion times
175 between 1.1 and 2.4 MYA, and the remaining 18 families do not exhibit subgenomic specificity.
176 Of these, 15 show no apparent burst in amplification, and three evidence of very recent (post-
177 polyploid) activity (Figure 4, Supplementary Figure 4, Supplemental Table 5).

178 Teff belongs to the Chloridoideae subfamily of grasses²⁷ which includes important
179 drought and heat tolerant C4 species such as the orphan grain crop finger millet and model
180 desiccation tolerant plants in the genera *Oropetium*, *Eragrostis*, *Tripogon*, *Sporobolus*, and
181 others. Most (~90%) of surveyed Chloridoideae species are polyploid, including many of the
182 aforementioned taxa, and this likely contributes to their diversity and stress tolerance¹⁶. We
183 utilized the wealth of genomic resources within Chloridoideae and more generally across
184 Poaceae to identify patterns associated with improved stress tolerance, polyploidy and genome
185 evolution in teff. The teff and *Oropetium* genomes have near complete collinearity, as
186 demonstrated by highly conserved gene content and order along each chromosome (Figure 3).
187 Teff and *Oropetium* show a clear 1:2 synteny pattern with 87% of teff genes having synteny to
188 one block in *Oropetium* and 85% of *Oropetium* genes having synteny to two blocks in the teff
189 genome (Figure 3a). This ratio corresponds to the A and B homoeologs of tetraploid teff and the
190 single orthologs of diploid *Oropetium*. Each *Oropetium* chromosome has clear collinearity to
191 two homoeologous teff chromosomes (Figure 3c). Three trios have no rearrangements (teff 3A,
192 3B, and *Oropetium* Chr3; 4A, 4B, Chr4; 6A, 6B, Chr8) six trios have one or more large-scale
193 inversions (1A, 1B, Chr1; 2A, 2B, Chr2; 5A, 5B, Chr7; 7A, 7B, Chr6; 8A, 8B, Chr9; 9A, 9B,
194 Chr5) and one trio has translocations (10A, 10B, Chr10). Of the 28,909 *Oropetium* genes, 74%
195 (21,293) have syntenic orthologs in both subgenomes of teff, 5% (1,503) are found in only one
196 subgenome, and 21% (6,113) have no syntenic orthologs in teff. Teff and the allotetraploid grain
197 crop finger millet have 2:2 synteny but only 69% of syntenic blocks are found in duplicate
198 because of the fragmented nature of the finger millet genome assembly²⁸ (Supplemental Figure

199 5). Only 56% (38,149) of the teff genes have two syntenic orthologs in finger millet and the
200 remaining 13 and 30% (9,228 and 20,878) have one or zero syntenic orthologs in finger millet
201 respectively.

202 Using Oropetium and teff syntenic orthologs, we calculated the ratio of nonsynonymous
203 (Ka) to synonymous substitutions (Ks) to identify genes putatively under selection during
204 domestication in teff. The top 10% of genes with the highest Ka/Ks ratios in teff (cutoff of 0.38)
205 are enriched in gene ontology (GO) terms related to somatic embryogenesis, pollen
206 differentiation, and reproductive phase transition among others (Supplemental Table 6). These
207 genes may have been intentional or inadvertent targets during domestication.

208 Following an allopolyploidy event, a dominant subgenome often emerges with
209 significantly more retained genes and higher homoeolog expression as the plant returns to a
210 diploid-like state¹³. This dominance is established immediately following the polyploidy event
211¹⁵, and patterns of biased fractionation have been observed in *Arabidopsis*¹³, maize²⁹, *Brassica*
212 *rapa*³⁰, and bread wheat³¹. Biased homoeolog loss (fractionation) is not universal, and other
213 allopolyploids such as *Capsella bursa-pastoris*³² and several *Cucurbita* species³³ display no
214 subgenome dominance. We searched for biased fractionation using syntenic orthologs from
215 Oropetium as anchors. The A and B subgenomes of teff have a near identical number of syntenic
216 orthologs to Oropetium (19,277 vs. 19,292 respectively) suggesting that there is little or no
217 biased fractionation in teff. Orthologs to 1,308 Oropetium genes are found as single copy loci in
218 teff, including 647 and 678 from the A and B subgenomes respectively. The remaining
219 orthologs are maintained in duplicate in teff compared to their single ortholog in Oropetium.
220 Together this suggests a general stability of gene content in *Eragrostis* after genome merger.

221

222 *Homoeolog expression patterns and subgenome dominance*

223 To test for patterns of sub-genome differentiation and dominance in teff, we surveyed
224 gene expression in eight developmentally distinct tissue types and two stages of progressive
225 drought stress. Sampled tissues include roots and shoots from seedlings and mature plants,
226 internodes, and two stages of developing seeds. Tissue from mature, well-watered leaves and two
227 time points of severe drought were also collected (leaf relative water content of 33% and 16%
228 respectively). Of the 23,303 syntenic gene pairs between the A and B subgenome, 15,325 have
229 homoeologous expression bias (HEB) in at least one tissue and 1,694 have biased expression in
230 all sampled tissues (Supplemental Figure 6). Pairwise comparisons between syntenic gene pairs
231 support a slight bias in transcript expression toward the B subgenome (Figure 4a). Roughly 56%
232 of the 207,873 pairwise comparisons across the ten tissues show biased expression toward
233 homoeologs in the B subgenome (Wilcoxon rank sum $P < 0.001$). This pattern is consistently
234 observed across all ten tissues and most chromosome pairs, but the difference is subtle when
235 robust cutoffs of differential expression are applied (Figure 4b and c; see methods). Individual
236 tissues have from 6,061 to 8,485 homoeologous gene pairs with significant differential
237 expression, including 52.3% biased toward the B subgenome (Kruskal–Wallis H test $P < 0.01$;
238 Figure 4b). Eight pairs of chromosomes show HEB toward the B subgenome, and chromosomes

239 1 and 8 have more dominant homoeologs from the A subgenome, but the difference is not
240 significant (Wilcoxon rank sum $P > 0.05$). Together this suggests that the B subgenome is
241 universally dominant over the A subgenome but when strict thresholds are applied, this
242 difference is minimal. Although we detected no evidence of recent homoeologous exchange, it is
243 possible that genes from the recessive genome were replaced with homoeologs from the
244 dominant subgenome, which would weaken patterns of subgenome dominance³⁴.

245 We tested whether gene pairs with HEB maintain patterns of dominance across all tissues
246 or whether dominant homoeologs are reversed in different tissues or under stress. The vast
247 majority of genes (86.9%; 13,322) with homoeologous expression bias maintain the same pattern
248 of dominance across all tissues, while 13.1% (2,002) of gene pairs have opposite dominance
249 patterns in different tissues. The remaining 7,675 gene pairs have no expression bias in any
250 tissues or both homoeologs have negligible expression. Severely dehydrated leaf tissue had the
251 most gene pairs with HEB (36%; 8,485) compared to seedling roots and shoots which each had
252 ~26% of pairs with HEB. These results are consistent with previous findings in allohexaploid
253 wheat³⁵ and allotetraploid *Tragopogon mirus*³⁶. We compared the ratio of nonsynonymous (Ka)
254 to synonymous substitution rates (Ks) in homoeologous gene pairs to test if genes with stronger
255 HEB are experiencing different patterns of selection. Gene pairs with stronger HEB had
256 significantly higher Ka/Ks than gene pairs with no HEB in any tissue (Supplemental Figure 7;
257 0.17 vs. 0.28; Mann-Whitney $P < 0.01$). We detected no difference in divergence (Ks) among
258 genes with varying degrees of HEB (Supplemental Figure 8). This suggests homoeologous gene
259 pairs with higher expression divergence are under more relaxed selective constraints than gene
260 pairs with balanced expression.

261

262 Discussion

263 Unlike the genomes of most polyploid grasses, the teff subgenomes are relatively small (~300
264 Mb), with high gene density and low transposable element content. The subgenomes are highly
265 syntenic along their length with no evidence of major inversions or structural rearrangements,
266 contrasting patterns observed in other similarly aged allopolyploids such as wheat³⁷, canola
267 (*Brassica napus*)³⁸, strawberry (*Fragaria ananassa*)³⁹, cotton⁴⁰, and proso millet⁴¹. The
268 general stability of the teff subgenomes may be attributed to low rates of homoeologous
269 exchange. An estimated 90% of Chloridoid grasses are polyploid and among the allopolyploid
270 species, multivalent pairing is rarely detected¹⁶. The twenty chromosome pairs in teff show
271 bivalent pairing in meiosis I²², and double reduction has not been observed in segregating
272 populations^{42,43}. Although homoeologous exchanges can result in advantageous emergent
273 phenotypes, they can also destabilize the karyotype, leading to reduced fertility and fitness⁴⁴. For
274 this reason, recent polyploids have long been considered “evolutionary dead ends”⁴⁵. Thus,
275 proper bivalent pairing (disomic inheritance) in natural allopolyploids may be favored, and the
276 near perfect synteny observed between teff subgenomes suggests that an underlying mechanism
277 may exist to prevent or reduce homoeologous exchanges in this species. We detected no
278 evidence of recent homoeologous exchange in teff based on Ks distribution, including exchanges
279 that would have happened at the inception of the polyploidy event 1.1 mya. Homeologous

280 exchanges are a common feature of allopolyploids ³⁴, and the lack of these events is a unique
281 feature of the teff genome.

282 The Teff A and B subgenomes, and *Oropetium genome* have high degrees of
283 chromosome level collinearity despite their distant divergence. This is particularly unusual as
284 polyploidy-rich lineages typically have high rates of chromosome evolution ⁴⁶. In contrast, our
285 analysis of the divergence dates of the diploid A and B genome ancestors (~5 mya) and the
286 formation of the tetraploid (~1.1 mya) indicates that the two genomes were so similar in structure
287 (i.e., gene content, gene order and chromosome size) that some tetrasomic pairing would have
288 been expected. Perhaps the status of the *PhI*-equivalent locus (loci) ⁴⁷ in *Eragrostis* is (are) so
289 dominant, that even low frequencies of homoeologous pairing are blocked. The high levels of
290 subgenome compatibility, genetic and chromosome stability, fidelity for chromosome pairing,
291 and low rates of homoeologous exchange allows polyploidy to dominate in the Chloridoideae
292 subfamily. This polyploidy in turn may have enabled the emergent resilience and robustness
293 observed in Chloridoid grasses.

294 Although we detected no biased fractionation between the teff subgenomes, we observed
295 a general subgenome dominance across tissues in the expression atlas. The B subgenome is
296 smaller and has fewer transposable elements, which may be contributing to the overall higher
297 homoeolog expression levels ¹⁵. Patterns of B subgenome dominance are relatively weak
298 compared to other allopolyploids ¹⁵, which may reflect the stability and lack of biased
299 fractionation in teff. The teff subgenomes have successfully partitioned their ancestral roles, and
300 most gene pairs display homoeolog expression bias. This bias is generally maintained across
301 tissues and treatments, and few gene pairs change bias in a tissue-specific manner. Severely
302 drought stressed leaf tissue has the highest proportion of genes with biased expression, which
303 may reflect adaptation to adverse environments. Extensive homoeolog expression bias is also
304 observed in hexaploid wheat ³⁵, octoploid sugarcane ⁴⁸, and tetraploid *Tragopogon mirus* ³⁶ and
305 may be a common feature of recent polyploid grasses.

306 The vast majority of genes in Teff are maintained as homoeologous gene pairs in the A
307 and B subgenomes, providing a significant obstacle for targeted breeding. Efforts to produce
308 semi-dwarf, lodging resistant teff using a mutagenesis approach have been more difficult
309 because of gene redundancy ⁴⁹. The resources provided here will help accelerate marker-assisted
310 selection and guide genome engineering-based approaches, which must take gene redundancy
311 into account. Most gene pairs have divergent expression profiles such that the subgenomes likely
312 contribute unequally to different agronomic traits. Teff is often described as an orphan grain crop
313 because of its limited investigation and improvement, resulting in relatively low yields under
314 ideal conditions compared to other cereals with intensive selection and breeding histories. Teff
315 and other grasses within Chloridoideae have high tolerance to abiotic stresses, and most of this
316 resilience was maintained during teff domestication. This may represent a historical alternative
317 selection scheme where maximum yield is exchanged for reliable harvest under poor
318 environmental conditions. Future efforts to improve food security should utilize the natural
319 resilience of these robust, stable, polyploid species.

320

321 **Methods**

322 *Plant materials*

323 The ‘Dabbi’ cultivar of teff (PI 524434, www.ars-grin.gov) was chosen for sequencing and for
324 constructing the expression atlas. Plant materials for High molecular weight (HMW) genomic
325 DNA extraction, Hi-C library construction and RNA were maintained in growth chambers under
326 a 12-hour photoperiod with day/night temperatures of 28°C and 22°C respectively and a light
327 intensity of 400 $\mu\text{E m}^{-2} \text{sec}^{-1}$. Tissue samples for the expression atlas were collected at ZT8
328 (Zeitgeber Time 8) to reduce issues associated with circadian oscillation. The tissue types used in
329 the expression atlas include shoots and roots from young seedlings, mature leaf, internode, root,
330 immature seeds and mature seeds. For the drought time points, mature teff plants were allowed
331 to dry slowly and leaf tissue was collected at subsequent days of extreme drought when the plant
332 tissues had 33% and 16% relative water content, as well as well-watered teff for comparison.
333 Three biological replicates were collected for each sample in the expression atlas. Leaf tissue
334 from seedlings was used for the HMW genomic DNA extraction and Hi-C library construction.
335 Tissues for HMW genomic DNA extraction and RNaseq were immediately frozen in liquid
336 nitrogen and stored at -80° C.

337

338 *DNA isolation, library construction, and sequencing*

339 HMW genomic DNA was isolated from young teff leaf tissue for both PacBio and Illumina
340 sequencing. A modified nuclei preparation⁵⁰ was used to extract HMW gDNA and residual
341 contaminants were removed using phenol chloroform purification. PacBio libraries were
342 constructed using the manufacturer’s protocol and were size selected for 30 kb fragments on the
343 BluePippen system (Sage Science) followed by subsequent purification using AMPure XP beads
344 (Beckman Coulter). The PacBio libraries were sequenced on a PacBio RSII system with P6C4
345 chemistry. In total, 5.5 million filtered PacBio reads were generated, collectively spanning 52.9
346 Gb or ~85x genome coverage (assuming a genome size of 622 Mb). The same batch of HMW
347 genomic DNA was used to construct Illumina DNaseq libraries for correcting residual errors in
348 the PacBio assembly. Libraries were constructed using the KAPA HyperPrep Kit (Kapa
349 Biosystems) followed by sequencing on an Illumina HiSeq4000 under paired-end mode (150
350 bp).

351

352 *RNA extraction and library construction*

353 RNA for the expression atlas was extracted using the Omega Biotek E.Z.N.A. ® Plant RNA kit
354 according to the manufacturer’s protocol. Roughly 200 mg of ground tissue was used for each
355 extraction. The RNA quality was validated using gel electrophoresis and the Qubit RNA IQ
356 Assay (ThermoFisher). Stranded RNaseq libraries were constructed using 2 μg of total RNA
357 quantified using the Qubit RNA HS assay kit (Invitrogen, USA) with the Illumina TruSeq
358 stranded total RNA LT sample prep kit (RS-122-2401 and RS-122-2402). Multiplexed libraries

359 were pooled and sequenced on an Illumina HiSeq4000 under paired-end 150nt mode. Three
360 replicates were sequenced for each timepoint/sample.

361

362 *Genome assembly*

363 The genome size of ‘Dabbi’ teff was estimated using flow cytometry as previously described⁵¹.
364 The estimated flow cytometry size was 622 Mb, which was consistent with kmer-based
365 estimations from Illumina data. The kmer plot had a unimodal distribution suggesting low within
366 genome heterozygosity and high differentiation from the teff A and B subgenomes. Raw PacBio
367 data was error corrected and assembled using Canu (V1.4)⁵² which produced accurate and
368 contiguous assembly for homozygous plant genomes. The following parameters were modified:
369 minReadLength=2000, GenomeSize=622Mb, minOverlapLength=1000. Assembly graphs were
370 visualized after each iteration of Canu in Bandage⁵³ to assess complexities related to repetitive
371 elements and homoeologous regions. The final Canu based PacBio assembly has a contig N50 of
372 1.55 Mb across 1,344 contigs with a total assembly size of 576 Mb. The raw PacBio contigs
373 were polished to remove residual errors with Pilon (V1.22)¹⁸ using 73x coverage of Illumina
374 paired-end 150 bp data. Illumina reads were quality-trimmed using Trimmomatic⁵⁴ followed by
375 aligning to the assembly with bowtie2 (V2.3.0)⁵⁵ under default parameters. Parameters for Pilon
376 were modified as follows: --flank 7, --K 49, and --mindepth 15. Pilon was run recursively three
377 times using the modified corrected assembly after each round. Ten full-length fosmids
378 (collectively spanning 351kb) were aligned to the final PacBio assembly to assess the quality.
379 The fosmids exhibited an average identity of 99.9% to the PacBio assembly, with individual
380 fosmids ranging from 99.3 to 100% nucleotide identity.

381

382 *Hi-C analysis and pseudomolecule construction*

383 The PacBio based teff contigs were anchored into a chromosome-scale assembly using a Hi-C
384 proximity-based assembly approach as previously described¹⁹. A Hi-C library was constructed
385 using 0.2 g of leaf tissue collected from newly emerged teff seedlings with the Proximo™ Hi-C
386 Plant kit (Phase Genomics) following the manufacturer’s protocol. After verifying quality, the
387 Hi-C library was size-selected for 300-600 bp fragments and sequenced on the Illumina HiSeq
388 4000 under paired-end 150 bp mode. 150 bp reads were used to avoid erroneous alignment in
389 highly similar homoeologous regions. In total, 226 million read pairs were used as input for the
390 Juicer and 3d-DNA Hi-C analysis and scaffolding pipelines^{21,56}. Illumina reads were quality-
391 trimmed using Trimmomatic⁵⁴ and aligned to the contigs using BWA (V0.7.16)²⁰ with strict
392 parameters (-n 0) to prevent mismatches and non-specific alignments in repetitive and
393 homoeologous regions. Contigs were ordered and oriented and assembly errors were identified
394 using the 3d-DNA pipeline with default parameters⁵⁶. The resulting hic contact matrix was
395 visualized using Juicebox, and misassemblies and misjoins were manually corrected based on
396 neighboring interactions. This approach identified 20 high-confidence clusters representing the
397 haploid chromosome number in Teff. The manually validated assembly was used to build
398 pseudomolecules using the finalize-output.sh script from 3d-DNA and chromosomes were

399 renamed and ordered by size and binned to the A and B subgenomes based on centromeric array
400 analysis (described in detail below).

401

402 *Identification of repetitive elements*

403 We first identified and masked the simple sequence repeats in the teff genome with GMATA⁵⁷,
404 and then conducted structure-based full-length transposable element (TE) identification using the
405 following bioinformatic tools: LTR_FINDER⁵⁸ and LTRharvest⁵⁹ to find LTR-RTs,
406 LTR_retriever⁶⁰ to acquire high-confidence full LTR retrotransposons, SINE-Finder⁶¹ to
407 identify SINEs, MGEscan-nonLTR (V2)⁶² to identify LINEs, MITE-Hunter⁶³ and MITE
408 Tracker⁶⁴ to identify TIRs, and HelitronScanner⁶⁵ to identify *Helitrons*. All TEs were classified
409 and manually checked according to the nomenclature system of transposons as described
410 previously⁶⁶ and against Repbase to validate their annotation⁶⁷. We used the newly identified
411 TEs as a custom library to identify full length and truncated TE elements through a homology-
412 based search with RepeatMasker (<http://www.repeatmasker.org>, version 4.0.7)⁶⁸ using the teff
413 pseudomolecules as input. Parameters for RepeatMasker were as described previously⁶⁹, and all
414 other parameters were left as default. The distribution of repeat sequences was then calculated.
415 Only LTR-RT families with at least 5 intact copies were used for analysis of subgenome
416 specificity. Within the 65 families having > 5 intact elements, we identified LTRs with
417 subgenomic specific activity. A family is considered as subgenomic specific if all intact elements
418 of this family are from the same subgenome. Subgenome specificity was verified through
419 BLAST of the element against the genome, and the distribution of matched sequences was
420 manually inspected for subgenome specificity. The approximate insertion dates of LTR-RTs
421 were calculated using the evolutionary distance between two LTR-RTs^{70,71} with the formula of
422 $T=K/2\mu$, where K is the divergence rate approximated by percent identity and μ is the neutral
423 mutation rate estimated as $\mu=1.3 \times 10^{-8}$ mutations per bp per year.

424 Centromeric repeat arrays were identified with the approach outlined in⁷² using Tandem
425 repeat finder (Version 4.07)⁷³. Parameters were modified as follows for Tandem repeat finder:
426 '1 1 2 80 5 200 2000 -d -h'. Centromere-specific repeats are often the most abundant tandem
427 repeats in the genome, and they were identified in teff by the following criteria: (1) copy number,
428 (2) sequence level conservation between chromosomes, (3) similarity to other grass repeats, and
429 (4) proximity to centromere-specific *gypsy* LTR-RTs. This approach identified two distinct
430 centromere-specific arrays (CenTA and CenTB) with a shared length of 159 bp yet distinct
431 sequence compositions. The consensus sequence of centromeric repeats from each chromosome
432 was used to construct a maximum likelihood phylogenetic tree implemented in MEGA5
433 (V10.0.5)⁷⁴. This approach separated centromeric repeats from the twenty chromosomes into
434 two distinct groups corresponding to the A and B subgenomes.

435

436 *Genome annotation*

437 Genes in the teff genome were annotated using the MAKER-P pipeline⁷⁵. The LTR-RT repeat
438 library from LTR retriever was used for repeat masking. Transcript-based evidence was
439 generated using RNAseq data from the ten tissues of the teff expression atlas. Quality trimmed
440 RNAseq reads were aligned to the unmasked teff genome using the splice aware alignment
441 program STAR (v2.6)⁷⁶ and transcripts were identified using StringTie (v1.3.4)⁷⁷ with default
442 parameters. The `-merge` flag was used to combine the output from individual libraries to
443 generate a representative set of non-redundant transcripts. Protein sequences from the
444 Arabidopsis⁷⁸, rice⁷⁹, and sorghum⁸⁰ genomes as well as proteins from the UniProtKB plant
445 databases⁸¹ were used as protein evidence. *Ab initio* gene prediction was conducted using SNAP⁸²
446 and Augustus (3.0.2)⁸³ with two rounds of iterative training. The resulting gene models were
447 filtered to remove any residual repetitive elements using BLAST with a non-redundant
448 transposase library. The annotation quality was assessed using the benchmarking universal
449 single-copy orthologs (BUSCO; v.2)⁸⁴ with the plant-specific dataset (embryophyta_odb9).

450

451 *RNAseq expression analysis and homoeolog expression bias*

452 Gene expression levels were quantified with the pseudo-aligner Kallisto (v 0.44.0)⁸⁵ using the
453 teff gene models as a reference. Paired-end Illumina reads from the ten tissues in the expression
454 atlas were quality trimmed using Trimmomatic (V0.33) with default parameters and pseudo-
455 aligned to the gene models with Kallisto under default parameters with 100 bootstraps per
456 sample. The teff A and B subgenomes have high sequence divergence (~7%) such that
457 misalignment between homoeologs was minimal. Expression levels were quantified as
458 Transcripts Per Million (TPM) and the three biological replicates were averaged for direct
459 homoeolog comparisons.

460

461 *Comparative genomics*

462 Homoeologous gene pairs between the teff A and B subgenomes and syntenic gene pairs across
463 select grasses were identified using the MCSCAN toolkit implemented in python
464 ([https://github.com/tanghaibao/jcvi/wiki/MCscan-\(Python-version\)](https://github.com/tanghaibao/jcvi/wiki/MCscan-(Python-version))). Teff homoeologs were
465 identified by all vs. all alignment using LAST, and hits were filtered using default parameters in
466 MCSCAN with a minimum block size of 5 genes. This approach identified 23,303
467 homoeologous, syntenic gene pairs between the A and B subgenome. Homoeologs gene pairs
468 with translocations were not identified using this syntenic approach and were thus excluded from
469 analysis. Tandem gene duplicates in teff were identified from the all vs. all LAST output with a
470 maximum gene distance of 10. Gene models from teff were aligned to the *Oropetium thomaeum*
471^{19,25} and *Sorghum bicolor*⁸⁰ genes as outlined above for comparative genomics analyses across
472 grasses. Macro and microsytentic dot plots, block depths, and karyotype comparisons were
473 generated in python using scripts from MSCAN.

474 Ka and Ks values were computed using a set of custom scripts available on GitHub:
475 https://github.com/Aeyocca/ka_ks_pipe/. The homoeologous gene pair list from the teff

476 subgenomes and syntenic orthologs between teff and *Oropetium* were used as input and the
477 protein sequences from each gene pair were aligned using MUSCLE v3.8.31⁸⁶. PAL2NAL
478 (v14)⁸⁷ was used to convert the peptide alignment to a nucleotide alignment and Ks values were
479 computed between gene pairs using codeml from PAML (V4.9h)⁸⁸ with parameters specified in
480 the control file found in the GitHub repository listed above.

481

482 **Data availability**

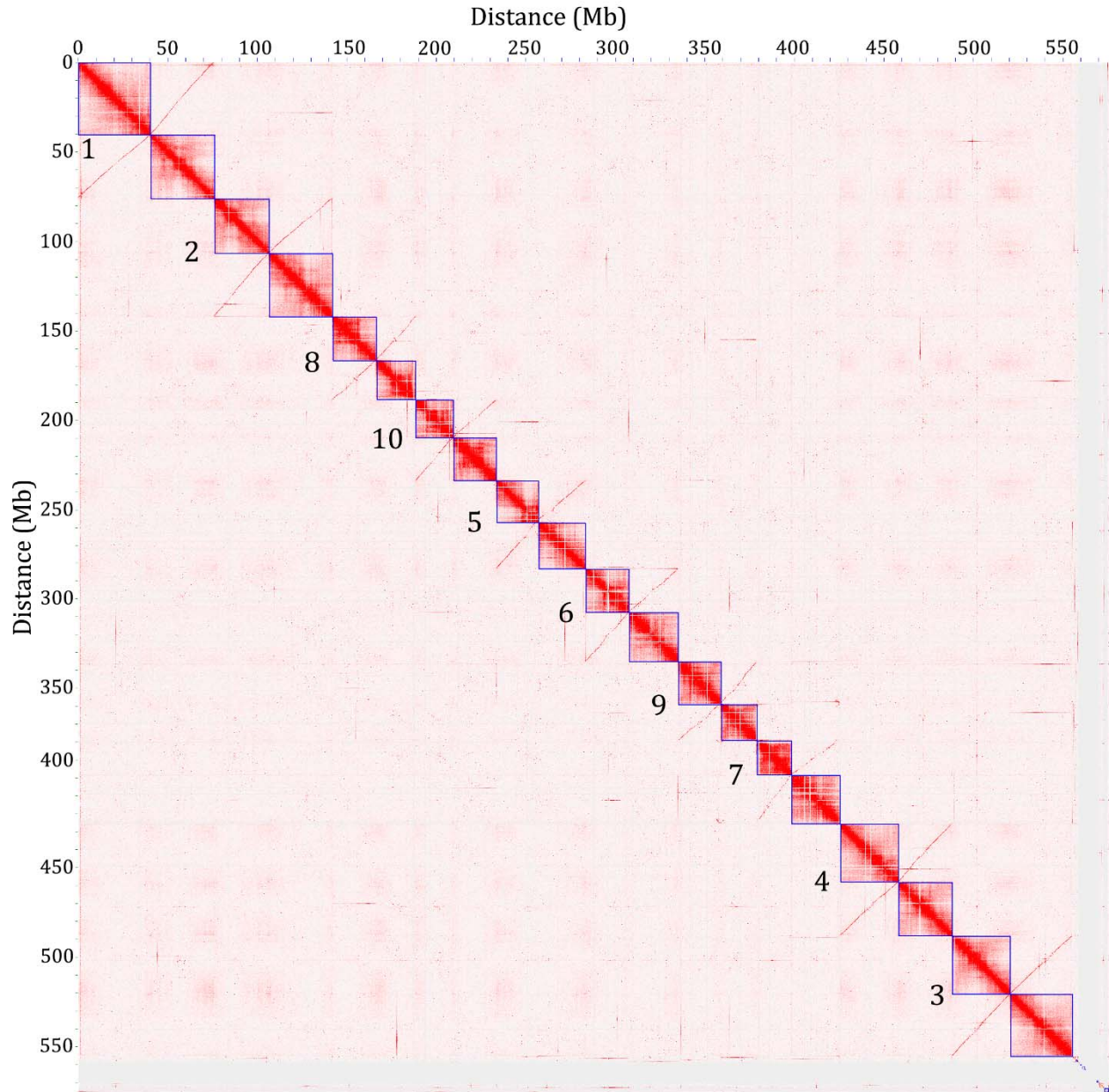
483 The raw PacBio data, Illumina DNaseq, and RNAseq data are available from the National
484 Center for Biotechnology Information Short Read Archive. RNAseq reads from the teff
485 expression atlas were deposited to the National Center for Biotechnology Information Short
486 Read Archive under bioproject PRJNA525065. The genome assembly and annotation for Tef is
487 available from CoGe under genome ID: id50954.

488

489 **Acknowledgments**

490 We are indebted to Tsegaye Dabi at the Salk Institute for Biological Studies for introducing us to
491 this amazing plant, and for inspiring generations of plant biologists. We thank Elliott Meer for
492 assistance with PacBio sequencing, and the Monsanto Genomics Team (Randy Kerstetter, Mitch
493 Sudkamp, Phil Latreille, Zijin Du and Joe Zhou) for full length sequenced fosmids. We thank
494 James Schnable for his helpful comments and suggestions on the manuscript. This work is
495 supported by funding from the National Science Foundation (MCB-1817347 to R.V.),
496 Department of Energy (DE-SC0012639 to T.C.M. and T.P.M.), and partial support from the Bill
497 & Melinda Gates Foundation (T.C.M. and D.B.).

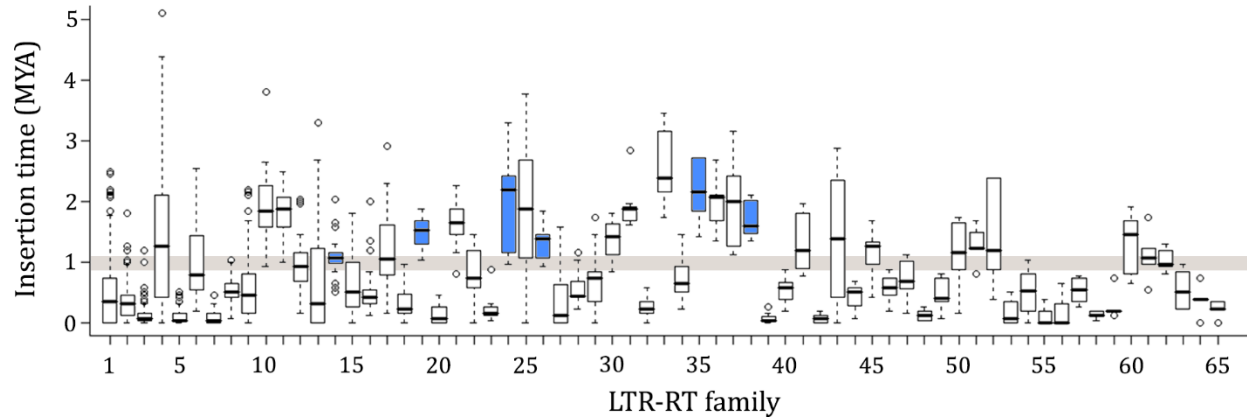
498



499

500 **Figure 1. Hi-C based clustering of the teff genome.** Heat map showing the density of Hi-C
501 interactions between contigs with red indicating high density of interactions. Distinct
502 chromosomes are highlighted by blue boxes and homoeologous chromosome pairs are
503 numbered.

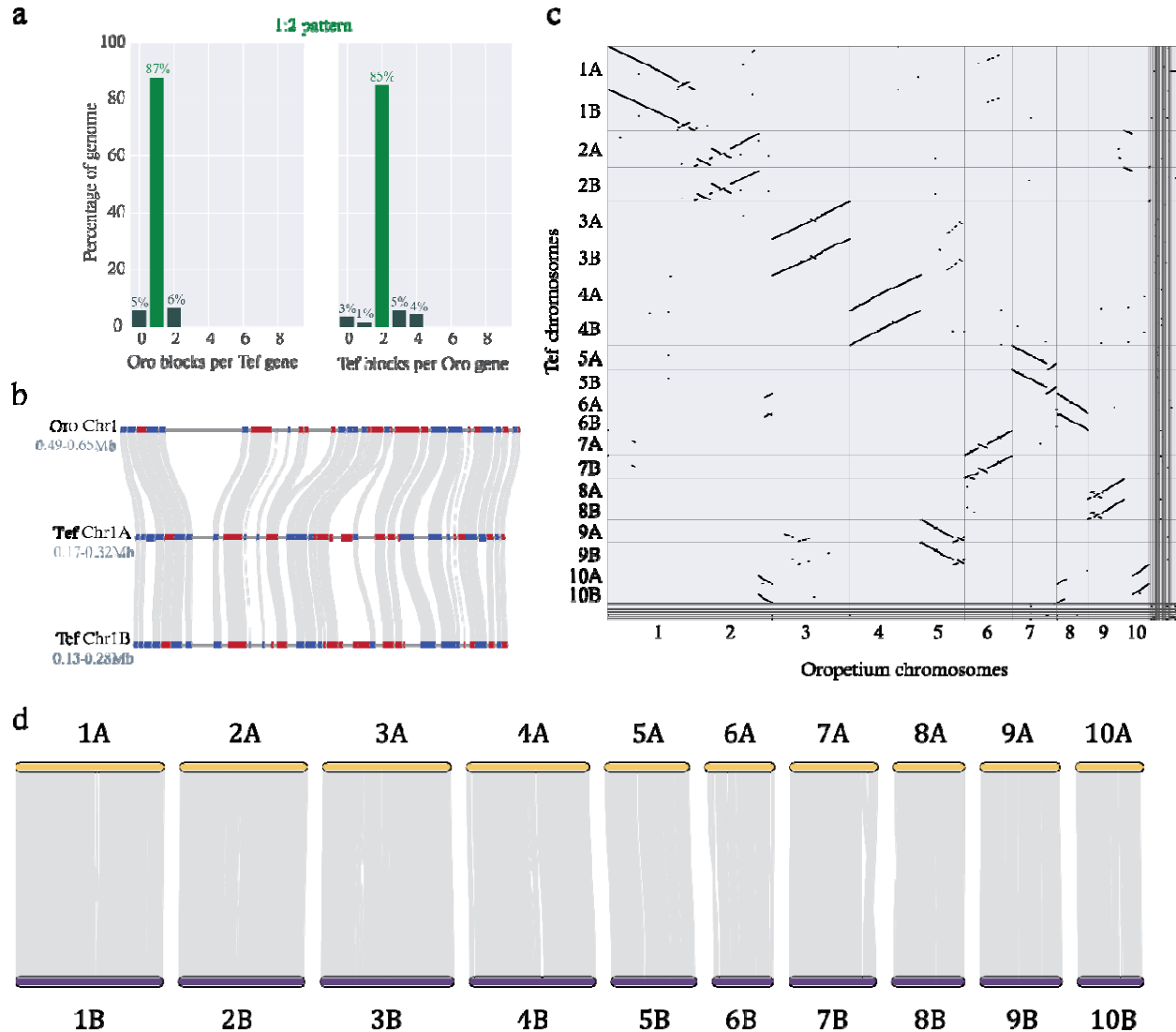
504



505

506 **Figure 2. Insertion dynamics of 65 LTR-RT families in teff.** Box plots of insertion time for
507 the 65 LTR-RT families having ≥ 5 intact LTR elements are plotted. Families 1-5 have ≥ 100
508 intact LTRs, 6-33 have ≥ 10 LTRs, and 34-65 have ≥ 5 LTRs. The six subgenome specific
509 families are highlighted in blue and the estimated range for the teff polyploidy event is shown in
510 brown. A substitution rate of $1.3e-8$ per site per year was used to infer the element insertion
511 times.

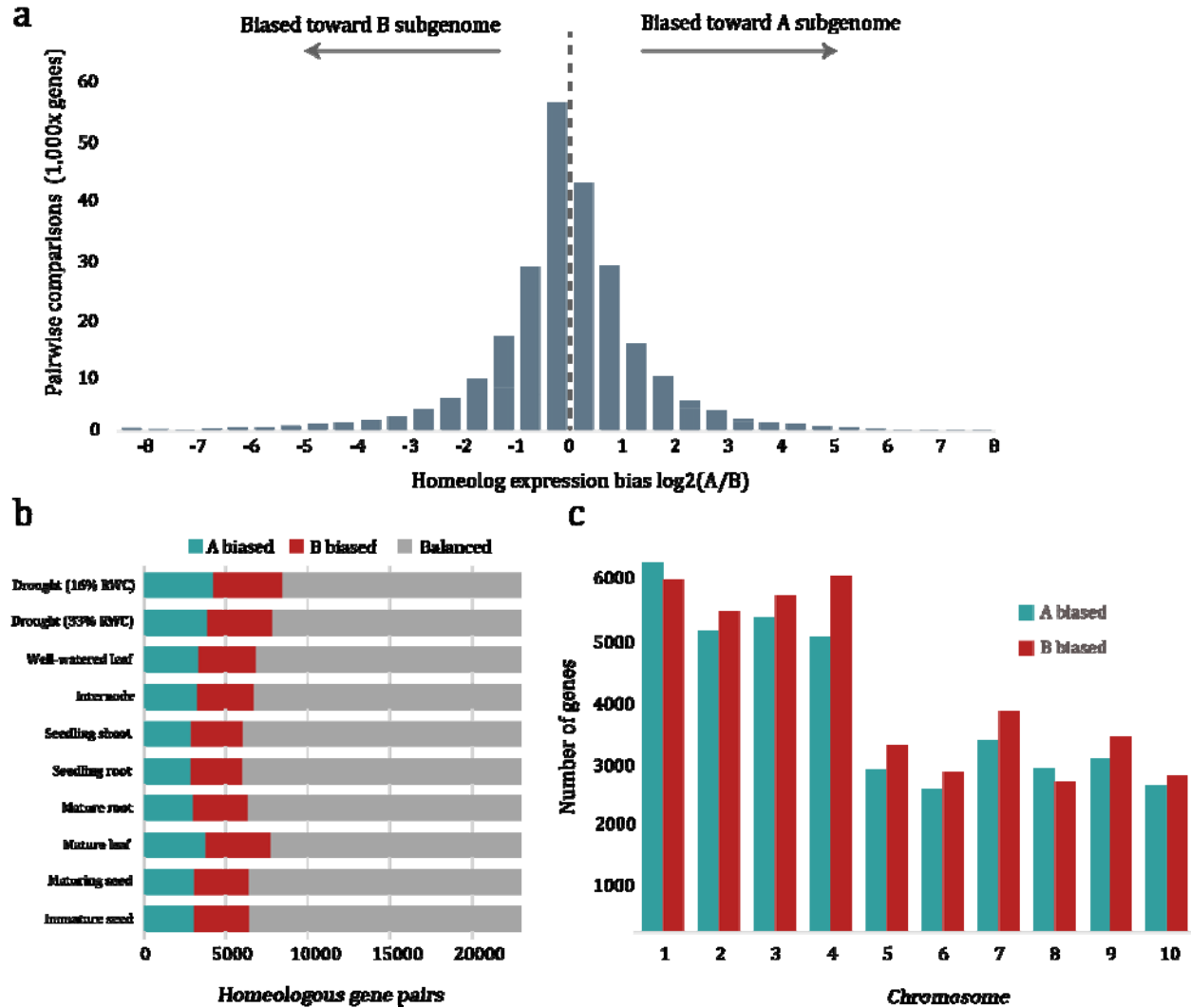
512



513

514 **Figure 3. Comparative genomics of the teff genome.** (a) Ratio of syntenic depth between
 515 Oropetium and teff. Syntenic blocks of Oropetium per teff gene (left) and syntenic blocks of teff
 516 per Oropetium gene (right) are shown indicating a clear 1:2 pattern of Oropetium to teff. (b)
 517 Microsynteny of the teff and Oropetium genomes. A region of the Oropetium Chromosome 1
 518 and the corresponding syntenic regions in homoeologous teff Chromosomes 1 A and B are
 519 shown. Genes are shown in red and blue (for forward and reverse orientation respectively) and
 520 syntenic gene pairs are connected by grey lines. (c) Macrosynteny of the teff and Oropetium
 521 genomes. Syntenic gene pairs are denoted by gray points. (d) Collinearity of the teff subgenomes.
 522 The ten chromosomes belonging to the teff A and B subgenomes are shown in yellow and purple
 523 respectively. Syntenic blocks between homoeologous regions are shown in grey.

524



525

526 **Figure 4. Homoeolog expression bias between the A and B subgenomes of teff.** (a) The
 527 distribution of homoeolog expression bias (HEB) between all gene pairs in all tissues. An HEB >
 528 0 indicates bias toward the A subgenome and a HEB < 0 indicates bias toward the B subgenome.
 529 (b) HEB across the ten tissues in the teff expression atlas. Gene pairs were classified as biased
 530 toward the A (blue) or B (red) subgenomes or balanced with no statistically significant
 531 differential expression (grey). (c) HEB in each of the ten pairs of chromosomes across all ten
 532 tissue types.

533

534 **Table 1.** Summary statistics of the teff genome

Chromosome	Size (bp)	Anchored contigs	Number of genes	Number of Tandem duplicates
1A	40,621,098	35	5,135	465
1B	35,710,944	32	4,829	469
2A	35,425,885	45	4,398	441
2B	30,633,641	23	4,112	382
3A	34,643,735	47	4,415	404
3B	32,575,812	43	4,370	417
4A	32,664,196	39	4,224	318
4B	29,936,223	32	4,127	294
5A	26,945,638	29	2,899	403
5B	24,206,550	36	2,785	385
6A	27,140,163	46	2,409	365
6B	19,415,607	31	1,992	225
7A	26,459,500	44	3,006	315
7B	23,383,462	34	2,843	307
8A	24,151,120	26	2,464	270
8B	21,147,804	28	2,373	239
9A	24,589,398	38	2,736	292
9B	21,940,566	23	2,673	270
10A	23,813,772	24	2,346	268
10B	20,101,091	32	2,151	227
unanchored	22,232,506	657	1,968	130
Total	577,738,711	1,344	68,255	6,886

535

536

537 **References:**

538

- 539 1. Mueller, N.G., Fritz, G.J., Patton, P., Carmody, S. & Horton, E.T. Growing the lost crops of eastern
540 North America's original agricultural system. *Nature plants* **3**, 17092 (2017).
- 541 2. Khoury, C.K. *et al.* Increasing homogeneity in global food supplies and the implications for food
542 security. *Proceedings of the National Academy of Sciences* **111**, 4001-4006 (2014).
- 543 3. Costanza, S., Dewet, J. & Harlan, J.R. Literature review and numerical taxonomy of *Eragrostis tef*
544 (*T'ef*). *Economic Botany* **33**, 413-424 (1979).
- 545 4. Demissie, A. *Tef* genetic resources in Ethiopia. in *Narrowing the Rift: Tef Research and*
546 *Development. Proceedings of the International Workshop on Tef Genetics and Improvement,*
547 *Debre Zeit, Ethiopia* 16-19 (2000).
- 548 5. D'Andrea, A.C. *T'ef* (*Eragrostis tef*) in ancient agricultural systems of highland Ethiopia. *Economic*
549 *Botany* **62**, 547-566 (2008).
- 550 6. Abraham, B. *et al.* The system of crop intensification: reports from the field on improving
551 agricultural production, food security, and resilience to climate change for multiple crops.
552 *Agriculture & Food Security* **3**, 4 (2014).
- 553 7. Cannarozzi, G. *et al.* Genome and transcriptome sequencing identifies breeding targets in the
554 orphan crop *tef* (*Eragrostis tef*). *BMC genomics* **15**, 581 (2014).
- 555 8. Gugsu, L. *et al.* The cytogenetics of *tef*. in *Narrowing the Rift: Tef Research and development.*
556 *Proceedings of the International Workshop on Tef Genetics and Improvement held at Debre Zeit,*
557 *Ethiopia* (2001).
- 558 9. Ingram, A.L. & Doyle, J.J. The origin and evolution of *Eragrostis tef* (Poaceae) and related
559 polyploids: evidence from nuclear *waxy* and plastid *rps16*. *American Journal of Botany* **90**, 116-
560 122 (2003).
- 561 10. Paterson, A.H. *et al.* Repeated polyploidization of *Gossypium* genomes and the evolution of
562 spinnable cotton fibres. *Nature* **492**, 423 (2012).
- 563 11. Osborn, T.C. The contribution of polyploidy to variation in Brassica species. *Physiologia*
564 *Plantarum* **121**, 531-536 (2004).
- 565 12. Ulrich, D. & Olbricht, K. Diversity of volatile patterns in sixteen *Fragaria vesca* L. accessions in
566 comparison to cultivars of *Fragaria x ananassa*. *Journal of Applied Botany and Food Quality*
567 **86**(2013).
- 568 13. Thomas, B.C., Pedersen, B. & Freeling, M. Following tetraploidy in an *Arabidopsis* ancestor,
569 genes were removed preferentially from one homoeolog leaving clusters enriched in dose-
570 sensitive genes. *Genome research* **16**, 934-946 (2006).
- 571 14. Freeling, M. *et al.* Fractionation mutagenesis and similar consequences of mechanisms removing
572 dispensable or less-expressed DNA in plants. *Current opinion in plant biology* **15**, 131-139 (2012).
- 573 15. Edger, P.P. *et al.* Subgenome dominance in an interspecific hybrid, synthetic allopolyploid, and a
574 140-year-old naturally established neo-allopolyploid monkeyflower. *The Plant Cell*, tpc.
575 00010.2017 (2017).
- 576 16. Roodt, R. & Spies, J.J. Chromosome studies in the grass subfamily Chloridoideae. II. An analysis
577 of polyploidy. *Taxon* **52**, 736-746 (2003).
- 578 17. Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and
579 repeat separation. *Genome research* **27**, 722-736 (2017).
- 580 18. Walker, B.J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection and
581 genome assembly improvement. *PLoS one* **9**, e112963 (2014).

- 582 19. VanBuren, R., Wai, C.M., Keilwagen, J. & Pardo, J. A chromosome-scale assembly of the model
583 desiccation tolerant grass *Oropetium thomaeum*. *Plant Direct* **2**, e00096 (2018).
- 584 20. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*
585 *preprint arXiv:1303.3997* (2013).
- 586 21. Durand, N.C. *et al.* Juicer provides a one-click system for analyzing loop-resolution Hi-C
587 experiments. *Cell systems* **3**, 95-98 (2016).
- 588 22. Tavassoli, A. University of London (1986).
- 589 23. Doyle, J.J. & Egan, A.N. Dating the origins of polyploidy events. *New Phytologist* **186**, 73-85
590 (2010).
- 591 24. SanMiguel, P., Gaut, B.S., Tikhonov, A., Nakajima, Y. & Bennetzen, J.L. The paleontology of
592 intergene retrotransposons of maize. *Nature genetics* **20**, 43 (1998).
- 593 25. VanBuren, R. *et al.* Single-molecule sequencing of the desiccation-tolerant grass *Oropetium*
594 *thomaeum*. *Nature* (2015).
- 595 26. Initiative, I.B. Genome sequencing and analysis of the model grass *Brachypodium distachyon*.
596 *Nature* **463**, 763 (2010).
- 597 27. Cotton, J.L. *et al.* Resolving deep relationships of PACMAD grasses: a phylogenomic approach.
598 *BMC plant biology* **15**, 178 (2015).
- 599 28. Hittalmani, S. *et al.* Genome and Transcriptome sequence of Finger millet (*Eleusine coracana* (L.)
600 Gaertn.) provides insights into drought tolerance and nutraceutical properties. *BMC genomics*
601 **18**, 465 (2017).
- 602 29. Schnable, J.C., Springer, N.M. & Freeling, M. Differentiation of the maize subgenomes by
603 genome dominance and both ancient and ongoing gene loss. *Proceedings of the National*
604 *Academy of Sciences* **108**, 4069-4074 (2011).
- 605 30. Wang, X. *et al.* The genome of the mesopolyploid crop species *Brassica rapa*. *Nature genetics* **43**,
606 1035 (2011).
- 607 31. Li, A. *et al.* mRNA and small RNA transcriptomes reveal insights into dynamic homoeolog
608 regulation of allopolyploid heterosis in nascent hexaploid wheat. *The Plant Cell*, tpc. 114.124388
609 (2014).
- 610 32. Douglas, G.M. *et al.* Hybrid origins and the earliest stages of diploidization in the highly
611 successful recent polyploid *Capsella bursa-pastoris*. *Proceedings of the National Academy of*
612 *Sciences*, 201412277 (2015).
- 613 33. Sun, H. *et al.* Karyotype stability and unbiased fractionation in the paleo-allotetraploid *Cucurbita*
614 genomes. *Molecular plant* **10**, 1293-1306 (2017).
- 615 34. Edger, P.P., McKain, M.R., Bird, K.A. & VanBuren, R. Subgenome assignment in allopolyploids:
616 challenges and future directions. *Current opinion in plant biology* **42**, 76-80 (2018).
- 617 35. Ramírez-González, R. *et al.* The transcriptional landscape of polyploid wheat. *Science* **361**,
618 eaar6089 (2018).
- 619 36. Buggs, R.J. *et al.* Tissue-specific silencing of homoeologs in natural populations of the recent
620 allopolyploid *Tragopogon mirus*. *New Phytologist* **186**, 175-183 (2010).
- 621 37. Appels, R. *et al.* Shifting the limits in wheat research and breeding using a fully annotated
622 reference genome. *Science* **361**, eaar7191 (2018).
- 623 38. Chalhoub, B. *et al.* Early allopolyploid evolution in the post-Neolithic *Brassica napus* oilseed
624 genome. *Science* **345**, 950-953 (2014).
- 625 39. Edger, P.P. *et al.* Origin and evolution of the octoploid strawberry genome. *Nature Genetics*
626 (2019).
- 627 40. Wang, M. *et al.* Reference genome sequences of two cultivated allotetraploid cottons,
628 *Gossypium hirsutum* and *Gossypium barbadense*. *Nature genetics*, 1 (2018).
- 629 41. Zou, C. *et al.* The genome of broomcorn millet. *Nature Communications* **10**, 436 (2019).

- 630 42. Bai, G., Tefera, H., Ayele, M. & Nguyen, H. A genetic linkage map of tef [*Eragrostis tef* (Zucc.)
631 Trotter] based on amplified fragment length polymorphism. *Theoretical and Applied Genetics*
632 **99**, 599-604 (1999).
- 633 43. Yu, J.-K., Graznak, E., Breseghello, F., Tefera, H. & Sorrells, M.E. QTL mapping of agronomic traits
634 in tef [*Eragrostis tef* (Zucc) Trotter]. *BMC plant biology* **7**, 30 (2007).
- 635 44. Gaeta, R.T. & Pires, J.C. Homoeologous recombination in allopolyploids: the polyploid ratchet.
636 *New Phytologist* **186**, 18-28 (2010).
- 637 45. Mayrose, I. *et al.* Recently formed polyploid plants diversify at lower rates. *Science* **333**, 1257-
638 1257 (2011).
- 639 46. Wendel, J.F. Genome evolution in polyploids. in *Plant molecular evolution* 225-249 (Springer,
640 2000).
- 641 47. Riley, R. & Chapman, V. Genetic control of the cytologically diploid behaviour of hexaploid
642 wheat. *Nature* **182**, 713 (1958).
- 643 48. Zhang, J. *et al.* Allele-defined genome of the autopolyploid sugarcane *Saccharum spontaneum* L.
644 *Nature genetics* **50**, 1565 (2018).
- 645 49. Zhu, Q. *et al.* High throughput discovery of mutations in tef semi-dwarfing genes by next
646 generation sequencing analysis. *Genetics, genetics*. 112.144436 (2012).
- 647 50. Zhang, H.B., Zhao, X., Ding, X., Paterson, A.H. & Wing, R.A. Preparation of megabase-size DNA
648 from plant nuclei. *The Plant Journal* **7**, 175-184 (1995).
- 649 51. Arumuganathan, K. & Earle, E. Estimation of nuclear DNA content of plants by flow cytometry.
650 *Plant molecular biology reporter* **9**, 229-241 (1991).
- 651 52. Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and
652 repeat separation. *bioRxiv*, 071282 (2017).
- 653 53. Wick, R.R., Schultz, M.B., Zobel, J. & Holt, K.E. Bandage: interactive visualization of de novo
654 genome assemblies. *Bioinformatics* **31**, 3350-3352 (2015).
- 655 54. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
656 data. *Bioinformatics*, btu170 (2014).
- 657 55. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nature methods* **9**,
658 357-359 (2012).
- 659 56. Dudchenko, O. *et al.* De novo assembly of the *Aedes aegypti* genome using Hi-C yields
660 chromosome-length scaffolds. *Science* **356**, 92-95 (2017).
- 661 57. Wang, X. & Wang, L. GMATA: An Integrated Software Package for Genome-Scale SSR Mining,
662 Marker Development and Viewing. *Frontiers in Plant Science* **7**(2016).
- 663 58. Xu, Z. & Wang, H. LTR_FINDER: an efficient tool for the prediction of full-length LTR
664 retrotransposons. *Nucleic Acids Res* **35**, W265-8 (2007).
- 665 59. Ellinghaus, D., Kurtz, S. & Willhoeft, U. LTRharvest, an efficient and flexible software for de novo
666 detection of LTR retrotransposons. *BMC Bioinformatics* **9**, 18 (2008).
- 667 60. Ou, S. & Jiang, N. LTR_retriever: A Highly Accurate and Sensitive Program for Identification of
668 Long Terminal Repeat Retrotransposons. *Plant Physiol* **176**, 1410-1422 (2018).
- 669 61. Wenke, T. *et al.* Targeted identification of short interspersed nuclear element families shows
670 their widespread existence and extreme heterogeneity in plant genomes. *Plant Cell* **23**, 3117-28
671 (2011).
- 672 62. Rho, M. & Tang, H. MGEScan-non-LTR: computational identification and classification of
673 autonomous non-LTR retrotransposons in eukaryotic genomes. *Nucleic Acids Res* **37**, e143
674 (2009).
- 675 63. Han, Y. & Wessler, S.R. MITE-Hunter: a program for discovering miniature inverted-repeat
676 transposable elements from genomic sequences. *Nucleic Acids Res* **38**, e199 (2010).

- 677 64. Crescente, J.M., Zavallo, D., Helguera, M. & Vanzetti, L.S. MITE Tracker: an accurate approach to
678 identify miniature inverted-repeat transposable elements in large genomes. *BMC Bioinformatics*
679 **19**, 348 (2018).
- 680 65. Xiong, W., He, L., Lai, J., Dooner, H.K. & Du, C. HelitronScanner uncovers a large overlooked
681 cache of Helitron transposons in many plant genomes. *Proc Natl Acad Sci U S A* **111**, 10263-8
682 (2014).
- 683 66. Wicker, T. *et al.* A unified classification system for eukaryotic transposable elements. *Nature*
684 *Reviews Genetics* **8**, 973 (2007).
- 685 67. Bao, W., Kojima, K.K. & Kohany, O. Repbase Update, a database of repetitive elements in
686 eukaryotic genomes. *Mob DNA* **6**, 11 (2015).
- 687 68. Tarailo-Graovac, M. & Chen, N. Using RepeatMasker to identify repetitive elements in genomic
688 sequences. *Current protocols in bioinformatics*, 4.10. 1-4.10. 14 (2009).
- 689 69. Luo, M.-C. *et al.* Genome sequence of the progenitor of the wheat D genome *Aegilops tauschii*.
690 *Nature* **551**, 498 (2017).
- 691 70. Dai, X. *et al.* Birth and Death of LTR-Retrotransposons in *Aegilops tauschii*. *Genetics* **210**, 1039-
692 1051 (2018).
- 693 71. Ma, J. & Bennetzen, J.L. Rapid recent growth and divergence of rice nuclear genomes.
694 *Proceedings of the National Academy of Sciences of the United States of America* **101**, 12404-
695 12410 (2004).
- 696 72. Melters, D.P. *et al.* Comparative analysis of tandem repeats from hundreds of species reveals
697 unique insights into centromere evolution. *Genome biology* **14**, R10 (2013).
- 698 73. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic acids research*
699 **27**, 573-580 (1999).
- 700 74. Tamura, K. *et al.* MEGA5: molecular evolutionary genetics analysis using maximum likelihood,
701 evolutionary distance, and maximum parsimony methods. *Molecular biology and evolution* **28**,
702 2731-2739 (2011).
- 703 75. Campbell, M.S. *et al.* MAKER-P: a tool kit for the rapid creation, management, and quality
704 control of plant genome annotations. *Plant physiology* **164**, 513-524 (2014).
- 705 76. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21 (2013).
- 706 77. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from RNA-seq
707 reads. *Nature biotechnology* **33**, 290-295 (2015).
- 708 78. Lamesch, P. *et al.* The Arabidopsis Information Resource (TAIR): improved gene annotation and
709 new tools. *Nucleic acids research* **40**, D1202-D1210 (2011).
- 710 79. Goff, S.A. *et al.* A draft sequence of the rice genome (*Oryza sativa* L. ssp. japonica). *Science* **296**,
711 92-100 (2002).
- 712 80. Paterson, A.H. *et al.* The Sorghum bicolor genome and the diversification of grasses. *Nature* **457**,
713 551 (2009).
- 714 81. Boutet, E., Lieberherr, D., Tognolli, M., Schneider, M. & Bairoch, A. Uniprotkb/swiss-prot. in
715 *Plant bioinformatics* 89-112 (Springer, 2007).
- 716 82. Korf, I. Gene finding in novel genomes. *BMC bioinformatics* **5**, 59 (2004).
- 717 83. Stanke, M. & Waack, S. Gene prediction with a hidden Markov model and a new intron
718 submodel. *Bioinformatics* **19**, ii215-ii225 (2003).
- 719 84. Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. & Zdobnov, E.M. BUSCO:
720 assessing genome assembly and annotation completeness with single-copy orthologs.
721 *Bioinformatics* **31**, 3210-3212 (2015).
- 722 85. Bray, N.L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq
723 quantification. *Nature biotechnology* **34**, 525 (2016).

- 724 86. Edgar, R.C. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
725 *Nucleic acids research* **32**, 1792-1797 (2004).
726 87. Suyama, M., Torrents, D. & Bork, P. PAL2NAL: robust conversion of protein sequence alignments
727 into the corresponding codon alignments. *Nucleic acids research* **34**, W609-W612 (2006).
728 88. Yang, Z. PAML: a program package for phylogenetic analysis by maximum likelihood.
729 *Bioinformatics* **13**, 555-556 (1997).
730