

1

2 **A chromosome scale assembly of the model desiccation tolerant grass *Oropetium thomaeum***

3

4

5 Robert VanBuren^{1,2*}, Ching Man Wai¹, Jens Keilwagen³, Jeremy Pardo⁴

6

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

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

9 ³Institute for Biosafety in Plant Biotechnology, Julius Kühn-Institut (JKI) - Federal Research Centre for
10 Cultivated Plants, Quedlinburg, Germany

11 ⁴Department of Plant Biology, Michigan State University, East Lansing, MI, 48824, USA

12 *corresponding author: bobvanburen@gmail.com

14 A preprint version of this manuscript is available at: <https://www.biorxiv.org/content/early/2018/07/31/378943>

16 **Abstract**

17 *Oropetium thomaeum* is an emerging model for desiccation tolerance and genome size evolution
18 in grasses. A draft genome of *Oropetium* was recently sequenced, but the lack of a chromosome
19 scale assembly has hindered comparative analyses and downstream functional genomics. Here,
20 we reassembled *Oropetium*, and anchored the genome into ten chromosomes using Hi-C based
21 chromatin interactions. A combination of high-resolution RNAseq data and homology-based
22 gene prediction identified thousands of new, conserved gene models that were absent from the
23 V1 assembly. This includes thousands of new genes with high expression across a desiccation
24 timecourse. Comparison between the sorghum and *Oropetium* genomes revealed a surprising
25 degree of chromosome-level collinearity, and several chromosome pairs have near perfect
26 synteny. Other chromosomes are collinear in the gene rich chromosome arms but have
27 experienced pericentric translocations. Together, these resources will be useful for the grass
28 comparative genomic community and further establish *Oropetium* as a model resurrection plant.

30 **Keywords:** grasses, comparative genomics, Hi-C, desiccation tolerance, chromosome-scale

31

32

33

34

35

36 **Introduction**

37 Desiccation tolerance evolved as an adaptation to extreme and prolonged drying, and
38 resurrection plants are among the most resilient plants on the planet. The molecular basis of
39 desiccation tolerance is still largely unknown, but a number of models have emerged to dissect
40 the genetic control of this trait (Hoekstra et al., 2001; Zhang and Bartels, 2018). The genomes of
41 several model resurrection plants have been sequenced including *Boea hygrometrica* (Xiao et al.,
42 2015), *Oropetium thomaeum* (VanBuren et al., 2015), *Xerophyta viscosa* (Costa et al., 2017),
43 *Lindernia brevidens* (VanBuren et al., 2018a), *Selaginella lepidophylla* (VanBuren et al., 2018b),
44 and *Selaginella tamariscina* (Xu et al., 2018). To date, no chromosome scale assemblies are
45 available for these species, limiting large-scale quantitative genetics and comparative genomics-
46 based approaches. Many resurrection plants are polyploid or have prohibitively large genomes
47 including those in the genera *Boea*, *Xerophyta*, *Eragrostis*, *Sporobolus*, and *Craterostigma*. This
48 complexity complicates genome assembly and gene redundancy in the polyploid species hinders
49 downstream functional genomics work.

50 *Oropetium thomaeum* (hereon referred to as Oropetium), is a diploid resurrection plant
51 with the smallest genome among the grasses (245 Mb) (Bartels and Mattar, 2002). Oropetium
52 plants are similar in size to *Arabidopsis*, but significantly smaller than the model grasses *Setaria*
53 *italica* (Li and Brutnell, 2011) and *Brachypodium distachyon* (Brkljacic et al., 2011), with a short
54 generation time of ~4 months. Oropetium is in the Chloridoideae subfamily of grasses and is
55 closely related to the orphan cereal crops tef (*Eragrostis tef*) and finger millet (*Eleusine*
56 *coracana*). Desiccation tolerance evolved independent several times within Chloridoideae (Gaff,
57 1977; Gaff and Latz, 1978; Gaff, 1987) making it a useful system for studying convergent
58 evolution. Together, these traits make Oropetium an attractive model for exploring the origin and
59 molecular basis of desiccation tolerance. Oropetium was one of the first plants to be sequenced
60 using the long reads of PacBio technology, and the assembly quality was comparable to early
61 Sanger sequencing based plant genomes such as rice and *Arabidopsis* (VanBuren et al., 2015).
62 Despite the high contiguity of Oropetium V1, the assembly has 625 contigs and the BioNano
63 based genome map was unable to produce chromosome-scale scaffolds. Furthermore, the V1
64 annotation was based on limited transcript evidence, and a high proportion of conserved plant
65 genes were missing (VanBuren et al., 2015). Here, we reassembled the Oropetium genome using
66 a more refined algorithm and generated a chromosome scale assembly using Hi-C based
67 chromatin interactions. The annotation quality was improved using high-resolution RNAseq data
68 and protein homology, facilitating detailed comparative genomics with other grasses.

69

70 **Results and Discussion**

71 The first version of the Oropetium genome (V1) was sequenced with high coverage PacBio data
72 (~72x) followed by error correction and assembly using the hierarchical genome assembly
73 process (HGAP) (VanBuren et al., 2015). We reassembled this PacBio data using the Canu
74 assembler (Koren et al., 2017), which can more accurately assemble and phase complex
75 repetitive regions. The resulting Canu based assembly (hereon referred to as V1.2) had fewer

76 contigs than the V1 HGAP assembly, but had otherwise similar assembly metrics (Table 1).
77 Draft contigs were polished using a two-step process to remove residual insertion/deletion (indel)
78 and single nucleotide errors. Contigs were first polished using the raw PacBio data with Quiver
79 (Chin et al., 2013), followed by four rounds of reiterative polishing with Pilon (Walker et al.,
80 2014) using high coverage Illumina paired end data. The final V1.2 assembly contains 436
81 contigs with an N50 of 2.0 Mb and total assembly size of 236 Mb. This is six megabases smaller
82 than the V1 assembly, with slightly lower contiguity. More intact long terminal repeat
83 retrotransposons (LTR-RTs) and centromere specific repeat arrays were identified in *Oropetium*
84 V1.2 compared to V1, suggesting the Canu assembler resolved these repetitive elements more
85 accurately. Thus, V1.2 was used for pseudomolecule construction.

86 The *Oropetium* V1.2 contigs were ordered and oriented into chromosome-scale
87 pseudomolecules using high-throughput chromatin conformation capture (Hi-C). Hi-C leverages
88 long-range interactions across distal regions of chromosomes to order and orient contigs. This
89 approach is similar to genetic map-based anchoring, but with higher resolution. Illumina data
90 generated from the Hi-C library was mapped to the V1.2 *Oropetium* genome using bwa (Li,
91 2013) and the proximity-based clustering matrix was generated using the Juicer and 3d-DNA
92 pipelines (Durand et al., 2016; Dudchenko et al., 2017). After filtering and manual curation, ten
93 high confidence clusters were identified (Figure 1). These ten clusters correspond to the haploid
94 chromosome number of *Oropetium*. Regions with low density interactions highlight the
95 centromeric and pericentromeric regions, and regions with higher than expected interactions
96 represent topologically associated domains. After splitting six chimeric PacBio contigs, 239
97 contigs were anchored and oriented into ten chromosomes spanning 226.5 Mb or 95.8 % of the
98 total assembled genome (Table 1). Chromosomes range in size from 11.0 to 34.7 Mb with an
99 average size of 22.6 Mb. Most of the unanchored contigs are small (average size 42kb), or are
100 entirely composed of rRNA, centromeric repeat arrays, or centromere specific LTR-RTs.
101 Telomeres were identified at both ends of Chromosomes 1, 2, 3, 4, 5, 7, and 9 and on one end of
102 Chromosomes 6, 8, and 10. Three unanchored contigs contain the remaining telomeres. This
103 supports the completeness and accuracy of the pseudomolecule construction.

104 The chromosome scale *Oropetium* genome (hereon referred to as V2) was reannotated
105 using the homology-based gene prediction program GeMoMa (Keilwagen et al., 2016;
106 Keilwagen et al., 2018). Protein coding sequences from 11 angiosperm genomes and RNAseq
107 data from *Oropetium* (VanBuren et al., 2017) were used as evidence. After filtering gene models
108 derived from transposases, the final annotation consists of 28,835 high-confidence gene models.
109 The annotation completeness was assessed using the Benchmarking Universal Single-Copy
110 Ortholog (BUSCO) embryophyta dataset. The V2 gene models have a BUSCO score of 98.9%,
111 suggesting the updated annotation is high-quality. In comparison, the *Oropetium* V1 annotation
112 has a BUSCO score of 72%, and many conserved gene models were likely missing or mis-
113 annotated. Nearly forty percent (11,227) of the gene models in V2 are new and were unannotated
114 in V1. In addition, 10,837 gene models from V1 were removed or substantially improved in the
115 V2 annotation. These discarded gene models either had little support based on protein homology
116 to other species and transcript evidence from *Oropetium*, or they were misannotated transposable
117 elements. In total, 94.3% of the gene models (27,216) were anchored to the ten chromosomes.

118 Among the newly annotated gene models are 3,525 tandem gene duplicates (Figure 2a). Tandem
119 duplicates span 3,062 arrays with 7,760 total genes. Of the arrays containing three or more
120 genes, only 49 are new to V2, and the majority contain genes previously identified in V1. The
121 boundaries of tandem duplicates are difficult to correctly annotate, resulting in fusions of two or
122 more gene copies. The homology-based annotation used in V2 was able to parse previously
123 fused gene models.

124 The expression pattern of newly annotated genes were surveyed using high-resolution
125 RNAseq expression data (VanBuren et al., 2017). This dataset consists of seven leaf samples
126 collected during desiccation and rehydration timecourses. Timepoints include well-watered, 7,
127 14, 21, and 30 days desiccated as well as 24 and 48 hours post rehydration. Differentially
128 expressed genes were identified based on comparisons of well-watered leaves with each
129 dehydration or rehydration timepoint. In addition, each timepoint was compared with the
130 timepoint immediately following it in the timecourse (i.e. day 7 dehydration vs day 14). In total,
131 17,204 gene models from the V1 annotation had detectable expression (count > 0 in at least one
132 sample) compared to 25,314 gene models in V2 (Figure 2b). Of the expressed genes, 9,149 V1
133 and 11,948 V2 gene models were classified as differentially expressed in at least one of the
134 comparisons. Most newly annotated genes (8,110) have detectable expression in at least one of
135 the seven timepoints, and the majority are expressed in all timepoints. In total, 2,799 new V2
136 gene models were differentially expressed, suggesting the newly annotated genes have important
137 and previously uncharacterized roles in desiccation tolerance.

138 We used the chromosome scale assembly of *Oroepetium* to survey patterns of genome
139 organization and evolution related to maintaining a small genome size. The proportion of LTR-
140 RTs in *Oropetium* V1 and V2 is similar, though V2 has more intact elements. LTR-RTs are the
141 most abundant repetitive elements in *Oropetium* and collectively span 27% (62 Mb) of the
142 genome. LTR-RTs are distributed non-randomly across the genome, and peaks of Gypsy LTR-
143 RTs are observed in each of the ten chromosomes (Figure 3). These peaks of Gypsy LTR-RTs
144 correspond to the pericentromeric regions. The pericentromeric regions show reduced
145 intrachromosomal interactions in the Hi-C matrix and contain arrays of centromeric repeats. The
146 *Oropetium* V2 genome contains 8,965 155 bp monomeric centromeric repeats; considerably
147 more than the 4,315 identified in the V1 assembly. The centromeric array sizes vary from 61 kb
148 in chromosome 10 to 1,598 kb in Chromosome 4 (Figure 3; Table 2). Array sizes are likely
149 underestimated, as only 52% of centromeric arrays were anchored to chromosomes, and 23
150 unanchored contigs contain centromeric repeat arrays. Gene density is low in the
151 pericentromeric regions, consistent with the rice, sorghum, maize, and *Brachypodium* genomes
152 (Paterson et al., 2009; Initiative, 2010; Du et al., 2017; Jiao et al., 2017). Collectively,
153 pericentromeric regions span 67.5 Mb or 29% of the genome, a much smaller proportion than
154 sorghum (62%; 460 Mb) (Paterson et al., 2009), but higher than rice (15%; 63 Mb) (Goff et al.,
155 2002). The majority of intact LTRs (86%; 628) have an insertion time of less than one million
156 years ago, with a steep drop off of insertion time after 0.4 MYA. This suggests LTRs are highly
157 active in *Oropetium* but rapidly fragmented and purged to maintain its small genome size.

158 Previous comparative genomics analyses supported a high degree of collinearity between
159 Oropetium and other grass genomes, but the draft assembly prevented detailed chromosome
160 level comparisons. To date, no chromosome scale assemblies are available for other
161 Chloridoideae grasses, though draft genomes are available for the orphan grain crops tef
162 (*Eragrostis tef*) (Cannarozzi et al., 2014) and finger millet (*Eleusine coracana*) (Hittalmani et
163 al., 2017). We compared the V2 Oropetium chromosomes to the high-quality BTX 623 Sorghum
164 genome (McCormick et al., 2018). Sorghum is in the Panicoideae subfamily of grasses which
165 diverged from the ancestors of Chloridoideae ~31 MYA (Cotton et al., 2015). Despite this
166 divergence, the ten chromosomes in Oropetium are largely collinear to the corresponding ten
167 chromosomes in Sorghum, though large-scale inversions and translocations were identified
168 (Figure 4a). Oropetium chromosomes 5, 6, and 8 are collinear along their length to sorghum
169 chromosomes 9, 6, and 5 respectively. Oropetium chromosomes 1, 2, 4, and 7, are collinear to
170 the arms of sorghum chromosomes 4, 10, 1, and 2, but the pericentric regions have translocated
171 to other chromosomes. Oropetium chromosome 9 and sorghum chromosome 7 are syntenic but
172 have two large-scale inversions, and Oropetium and sorghum chromosome 3 are syntenic with
173 one inversion.

174 The sorghum genome is roughly three fold larger than Oropetium, and genome size
175 dynamics in grasses are driven by purge and accumulation of retrotransposons (Wicker et al.,
176 2010). Gene rich regions of Oropetium are 2-3x more compact than orthologous regions in
177 sorghum, and much of this expansion in sorghum is caused by intergenic blocks of LTR-RTs
178 (Figure 4b), consistent with patterns observed in the V1 assembly (VanBuren et al., 2015). The
179 chromosome-scale nature of Oropetium V2 allowed us to survey patterns of collinearity in the
180 pericentromeric regions. These regions have a lower degree of synteny with sorghum compared
181 to gene rich euchromatin, consistent with retrotransposon-mediated rearrangements (Figure 4b).
182 Pericentromeres are greatly expanded in Oropetium compared to the gene rich euchromatic
183 blocks, similar to patterns observed in sorghum.

184 The read lengths of third generation sequencing technologies enable near gapless
185 assemblies with high contiguity for virtually any plant genome. PacBio and Nanopore based
186 genomes have a better representation of gene and regulatory sequences, but often lack the
187 chromosome-scale scaffolding required for comparative genomics and quantitative genetics. The
188 pseudomolecules in Oropetium V2 allowed us to more accurately identify syntenic orthologs in
189 other grasses and make detailed comparisons of chromosome evolution. The V2 chromosome-
190 scale assembly will serve as a reference for future population genomics work and positional
191 cloning of desiccation related genes. Together, this highlights the need to improve and scaffold
192 existing high-quality reference genomes.

193

194 **Methods**

195 *Genome reassembly*

196 The raw PacBio reads from the Oropetium V1 release (VanBuren et al., 2015) were reassembled
197 with improved algorithms to better resolve highly complex and repetitive regions. PacBio data

198 was error corrected and assembled using Canu (V1.4)(Koren et al., 2017) with the following
199 modifications: minReadLength=1500, GenomeSize=245Mb, minOverlapLength=1000. Other
200 parameters were left as default. The resulting assembly graph was visualized in Bandage (Wick
201 et al., 2015). The assembly graph was free of heterozygosity related bubbles, but many nodes
202 (contigs) were interconnected by a high copy number retrotransposon. The Canu based contigs
203 (assembly V1.2) were first polished using Quiver(Chin et al., 2013) with the raw PacBio data
204 and default parameters. Contigs were further polished with Pilon (V1.22)(Walker et al., 2014)
205 using ~120x coverage of paired-end 150 bp Illumina data. Quality-trimmed Illumina reads were
206 aligned to the draft contigs using bowtie2 (V2.3.0) (Langmead and Salzberg, 2012) with default
207 parameters. The overall alignment rate was 95.5%, which was slightly higher than alignment
208 against the HGAP V1 assembly (94.5%). The following parameters for Pilon were modified: --
209 flank 7, --K 49, and --mindepth 25. Other parameters were left as default. Pilon was run four
210 times with an updated reference and realignment of Illumina data after each iteration. Indel
211 corrections plateaued after the third iteration, suggesting polishing removed most residual
212 assembly errors.

213

214 *HiC library construction analysis, and genome anchoring*

215 Oropetium descended from the original plants used for PacBio sequencing were collected for Hi-
216 C library construction and RNAseq. Oropetium is highly selfing with low heterozygosity, and we
217 expect minimal differences to be introduced in the new version. Oropetium seeds are available
218 upon request. Oropetium plants were maintained under day/night temperatures of 26 and 22°C
219 respectively, with a light intensity of 200 $\mu\text{E m}^{-2} \text{ sec}^{-1}$ and 16/8 hr photoperiod. Young leaf
220 tissue was used for HiC library construction with the Proximo™ Hi-C Plant kit (Phase
221 Genomics) following the manufactures protocol. Briefly, 0.2 grams of fresh, young leaf tissue
222 was finely chopped and the chromatin was immediately crosslinked. The chromatin was
223 fragmented and proximity ligated, followed by library construction. The final library was size
224 selected for 300-600 bp and sequenced on the Illumina HiSeq 4000 under paired-end 150 bp
225 mode. Adapters were trimmed and low-quality sequences were removed using Trimmomatic
226 (V0.36) (Bolger et al., 2014). Read pairs were aligned to the Oropetium contigs using bwa
227 (V0.7.16)(Li, 2013) with strict parameters (-n 0) to prevent mismatches and non-specific
228 alignments in duplicated and repetitive regions. SAM files from bwa were used as input in the
229 Juicer pipeline, and PCR duplicates with the same genome coordinates were filtered prior to
230 constructing the interaction based distance matrix. In total, 101 filtered read pairs were used as
231 input for the Juicer and 3d-DNA HiC analysis and scaffolding pipelines (Durand et al., 2016;
232 Dudchenko et al., 2017). Contig ordering, orientation, and chimera splitting was done using the
233 3d-DNA pipeline(Dudchenko et al., 2017) under default parameters. Contig misassemblies and
234 scaffold misjoins were manually detected and corrected based on interaction densities from
235 visualization in Juicebox. In total, six chimeric contigs were identified and split at the junction
236 with closest interaction data. The manually validated assembly was used as input to build the ten
237 scaffolds (chromosomes) using the finalize-output.sh script from 3d-DNA. Chromosomes and
238 unanchored contigs were renamed by size, producing the V2 assembly.

239

240 *Genome annotation*

241 The Oropetum V2 assembly was reannotated using the homology-based gene prediction program
242 Gene Model Mapper (GeMoMa: V 1.5.2) (Keilwagen et al., 2016; Keilwagen et al., 2018).
243 GeMoMa uses protein homology and RNAseq evidence to predict gene models. Genome
244 assemblies and gene annotation for the following 11 species were downloaded from Phytozome
245 (V12) and used as homology based evidence: *Arabidopsis thaliana*, *Brachypodium distachyon*,
246 *Glycine max*, *Oryza sativa*, *Panicum hallii*, *Populus trichocarpa*, *Prunus persica*, *Setaria italica*,
247 *Solanum lycopersicum*, *Sorghum bicolor*, *Theobroma cacao*. Translated coding exons and
248 proteins from the reference gene annotations and genome assemblies were extracted using the
249 module Extractor function of GeMoMa (module Extractor: Ambiguity=AMBIGUOUS, r=true).
250 RNAseq data from Oropetium desiccation and rehydration timecourses (VanBuren et al., 2017)
251 was aligned to the V2 Oropetium genome using HISAT2 (Kim et al., 2015) with default
252 parameters. The resulting BAM files were used to extract intron and exon boundaries using the
253 module ERE (module ERE: s=FR_FIRST_STRAND, c=true). translated coding exons from
254 other species were aligned to the Oropetium genome using tblastn and transcripts were predicted
255 based on each reference species independently using the extracted introns and coverage (module
256 GeMoMa). Finally, the predictions based on the 11 reference species were combined to obtain a
257 final prediction using the module GAF. Gene models containing transposases were filtered,
258 resulting in a final annotation of 28,835 gene models. The annotation completeness was assessed
259 using the plant specific Benchmarking Universal Single-Copy Ortholog (BUSCO) dataset
260 (version 3.0.2, embryophyta_odb9) (Simão et al., 2015). The following report was obtained from
261 BUSCO: 98.9% overall, 95.4% single copy, 3.5% duplicated, 0.6% fragmented, 0.5% missing.
262 Gene model names from V1 were conserved where possible, and new gene models received new
263 names.

264

265 *Expression analysis*

266 Oropetium RNAseq data from desiccation and rehydration timecourses was reanalyzed using the
267 updated gene model annotations (VanBuren et al., 2017). Four time points during dehydration
268 (days 7, 14, 21, and 30), two during rehydration (24 and 48 hours), and one well-watered sample
269 were analyzed. Based on principle component analysis, replicate 2 of the ‘well-watered and
270 ‘D21’ samples were excluded from the analysis. Each other timepoint had three replicates. Gene
271 expression was quantified on a transcript level using salmon (v 0.9.1) in quasi-mapping mode
272 (Patro et al., 2017). Default parameters were used with the internal GC bias correction in salmon.
273 The R package tximport (v 1.2.0) was used to map transcript level quantifications to gene level
274 counts (Team, 2013; Soneson et al., 2015). We conducted differential expression analysis with
275 the remaining samples using the R package DESeq2 (v 1.14.1) set to default parameters [3,4].

276

277 *Identification of LTR-RTs*

278 A preliminary list of candidate long terminal repeat retrotransposons (LTR-RTs) from
279 Oropetium were identified using LTR_Finder (V1.02) (Xu and Wang, 2007) and LTRharvest
280 (Ellinghaus et al., 2008). The following parameters for LTRharvest were modified: -similar 90 –
281 vic 10 –seed 20 –minlenltr 100 –maxlenltr 7000 –mintsd 4 –maxtsd 6 –motif TGCA –motifmis
282 1. LTR_Finder parameters were: -D 15000 –d 1000 –L 7000 –l 100 –p 20 –C –M 0.9.
283 LTR_retriever(Ou and Jiang, 2017) was used to filter out false LTR retrotransposons using the
284 target site duplications, terminal motifs, and Pfam domains. Default parameters were used for
285 LTRretriever. LTRretirever produced a list of full length, high-quality LTRs. LTRs were
286 annotated across the genome using RepeatMasker (<http://www.repeatmasker.org/>)(Smit et al.,
287 1996) and the non-redundant LTR-RT library constructed by LTR_retriever. The insertion time
288 of intact LTRs was calculated in LTR_retriever using the formula $T=K/2\mu$ with a neutral
289 mutation rate of $\mu=1 \times 10^{-8}$ mutations per bp per year.

290

291 *Comparative genomics*

292 Syntenic gene pairs between the Oropetium and Sorghum genomes were identified using the
293 MCSCAN toolkit (V1.1) (Wang et al., 2012) implemented in python
294 ([https://github.com/tanghaibao/icvi/wiki/MCscan-\(Python-version\)](https://github.com/tanghaibao/icvi/wiki/MCscan-(Python-version))). Default parameters were
295 used. Gene models were aligned using LAST and hits were filtered to find the best 1:1 syntenic
296 blocks. Macrosyntenic dotplots were constructed in MCScan.

297

298 **Availability of supporting data:**

299 The V2 Oropetium genome assembly and updated annotation can be downloaded from CoGe
300 (<https://genomevolution.org/coge>) under Genome ID 51527 and from Phytozome
301 (<https://phytozome.jgi.doe.gov/pz/portal.html>). The raw Hi-C Illumina data has been deposited
302 on the Short Read Archive (SRA) under NCBI BioProject ID PRJNA481965.

303

304 **Acknowledgements:**

305 This work is supported by funding from the National Science Foundation (MCB-1817347 to
306 R.V.).

307

308 **Author contributions:**

309 R.V. designed research; R.V., C.M.W, J.K. and J.P. performed research and/or analyzed data;
310 and R.V. wrote the paper. All authors reviewed the manuscript.

311

312 A preprint version of this manuscript is available at:
313 <https://www.biorxiv.org/content/early/2018/07/31/378943>

- 314 **References:**
- 315
- 316 **Bartels, D., and Mattar, M.** (2002). *Oropetium thomaeum*: A resurrection grass with a diploid genome.
317 *Maydica* **47**, 185-192.
- 318 **Bolger, A.M., Lohse, M., and Usadel, B.** (2014). Trimmomatic: a flexible trimmer for Illumina sequence
319 data. *Bioinformatics*, btu170.
- 320 **Brkljacic, J., Grotewold, E., Scholl, R., Mockler, T., Garvin, D.F., Vain, P., Brutnell, T., Sibout, R., Bevan,**
321 **M., and Budak, H.** (2011). Brachypodium as a model for the grasses: today and the future. *Plant
322 Physiology*, pp. 111.179531.
- 323 **Cannarozzi, G., Plaza-Wüthrich, S., Esfeld, K., Larti, S., Wilson, Y.S., Girma, D., de Castro, E., Chanyalew,**
324 **S., Blösch, R., and Farinelli, L.** (2014). Genome and transcriptome sequencing identifies breeding
325 targets in the orphan crop tef (*Eragrostis tef*). *BMC genomics* **15**, 581.
- 326 **Chin, C.-S., Alexander, D.H., Marks, P., Klammer, A.A., Drake, J., Heiner, C., Clum, A., Copeland, A.,
327 Huddleston, J., and Eichler, E.E.** (2013). Nonhybrid, finished microbial genome assemblies from
328 long-read SMRT sequencing data. *Nature methods* **10**, 563-569.
- 329 **Costa, M., Artur, M., Maia, J., Jonkheer, E., Derkx, M., Nijveen, H., Williams, B., Mundree, S.G.,
330 Jiménez-Gómez, J.M., and Hesselink, T.** (2017). A footprint of desiccation tolerance in the
331 genome of *Xerophyta viscosa*. *Nature plants* **3**, 17038.
- 332 **Cotton, J.L., Wysocki, W.P., Clark, L.G., Kelchner, S.A., Pires, J.C., Edger, P.P., Mayfield-Jones, D., and
333 Duvall, M.R.** (2015). Resolving deep relationships of PACMAD grasses: a phylogenomic
334 approach. *BMC plant biology* **15**, 178.
- 335 **Du, H., Yu, Y., Ma, Y., Gao, Q., Cao, Y., Chen, Z., Ma, B., Qi, M., Li, Y., and Zhao, X.** (2017). Sequencing
336 and de novo assembly of a near complete indica rice genome. *Nature Communications* **8**, 15324.
- 337 **Dudchenko, O., Batra, S.S., Omer, A.D., Nyquist, S.K., Hoeger, M., Durand, N.C., Shamim, M.S.,
338 Machol, I., Lander, E.S., and Aiden, A.P.** (2017). De novo assembly of the *Aedes aegypti* genome
339 using Hi-C yields chromosome-length scaffolds. *Science* **356**, 92-95.
- 340 **Durand, N.C., Shamim, M.S., Machol, I., Rao, S.S., Huntley, M.H., Lander, E.S., and Aiden, E.L.** (2016).
341 Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell systems* **3**,
342 95-98.
- 343 **Ellinghaus, D., Kurtz, S., and Willhoeft, U.** (2008). LTRharvest, an efficient and flexible software for de
344 novo detection of LTR retrotransposons. *BMC bioinformatics* **9**, 18.
- 345 **Gaff, D.** (1977). Desiccation tolerant vascular plants of Southern Africa. *Oecologia* **31**, 95-109.
- 346 **Gaff, D.** (1987). Desiccation tolerant plants in South America. *Oecologia* **74**, 133-136.
- 347 **Gaff, D., and Latz, P.** (1978). The occurrence of resurrection plants in the Australian flora. *Australian
348 Journal of Botany* **26**, 485-492.
- 349 **Goff, S.A., Ricke, D., Lan, T.-H., Presting, G., Wang, R., Dunn, M., Glazebrook, J., Sessions, A., Oeller, P.,
350 and Varma, H.** (2002). A draft sequence of the rice genome (*Oryza sativa* L. ssp. *japonica*).
351 *Science* **296**, 92-100.
- 352 **Hittalmani, S., Mahesh, H., Shirke, M.D., Biradar, H., Uday, G., Aruna, Y., Lohithaswa, H., and
353 Mohanrao, A.** (2017). Genome and Transcriptome sequence of Finger millet (*Eleusine coracana*
354 (L.) Gaertn.) provides insights into drought tolerance and nutraceutical properties. *BMC
355 genomics* **18**, 465.
- 356 **Hoekstra, F.A., Golovina, E.A., and Buitink, J.** (2001). Mechanisms of plant desiccation tolerance. *Trends
357 in plant science* **6**, 431-438.
- 358 **Initiative, I.B.** (2010). Genome sequencing and analysis of the model grass *Brachypodium distachyon*.
359 *Nature* **463**, 763.

- 360 **Jiao, Y., Peluso, P., Shi, J., Liang, T., Stitzer, M.C., Wang, B., Campbell, M.S., Stein, J.C., Wei, X., and**
361 **Chin, C.-S.** (2017). Improved maize reference genome with single-molecule technologies.
362 *Nature*.
- 363 **Keilwagen, J., Hartung, F., Paulini, M., Twardziok, S.O., and Grau, J.** (2018). Combining RNA-seq data
364 and homology-based gene prediction for plants, animals and fungi. *BMC bioinformatics* **19**, 189.
- 365 **Keilwagen, J., Wenk, M., Erickson, J.L., Schattat, M.H., Grau, J., and Hartung, F.** (2016). Using intron
366 position conservation for homology-based gene prediction. *Nucleic acids research* **44**, e89-e89.
- 367 **Kim, D., Langmead, B., and Salzberg, S.L.** (2015). HISAT: a fast spliced aligner with low memory
368 requirements. *Nature methods* **12**, 357.
- 369 **Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H., and Phillippy, A.M.** (2017). Canu: scalable
370 and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome*
371 *research* **27**, 722-736.
- 372 **Langmead, B., and Salzberg, S.L.** (2012). Fast gapped-read alignment with Bowtie 2. *Nature methods* **9**,
373 357-359.
- 374 **Li, H.** (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*
375 preprint arXiv:1303.3997.
- 376 **Li, P., and Brutnell, T.P.** (2011). *Setaria viridis* and *Setaria italica*, model genetic systems for the Panicoideae
377 grasses. *Journal of experimental botany* **62**, 3031-3037.
- 378 **McCormick, R.F., Truong, S.K., Sreedasyam, A., Jenkins, J., Shu, S., Sims, D., Kennedy, M.,**
379 **Amirebrahimi, M., Weers, B.D., and McKinley, B.** (2018). The Sorghum bicolor reference
380 genome: improved assembly, gene annotations, a transcriptome atlas, and signatures of
381 genome organization. *The Plant Journal* **93**, 338-354.
- 382 **Ou, S., and Jiang, N.** (2017). LTR_retriever: A Highly Accurate And Sensitive Program For Identification Of
383 LTR Retrotransposons. *bioRxiv*.
- 384 **Paterson, A.H., Bowers, J.E., Bruggmann, R., Dubchak, I., Grimwood, J., Gundlach, H., Haberer, G.,**
385 **Hellsten, U., Mitros, T., and Poliakov, A.** (2009). The Sorghum bicolor genome and the
386 diversification of grasses. *Nature* **457**, 551-556.
- 387 **Patro, R., Duggal, G., Love, M.I., Irizarry, R.A., and Kingsford, C.** (2017). Salmon provides fast and bias-
388 aware quantification of transcript expression. *Nature methods* **14**, 417.
- 389 **Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V., and Zdobnov, E.M.** (2015). BUSCO:
390 assessing genome assembly and annotation completeness with single-copy orthologs.
391 *Bioinformatics* **31**, 3210-3212.
- 392 **Smit, A.F., Hubley, R., and Green, P.** (1996). RepeatMasker Open-3.0.
- 393 **Soneson, C., Love, M.I., and Robinson, M.D.** (2015). Differential analyses for RNA-seq: transcript-level
394 estimates improve gene-level inferences. *F1000Research* **4**.
- 395 **Team, R.C.** (2013). R: A language and environment for statistical computing.
- 396 **VanBuren, R., Wai, C.M., Pardo, J., Giarola, V., Ambrosini, S., Song, X., and Bartels, D.** (2018a).
397 Desiccation Tolerance Evolved through Gene Duplication and Network Rewiring in *Lindernia*.
398 *Plant Cell*.
- 399 **VanBuren, R., Wai, J., Zhang, Q., Song, X., Edger, P.P., Bryant, D., Michael, T.P., Mockler, T.C., and**
400 **Bartels, D.** (2017). Seed desiccation mechanisms co-opted for vegetative desiccation in the
401 resurrection grass *Oropetium thomeaeum*. *Plant, cell & environment*.
- 402 **VanBuren, R., Wai, C.M., Ou, S., Pardo, J., Bryant, D., Jiang, N., Mockler, T.C., Edger, P., and Michael,**
403 **T.P.** (2018b). Extreme haplotype variation in the desiccation-tolerant clubmoss *Selaginella*
404 *lepidophylla*. *Nature communications* **9**, 13.
- 405 **VanBuren, R., Bryant, D., Edger, P.P., Tang, H., Burgess, D., Challabathula, D., Spittle, K., Hall, R., Gu, J.,**
406 **and Lyons, E.** (2015). Single-molecule sequencing of the desiccation-tolerant grass *Oropetium*
407 *thomaeum*. *Nature* **527**, 508-511.

- 408 **Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C.A., Zeng, Q.,**
409 **Wortman, J., and Young, S.K.** (2014). Pilon: an integrated tool for comprehensive microbial
410 variant detection and genome assembly improvement. *PLoS one* **9**, e112963.
- 411 **Wang, Y., Tang, H., DeBarry, J.D., Tan, X., Li, J., Wang, X., Lee, T.-h., Jin, H., Marler, B., and Guo, H.**
412 (2012). MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and
413 collinearity. *Nucleic acids research* **40**, e49-e49.
- 414 **Wick, R.R., Schultz, M.B., Zobel, J., and Holt, K.E.** (2015). Bandage: interactive visualization of de novo
415 genome assemblies. *Bioinformatics* **31**, 3350-3352.
- 416 **Wicker, T., Buchmann, J.P., and Keller, B.** (2010). Patching gaps in plant genomes results in gene
417 movement and erosion of colinearity. *Genome research*, gr. 107284.107110.
- 418 **Xiao, L., Yang, G., Zhang, L., Yang, X., Zhao, S., Ji, Z., Zhou, Q., Hu, M., Wang, Y., and Chen, M.** (2015).
419 The resurrection genome of *Boea hygrometrica*: A blueprint for survival of dehydration.
420 *Proceedings of the National Academy of Sciences* **112**, 5833-5837.
- 421 **Xu, Z., and Wang, H.** (2007). LTR_FINDER: an efficient tool for the prediction of full-length LTR
422 retrotransposons. *Nucleic Acids Research* **35**, W265-W268.
- 423 **Xu, Z., Xin, T., Bartels, D., Li, Y., Gu, W., Yao, H., Liu, S., Yu, H., Pu, X., and Zhou, J.** (2018). Genome
424 analysis of the ancient tracheophyte *Selaginella tamariscina* reveals evolutionary features
425 relevant to the acquisition of desiccation tolerance. *Molecular plant*.
- 426 **Zhang, Q., and Bartels, D.** (2018). Molecular responses to dehydration and desiccation in desiccation-
427 tolerant angiosperm plants. *Journal of experimental botany* **69**, 3211-3222.

428

429

430 **Figure legends:**

431

432 **Figure 1. Hi-C based contig anchoring.** Post-clustering heat map showing density of Hi-C interactions
433 between contigs from the Juicer and 3d-DNA pipeline. The ten Oropetium chromosomes are highlighted
434 by blue squares.

435

436 **Figure 2. Characterization of the updated V2 Oropetium annotation.** (a) Tandem gene array size
437 comparison of the V1 and V2 annotation. Tandem genes identified in V1 are shown in blue and tandem
438 genes newly annotated in V2 are shown in gold. (b) Comparison of expression patterns from the V1 and
439 V2 annotation. The total number of genes with detectable expression and differential expression (DE) in
440 the Oropetium desiccation/rehydration timecourse are plotted.

441

442 **Figure 3. Landscape of the Oropetium genome.** *Gypsy* and *Copia* long terminal repeat retrotransposons
443 (LTR-RT) and CDS density are plotted for the ten Oropetium chromosomes. Features are plotted in
444 sliding windows of 50kb with 25kb step size. The location of centromere specific tandem arrays is
445 highlighted by red bars. The heatmaps below each landscape show relative density with red indicating
446 high density and blue indicating low density for each feature.

447

448 **Figure 4. Comparative genomics between Oropetium and Sorghum.** (a) Macrosyntenic dotplot of the
449 Oropetium and Sorghum chromosomes based on 18,889 gene pairs. Each black dot represents a synteny
450 region between the two genomes. (b) Microsynteny of a typical genic region of Sorghum and Oropetium
451 (top) and the pericentromeric region of Chromosome 6 of Oropetium and Sorghum (bottom). LTR-RTs
452 are shown in yellow and genes are shown in blue. Syntenic orthologs are connected by gray lines. The
453 centromeric repeat array in Oropetium is shown in red.

454

455

456 **Table 1:** Comparison of the Oropetium V1 and V2 assembly and annotation statistics

Statistics	V1	V2
# of contigs	625	436
Contig N50	2.38 Mb	2.02 Mb
Scaffold N50	NA	20.5 Mb
Total assembly size	243 Mb	236 Mb
Gene models	28,446	28,835
BUSCO	72.1%	98.9%

457

458

459 **Table 2:** Centromeric repeat array composition

Chromosome	Start Cent. Array (bp)	End Cent. Array (bp)	Number of Cent. Repeats	Cent. Size (bp)
Chr_1	18,899,082	19,114,162	154	215,080
Chr_2	18,277,215	18,463,229	786	186,014
Chr_3	18,882,303	18,993,598	308	111,295
Chr_4	11,739,636	13,338,554	176	1,598,918
Chr_5	10,361,368	10,828,355	800	466,987
Chr_6	3,649,010	3,746,417	513	97,407
Chr_7	12,434,273	12,559,564	272	125,291
Chr_8	8,288,262	9,010,114	306	721,852
Chr_9	6,142,739	7,433,209	1,044	1,290,470
Chr_10	3,147,692	3,209,432	155	61,740
Unanchored			4,258	982,774

460