

ORIGINAL RESEARCH

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

 Robert VanBuren^{1,2} | Ching Man Wai¹ | Jens Keilwagen³ | Jeremy Pardo⁴
¹Department of Horticulture, Michigan State University, East Lansing, Michigan

²Plant Resilience Institute, Michigan State University, East Lansing, Michigan

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

⁴Department of Plant Biology, Michigan State University, East Lansing, Michigan

Correspondence

 Robert VanBuren, Department of Horticulture, Michigan State University, East Lansing, MI.
 Email: bobvanburen@gmail.com

Funding information

National Science Foundation, Grant/Award Number: MCB-1817347

Abstract

Oropetium thomaeum is an emerging model for desiccation tolerance and genome size evolution in grasses. A draft genome of *Oropetium* was recently sequenced, but the lack of a chromosome-scale assembly has hindered comparative analyses and downstream functional genomics. Here, we reassembled *Oropetium*, and anchored the genome into 10 chromosomes using high-throughput chromatin conformation capture (Hi-C) based chromatin interactions. A combination of high-resolution RNAseq data and homology-based gene prediction identified thousands of new, conserved gene models that were absent from the V1 assembly. This includes thousands of new genes with high expression across a desiccation timecourse. Comparison between the *Sorghum* and *Oropetium* genomes revealed a surprising degree of chromosome-level collinearity, and several chromosome pairs have near perfect synteny. Other chromosomes are collinear in the gene rich chromosome arms but have experienced pericentric translocations. Together, these resources will be useful for the grass-comparative genomic community and further establish *Oropetium* as a model resurrection plant.

KEYWORDS

chromosome-scale, comparative genomics, desiccation tolerance, grasses, Hi-C

1 | INTRODUCTION

Desiccation tolerance evolved as an adaptation to extreme and prolonged drying, and resurrection plants are among the most resilient plants on the planet. The molecular basis of desiccation tolerance is still largely unknown, but a number of models have emerged to dissect the genetic control of this trait (Hoekstra, Golovina, & Buitink, 2001; Zhang & Bartels, 2018). The genomes of several model resurrection plants have been sequenced including *Boea hygrometrica* (Xiao et al., 2015), *Oropetium thomaeum* (VanBuren et al., 2015), *Xerophyta viscosa* (Costa et al., 2017), *Lindernia brevidens* (VanBuren, Wai, Pardo, et al., 2018), *Selaginella lepidophylla* (VanBuren, Wai, Ou, et al., 2018), and *Selaginella tamariscina* (Xu et al., 2018). To date, no

chromosome scale assemblies are available for these species, limiting large-scale quantitative genetics and comparative genomics-based approaches. Many resurrection plants are polyploid or have prohibitively large genomes including those in the genera *Boea*, *Xerophyta*, *Eragrostis*, *Sporobolus*, and *Craterostigma*. This complexity complicates genome assembly and gene redundancy in the polyploid species hinders downstream functional genomics work.

Oropetium thomaeum (hereon referred to as *Oropetium*), is a diploid resurrection plant with the smallest genome among the grasses (245 Mb) (Bartels & Mattar, 2002). *Oropetium* plants are similar in size to *Arabidopsis*, but significantly smaller than the model grasses *Setaria italica* (Li & Brutnell, 2011) and *Brachypodium distachyon* (Brkljacic et al., 2011), with a short generation time of ~4 months. *Oropetium* is in the Chloridoideae subfamily of grasses and is closely related to the orphan cereal crops *tef* (*Eragrostis tef*)

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

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2018 The Authors. *Plant Direct* published by American Society of Plant Biologists, Society for Experimental Biology and John Wiley & Sons Ltd.



and finger millet (*Eleusine coracana*). Desiccation tolerance evolved independently several times within Chloridoideae (Gaff, 1977, 1987; Gaff & Latz, 1978) making it a useful system for studying convergent evolution. Together, these traits make Oropetium an attractive model for exploring the origin and molecular basis of desiccation tolerance. Oropetium was one of the first plants to be sequenced using the long reads of PacBio technology, and the assembly quality was comparable to early Sanger sequencing-based plant genomes such as rice and *Arabidopsis* (VanBuren et al., 2015). Despite the high contiguity of Oropetium V1, the assembly has 625 contigs and the BioNano based genome map was unable to produce chromosome-scale scaffolds. Furthermore, the V1 annotation was based on limited transcript evidence, and a high proportion of conserved plant genes were missing (VanBuren et al., 2015). Here, we reassembled the Oropetium genome using a more refined algorithm and generated a chromosome-scale assembly using Hi-C based chromatin interactions. The annotation quality was improved using high-resolution RNAseq data and protein homology, facilitating detailed comparative genomics with other grasses.

2 | METHODS

2.1 | Genome reassembly

The raw PacBio reads from the Oropetium V1 release (VanBuren et al., 2015) were reassembled with improved algorithms to better resolve highly complex and repetitive regions. PacBio data were error corrected and assembled using Canu (V1.4) (Koren et al., 2017) with the following modifications: `minReadLength = 1,500`, `GenomeSize = 245 Mb`, and `minOverlapLength = 1,000`. Other parameters were left as default. The resulting assembly graph was visualized in Bandage (Wick, Schultz, Zobel, & Holt, 2015). The assembly graph was free of heterozygosity related bubbles, but many nodes (contigs) were interconnected by a high copy number retrotransposon. The Canu based contigs (assembly V1.2) were first polished using Quiver (Chin et al., 2013) with the raw PacBio data and default parameters. Contigs were further polished with Pilon (V1.22) (Walker et al., 2014) using $\sim 120\times$ coverage of paired-end 150 bp Illumina data. Quality-trimmed Illumina reads were aligned to the draft contigs using bowtie2 (V2.3.0) (Langmead & Salzberg, 2012) with default parameters. The overall alignment rate was 95.5%, which was slightly higher than alignment against the HGAP V1 assembly (94.5%). The following parameters for Pilon were modified: `-flank 7`, `-K 49`, and `-mindepth 25`. Other parameters were left as default. Pilon was run four times with an updated reference and realignment of Illumina data after each iteration. Indel corrections plateaued after the third iteration, suggesting polishing removed most residual assembly errors.

2.2 | Hi-C library construction analysis and genome anchoring

Oropetium descended from the original plants used for PacBio sequencing was collected for Hi-C library construction and RNAseq.

Oropetium is highly selfing with low heterozygosity, and we expect minimal differences to be introduced in the new version. Oropetium seeds are available upon request. Oropetium plants were maintained under day/night temperatures of 26 and 22°C, respectively, with a light intensity of 200 E m⁻² s⁻¹ and 16/8 hr photoperiod. Young leaf tissue was used for Hi-C library construction with the Proximo™ Hi-C Plant kit (Phase Genomics) following the manufacturer's protocol. Briefly, 0.2 g of fresh, young leaf tissue was finely chopped and the chromatin was immediately crosslinked. The chromatin was fragmented and proximity ligated, followed by library construction. The final library was size selected for 300–600 bp and sequenced on the Illumina HiSeq 4000 under paired-end 150 bp mode. Adapters were trimmed and low-quality sequences were removed using Trimmomatic (V0.36) (Bolger, Lohse, & Usadel, 2014). Read pairs were aligned to the Oropetium contigs using bwa (V0.7.16) (Li, 2013) with strict parameters (`-n 0`) to prevent mismatches and non-specific alignments in duplicated and repetitive regions. SAM files from bwa were used as input in the Juicer pipeline, and PCR duplicates with the same genome coordinates were filtered prior to constructing the interaction based distance matrix. In total, 101 filtered read pairs were used as input for the Juicer and 3d-DNA Hi-C analysis and scaffolding pipelines (Dudchenko et al., 2017; Durand et al., 2016). Contig ordering, orientation, and chimera splitting were done using the 3d-DNA pipeline (Dudchenko et al., 2017) under default parameters. Contig misassemblies and scaffold misjoins were manually detected and corrected based on interaction densities from visualization in Juicebox. In total, six chimeric contigs were identified and split at the junction with closest interaction data. The manually validated assembly was used as input to build the 10 scaffolds (chromosomes) using the `finalize-output.sh` script from 3d-DNA. Chromosomes and unanchored contigs were renamed by size, producing the V2 assembly.

2.3 | Genome annotation

The Oropetium V2 assembly was reannotated using the homology-based gene prediction program Gene Model Mapper (GeMoMa: V 1.5.2) (Keilwagen et al., 2016, 2018). GeMoMa uses protein homology and RNAseq evidence to predict gene models. Genome assemblies and gene annotation for the following 11 species were downloaded from Phytozome (V12) and used as homology-based evidence: *Arabidopsis thaliana*, *Brachypodium distachyon*, *Glycine max*, *Oryza sativa*, *Panicum hallii*, *Populus trichocarpa*, *Prunus persica*, *Setaria italica*, *Solanum lycopersicum*, *Sorghum bicolor*, and *Theobroma cacao*. Translated coding exons and proteins from the reference gene annotations and genome assemblies were extracted using the module Extractor function of GeMoMa (module Extractor: `Ambiguity = AMBIGUOUS`, `r = true`). RNAseq data from Oropetium desiccation and rehydration timecourses (VanBuren et al., 2017) were aligned to the V2 Oropetium genome using HISAT2 (Kim, Langmead, & Salzberg, 2015) with default parameters. The resulting BAM files were used to extract intron and exon boundaries using the module ERE (module ERE: `s = FR_FIRST_STRAND`, `c = true`). Translated coding exons from other species were aligned to the Oropetium

genome using tblastn and transcripts were predicted based on each reference species independently using the extracted introns and coverage (module GeMoMa). Finally, the predictions based on the 11 reference species were combined to obtain a final prediction using the module GAF. Gene models containing transposases were filtered, resulting in a final annotation of 28,835 gene models. The annotation completeness was assessed using the plant specific Benchmarking Universal Single-Copy Ortholog (BUSCO) dataset (version 3.0.2, embryophyta_odb9) (Simão, Waterhouse, Ioannidis, Kriventseva, & Zdobnov, 2015). The following report was obtained from BUSCO: 98.9% overall, 95.4% single copy, 3.5% duplicated, 0.6% fragmented, and 0.5% missing. Gene model names from V1 were conserved where possible, and new gene models received new names.

2.4 | Expression analysis

Oropetium RNAseq data from desiccation and rehydration time-courses were reanalyzed using the updated gene model annotations (VanBuren et al., 2017). Four timepoints during dehydration (days 7, 14, 21, and 30), two during rehydration (24 and 48 hr), and one well-watered sample were analyzed. Based on principle component analysis, replicate two of the “well-watered” and “D21” samples were excluded from the analysis. Each other timepoint had three replicates. Gene expression was quantified on a transcript level using salmon (v 0.9.1) in quasi-mapping mode (Patro, Duggal, Love, Irizarry, & Kingsford, 2017). Default parameters were used with the internal GC bias correction in salmon. The R package tximport (v 1.2.0) was used to map transcript level quantifications to gene level counts (R Core Team, 2013; Sonesson, Love, & Robinson, 2015). We conducted differential expression analysis with the remaining samples using the R package DESeq2 (v 1.14.1) set to default parameters [3,4].

2.5 | Identification of LTR-RTs

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

2.6 | Comparative genomics

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

3 | RESULTS AND DISCUSSION

The first version of the Oropetium genome (V1) was sequenced with high coverage PacBio data (~72x) followed by error correction and assembly using the hierarchical genome assembly process (HGAP) (VanBuren et al., 2015). We reassembled this PacBio data using the Canu assembler (Koren et al., 2017), which can more accurately assemble and phase complex repetitive regions. The resulting Canu based assembly (hereon referred to as V1.2) had fewer contigs than the V1 HGAP assembly, but had otherwise similar assembly metrics (Table 1). Draft contigs were polished using a two-step process to remove residual insertion/deletion (Indel) and single nucleotide errors. Contigs were first polished using the raw PacBio data with Quiver (Chin et al., 2013), followed by four rounds of reiterative polishing with Pilon (Walker et al., 2014) using high coverage Illumina paired-end data. The final V1.2 assembly contains 436 contigs with an N50 of 2.0 Mb and total assembly size of 236 Mb. This is six megabases smaller than the V1 assembly, with slightly lower contiguity. More intact LTR-RTs and centromere specific repeat arrays were identified in Oropetium V1.2 compared to V1, suggesting the Canu assembler resolved these repetitive elements more accurately. Thus, V1.2 was used for pseudomolecule construction.

The Oropetium V1.2 contigs were ordered and oriented into chromosome-scale pseudomolecules using Hi-C. Hi-C leverages long-range interactions across distal regions of chromosomes to order and orient contigs. This approach is similar to genetic map-based anchoring, but with higher resolution. Illumina data generated from the Hi-C library were mapped to the V1.2 Oropetium genome using bwa (Li, 2013) and the proximity-based clustering matrix was generated using the Juicer and 3d-DNA pipelines (Dudchenko et al.,

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%

2017; Durand et al., 2016). After filtering and manual curation, 10 high-confidence clusters were identified (Figure 1). These 10 clusters correspond to the haploid chromosome number of *Oropetium*. Regions with low density interactions highlight the centromeric and pericentromeric regions, and regions with higher than expected interactions represent topologically associated domains. After splitting six chimeric PacBio contigs, 239 contigs were anchored and oriented into 10 chromosomes spanning 226.5 Mb or 95.8% of the total assembled genome (Table 1). Chromosomes range in size from 11.0 to 34.7 Mb with an average size of 22.6 Mb. Most of the unanchored contigs are small (average size 42 kb), or are entirely composed of rRNA, centromeric repeat arrays, or centromere specific LTR-RTs. Telomeres were identified at both ends of Chromosomes 1, 2, 3, 4, 5, 7, and 9 and on one end of Chromosomes 6, 8, and 10. Three unanchored contigs contain the remaining telomeres. This supports the completeness and accuracy of the pseudomolecule construction.

The chromosome scale *Oropetium* genome (hereon referred to as V2) was reannotated using the homology-based gene prediction program GeMoMa (Keilwagen, Hartung, Paulini, Twardziok, & Grau, 2018; Keilwagen et al., 2016). Protein coding sequences from 11 angiosperm genomes and RNAseq data from *Oropetium* (VanBuren et al., 2017) were used as evidence. After filtering gene models derived from transposases, the final annotation consists of 28,835 high-confidence gene models. The annotation completeness was

assessed using the Benchmarking Universal Single-Copy Ortholog (BUSCO) embryophyta dataset. The V2 gene models have a BUSCO score of 98.9%, suggesting that the updated annotation is high-quality. In comparison, the *Oropetium* V1 annotation has a BUSCO score of 72%, and many conserved gene models were likely missing or misannotated. Nearly 40% (11,227) of the gene models in V2 are new and were unannotated in V1. In addition, 10,837 gene models from V1 were removed or substantially improved in the V2 annotation. These discarded gene models either had little support based on protein homology to other species and transcript evidence from *Oropetium*, or they were misannotated transposable elements. In total, 94.3% of the gene models (27,216) were anchored to the 10 chromosomes. Among the newly annotated gene models are 3,525 tandem gene duplicates (Figure 2a). Tandem duplicates span 3,062 arrays with 7,760 total genes. Of the arrays containing three or more genes, only 49 are new to V2, and the majority contain genes previously identified in V1. The boundaries of tandem duplicates are difficult to correctly annotate, resulting in fusions of two or more gene copies. The homology-based annotation used in V2 was able to parse the previously fused gene models.

The expression patterns of newly annotated genes were surveyed using high-resolution RNAseq expression data (VanBuren et al., 2017). This dataset consists of seven leaf samples collected during desiccation and rehydration timecourses. Timepoints include well-watered, 7, 14, 21, and 30 days desiccated as well as 24 and

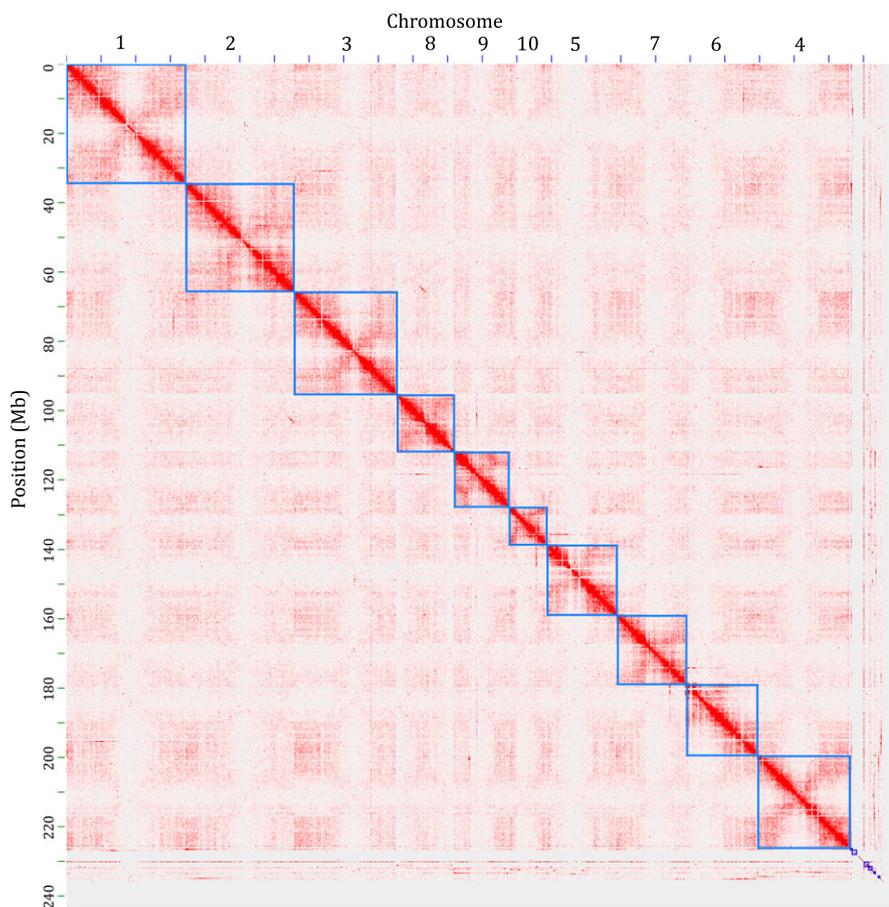


FIGURE 1 Hi-C based contig anchoring. Post-clustering heat map showing density of Hi-C interactions between contigs from the Juicer and 3d-DNA pipeline. The 10 *Oropetium* chromosomes are highlighted by blue squares

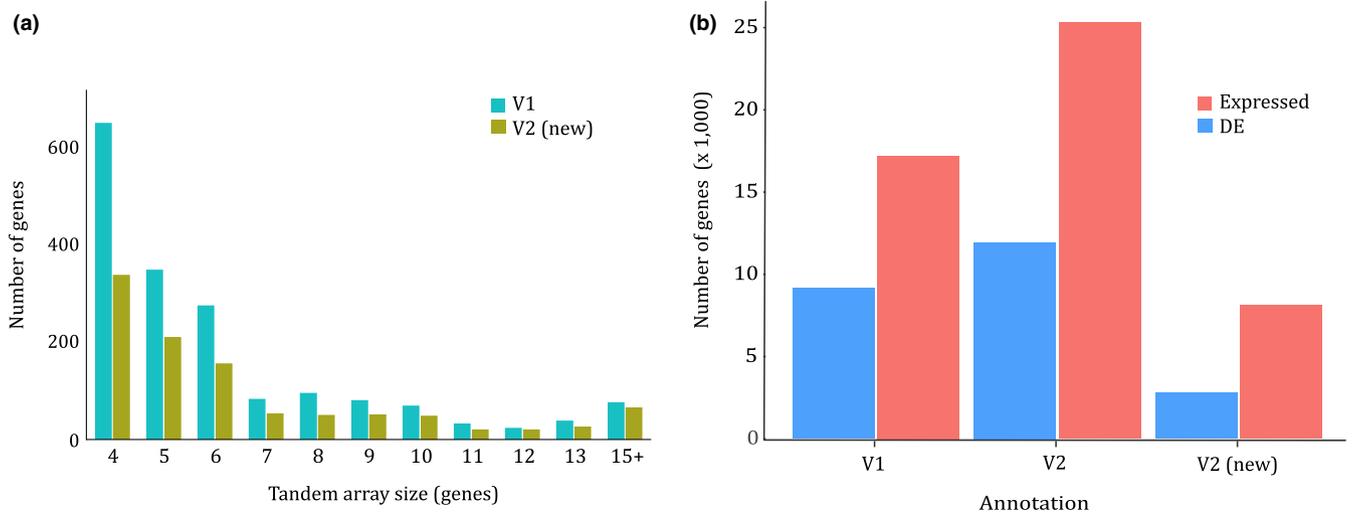


FIGURE 2 Characterization of the updated V2 *Oropetium* annotation. (a) Tandem gene array size comparison of the V1 and V2 annotation. Tandem genes identified in V1 are shown in blue and tandem genes newly annotated in V2 are shown in gold. (b) Comparison of expression patterns from the V1 and V2 annotation.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

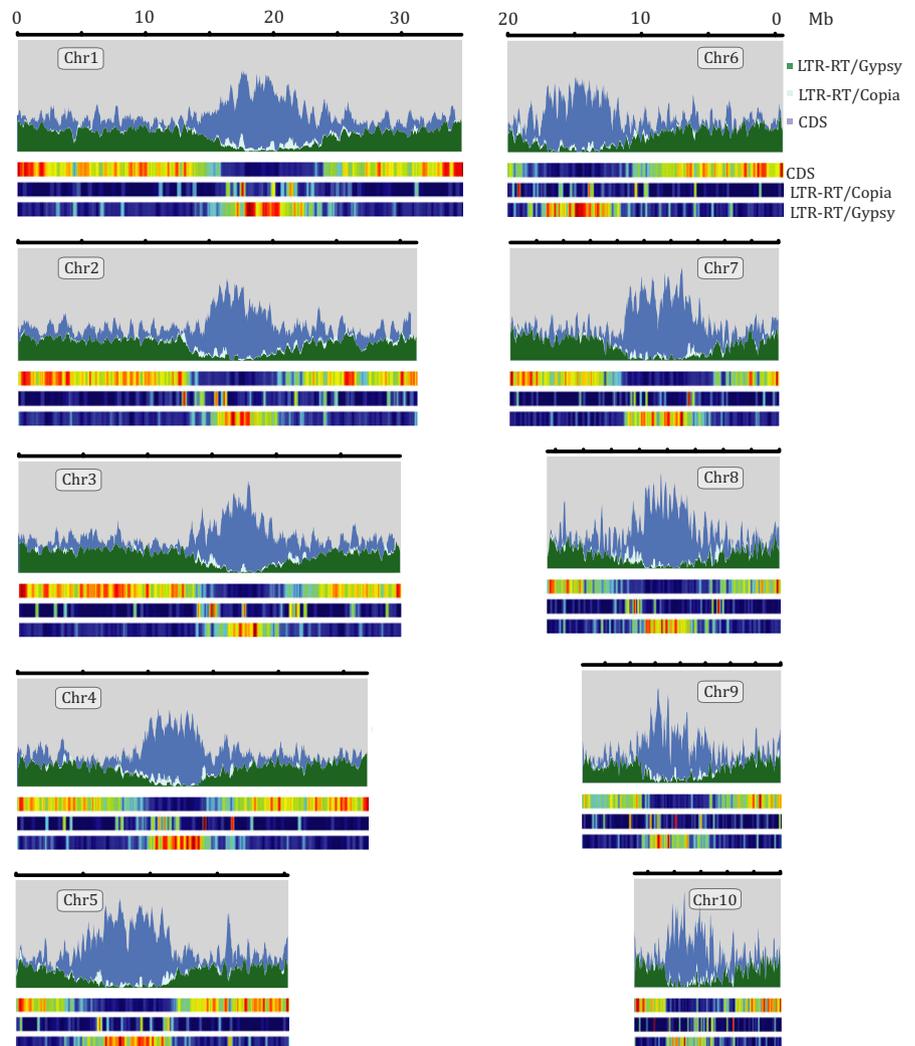


FIGURE 3 Landscape of the *Oropetium* genome. *Gypsy* and *Copia* long terminal repeat retrotransposons (LTR-RT) and CDS density are plotted for the 10 *Oropetium* chromosomes. Features are plotted in sliding windows of 50 kb with 25 kb step size. The location of centromere specific tandem arrays is highlighted by red bars. The heat maps below each landscape show relative density with red indicating high density and blue indicating low density for each feature

48 hr post rehydration. Differentially expressed genes were identified based on comparisons of well-watered leaves with each dehydration or rehydration timepoint. In addition, each timepoint was compared with the timepoint immediately following it in the time-course (i.e., day 7 dehydration vs. day 14). In total, 17,204 gene models from the V1 annotation had detectable expression (count > 0 in at least one sample) compared to 25,314 gene models in V2 (Figure 2b). Of the expressed genes, 9,149 V1 and 11,948 V2 gene models were classified as differentially expressed in at least one of

the comparisons. Most newly annotated genes (8,110) have detectable expression in at least one of the seven timepoints, and the majority are expressed in all timepoints. In total, 2,799 new V2 gene models were differentially expressed, suggesting the newly annotated genes have important and previously uncharacterized roles in desiccation tolerance.

We used the chromosome-scale assembly of *Oropetium* to survey patterns of genome organization and evolution related to maintaining a small genome size. The proportion of LTR-RTs in *Oropetium* V1 and V2 is similar, though V2 has more intact elements. LTR-RTs are the most abundant repetitive elements in *Oropetium* and collectively span 27% (62 Mb) of the genome. LTR-RTs are distributed non-randomly across the genome, and peaks of Gypsy LTR-RTs are observed in each of the 10 chromosomes (Figure 3). These peaks of Gypsy LTR-RTs correspond to the pericentromeric regions. The pericentromeric regions show reduced intra-chromosomal interactions in the Hi-C matrix and contain arrays of centromeric repeats. The *Oropetium* V2 genome contains 8,965 155 bp monomeric centromeric repeats; considerably more than the 4,315 identified in the V1 assembly. The centromeric array sizes vary from 61 kb in chromosome 10 to 1,598 kb in Chromosome 4 (Figure 3; Table 2). Array sizes are likely underestimated, as only 52% of centromeric arrays were anchored to chromosomes, and 23 unanchored contigs contain centromeric repeat arrays. Gene density is low in the pericentromeric regions, consistent with the rice, sorghum, maize, and *Brachypodium* genomes (Du et al., 2017; Initiative,

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

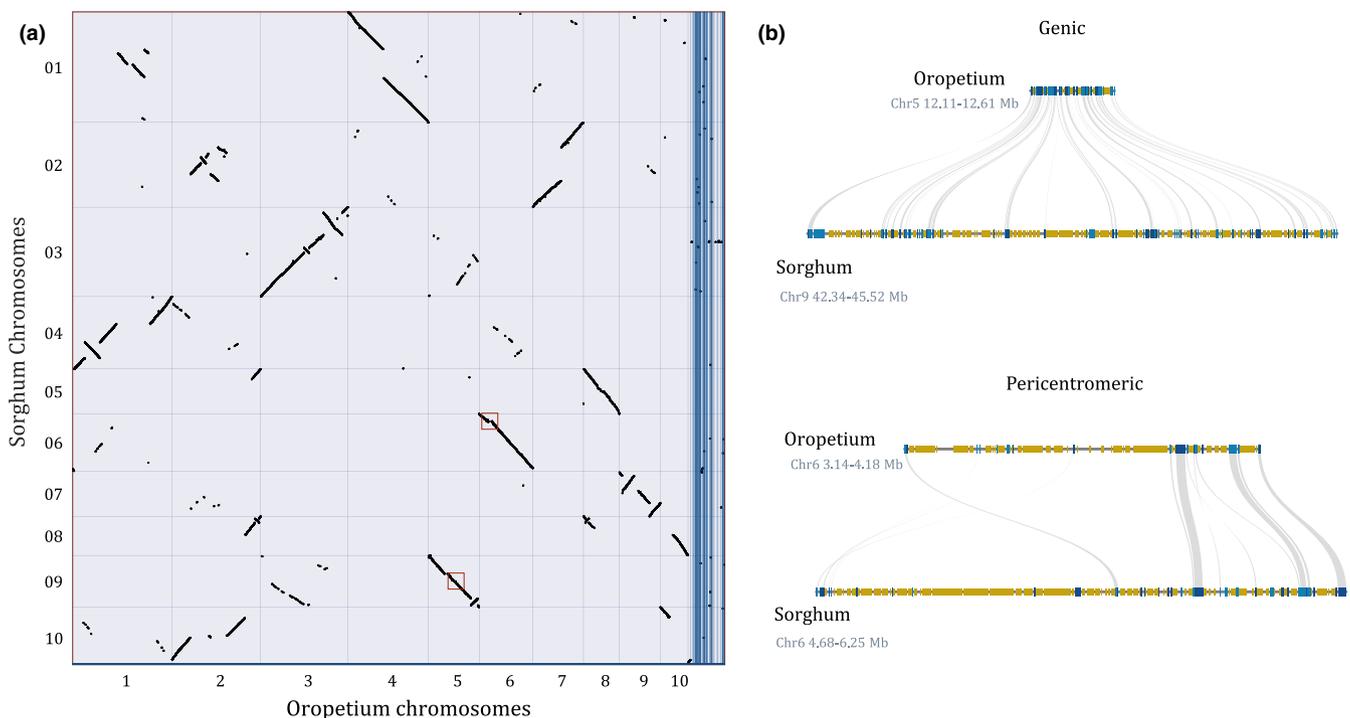


FIGURE 4 Comparative genomics between *Oropetium* and *Sorghum*. (a) Macrosynteny dotplot of the *Oropetium* and *Sorghum* chromosomes based on 18,889 gene pairs. Each black dot represents a syntenic region between the two genomes. (b) Microsynteny of a typical genic region of *Sorghum* and *Oropetium*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



2010; Jiao et al., 2017; Paterson et al., 2009). Collectively, pericentromeric regions span 67.5 Mb or 29% of the genome, a much smaller proportion than Sorghum (62%; 460 Mb) (Paterson et al., 2009), but higher than rice (15%; 63 Mb) (Goff et al., 2002). The majority of intact LTRs (86%; 628) have an insertion time of less than one million years ago, with a steep drop off of insertion time after 0.4 MYA. This suggests that LTRs are highly active in *Oropetium* but rapidly fragmented and purged to maintain its small genome size.

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

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

The read lengths of third generation sequencing technologies enable near gapless assemblies with high contiguity for virtually any plant genome. PacBio and Nanopore based genomes have a better representation of gene and regulatory sequences, but often lack the chromosome-scale scaffolding required for comparative genomics and quantitative genetics. The pseudomolecules in *Oropetium* V2 allowed us to more accurately identify syntenic orthologs in other grasses and make detailed comparisons of chromosome evolution. The V2 chromosome-scale assembly will serve as a reference for future population genomics work and positional cloning of

desiccation related genes. Together, this highlights the need to improve and scaffold existing high-quality reference genomes.

ACKNOWLEDGMENTS

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

AUTHOR CONTRIBUTIONS

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

AVAILABILITY OF SUPPORTING DATA

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

REFERENCES

- Bartels, D., & Mattar, M. (2002). *Oropetium thomaeum*: A resurrection grass with a diploid genome. *Maydica*, 47, 185–192.
- Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics*, 30, 2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>
- Brkljacic, J., Grotewold, E., Scholl, R., Mockler, T., Garvin, D. F., Vain, P., ... Budak, H. (2011). Brachypodium as a model for the grasses: Today and the future. *Plant Physiology*, 157, 3–13. <https://doi.org/10.1104/pp.111.179531>
- Cannarozzi, G., Plaza-Wüthrich, S., Esfeld, K., Larti, S., Wilson, Y. S., Girma, D., ... Farinelli, L. (2014). Genome and transcriptome sequencing identifies breeding targets in the orphan crop tef (*Eragrostis tef*). *BMC Genomics*, 15, 581. <https://doi.org/10.1186/1471-2164-15-581>
- Chin, C.-S., Alexander, D. H., Marks, P., Klammer, A. A., Drake, J., Heiner, C., ... Eichler, E. E. (2013). Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nature Methods*, 10, 563–569. <https://doi.org/10.1038/nmeth.2474>
- Costa, M., Artur, M., Maia, J., Jonkheer, E., Derks, M., Nijveen, H., ... Hesselink, T. (2017). A footprint of desiccation tolerance in the genome of *Xerophyta viscosa*. *Nature Plants*, 3, 17038. <https://doi.org/10.1038/nplants.2017.38>
- Cotton, J. L., Wysocki, W. P., Clark, L. G., Kelchner, S. A., Pires, J. C., Edger, P. P., ... Duvall, M. R. (2015). Resolving deep relationships of PACMAD grasses: A phylogenomic approach. *BMC Plant Biology*, 15, 178. <https://doi.org/10.1186/s12870-015-0563-9>
- Du, H., Yu, Y., Ma, Y., Gao, Q., Cao, Y., Chen, Z., ... Zhao, X. (2017). Sequencing and de novo assembly of a near complete indica rice genome. *Nature Communications*, 8, 15324. <https://doi.org/10.1038/ncomms15324>
- Dudchenko, O., Batra, S. S., Omer, A. D., Nyquist, S. K., Hoeger, M., Durand, N. C., ... Aiden, A. P. (2017). De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds. *Science*, 356, 92–95. <https://doi.org/10.1126/science.123327>

- Durand, N. C., Shamim, M. S., Machol, I., Rao, S. S., Huntley, M. H., Lander, E. S., & Aiden, E. L. (2016). Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Systems*, 3, 95–98. <https://doi.org/10.1016/j.cels.2016.07.002>
- Ellinghaus, D., Kurtz, S., & Willhoeft, U. (2008). LTRharvest, an efficient and flexible software for de novo detection of LTR retrotransposons. *BMC Bioinformatics*, 9, 18. <https://doi.org/10.1186/1471-2105-9-18>
- Gaff, D. (1977). Desiccation tolerant vascular plants of Southern Africa. *Oecologia*, 31, 95–109. <https://doi.org/10.1007/BF00348713>
- Gaff, D. (1987). Desiccation tolerant plants in South America. *Oecologia*, 74, 133–136. <https://doi.org/10.1007/BF00377357>
- Gaff, D., & Latz, P. (1978). The occurrence of resurrection plants in the Australian flora. *Australian Journal of Botany*, 26, 485–492. <https://doi.org/10.1071/BT9780485>
- Goff, S. A., Ricke, D., Lan, T.-H., Presting, G., Wang, R., Dunn, M., ... Varma, H. (2002). A draft sequence of the rice genome (*Oryza sativa* L. ssp. *japonica*). *Science*, 296, 92–100. <https://doi.org/10.1126/science.1068275>
- Hittalmani, S., Mahesh, H., Shirke, M. D., Biradar, H., Uday, G., Aruna, Y., ... Mohanrao, A. (2017). Genome and Transcriptome sequence of Finger millet (*Eleusine coracana* (L.) Gaertn.) provides insights into drought tolerance and nutraceutical properties. *BMC Genomics*, 18, 465. <https://doi.org/10.1186/s12864-017-3850-z>
- Hoekstra, F. A., Golovina, E. A., & Buitink, J. (2001). Mechanisms of plant desiccation tolerance. *Trends in Plant Science*, 6, 431–438. [https://doi.org/10.1016/S1360-1385\(01\)02052-0](https://doi.org/10.1016/S1360-1385(01)02052-0)
- Initiative, I. B. (2010). Genome sequencing and analysis of the model grass *Brachypodium distachyon*. *Nature*, 463, 763.
- Jiao, Y., Peluso, P., Shi, J., Liang, T., Stitzer, M. C., Wang, B., ... Chin, C.-S. (2017). Improved maize reference genome with single-molecule technologies. *Nature*, 546, 524–527.
- Keilwagen, J., Hartung, F., Paulini, M., Twardziok, S. O., & Grau, J. (2018). Combining RNA-seq data and homology-based gene prediction for plants, animals and fungi. *BMC Bioinformatics*, 19, 189. <https://doi.org/10.1186/s12859-018-2203-5>
- Keilwagen, J., Wenk, M., Erickson, J. L., Schattat, M. H., Grau, J., & Hartung, F. (2016). Using intron position conservation for homology-based gene prediction. *Nucleic Acids Research*, 44, e89. <https://doi.org/10.1093/nar/gkw092>
- Kim, D., Langmead, B., & Salzberg, S. L. (2015). HISAT: A fast spliced aligner with low memory requirements. *Nature Methods*, 12, 357. <https://doi.org/10.1038/nmeth.3317>
- Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., & Phillippy, A. M. (2017). Canu: Scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Research*, 27, 722–736. <https://doi.org/10.1101/gr.215087.116>
- Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods*, 9, 357–359. <https://doi.org/10.1038/nmeth.1923>
- Li, H. (2013). *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM*. arXiv preprint arXiv:1303.3997.
- Li, P., & Brutnell, T. P. (2011). *Setaria viridis* and *Setaria italica*, model genetic systems for the Panicoid grasses. *Journal of Experimental Botany*, 62, 3031–3037. <https://doi.org/10.1093/jxb/err096>
- McCormick, R. F., Truong, S. K., Sreedasyam, A., Jenkins, J., Shu, S., Sims, D., ... McKinley, B. (2018). The *Sorghum bicolor* reference genome: Improved assembly, gene annotations, a transcriptome atlas, and signatures of genome organization. *The Plant Journal*, 93, 338–354. <https://doi.org/10.1111/tpj.13781>
- Ou, S., & Jiang, N. (2017). LTR_retriever: A highly accurate and sensitive program for identification of LTR retrotransposons. *Plant Physiology*, 176, 1410–1422.
- Paterson, A. H., Bowers, J. E., Bruggmann, R., Dubchak, I., Grimwood, J., Gundlach, H., ... Poliakov, A. (2009). The *Sorghum bicolor* genome and the diversification of grasses. *Nature*, 457, 551–556. <https://doi.org/10.1038/nature07723>
- Patro, R., Duggal, G., Love, M. I., Irizarry, R. A., & Kingsford, C. (2017). Salmon provides fast and bias-aware quantification of transcript expression. *Nature Methods*, 14, 417. <https://doi.org/10.1038/nmeth.4197>
- R Core Team (2013). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., & Zdobnov, E. M. (2015). BUSCO: Assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*, 31, 3210–3212. <https://doi.org/10.1093/bioinformatics/btv351>
- Smit, A. F., Hubley, R., & Green, P. (1996). *RepeatMasker Open-3.0*. Retrieved from <http://www.repeatmasker.org/>
- Soneson, C., Love, M. I., & Robinson, M. D. (2015). Differential analyses for RNA-seq: Transcript-level estimates improve gene-level inferences. *F1000Research*, 4, 1521. <https://doi.org/10.12688/f1000research>
- VanBuren, R., Bryant, D., Edger, P. P., Tang, H., Burgess, D., Challabathula, D., ... Lyons, E. (2015). Single-molecule sequencing of the desiccation-tolerant grass *Oropetium thomaeum*. *Nature*, 527, 508–511. <https://doi.org/10.1038/nature15714>
- VanBuren, R., Wai, C. M., Ou, S., Pardo, J., Bryant, D., Jiang, N., ... Michael, T. P. (2018). Extreme haplotype variation in the desiccation-tolerant clubmoss *Selaginella lepidophylla*. *Nature Communications*, 9, 13. <https://doi.org/10.1038/s41467-017-02546-5>
- VanBuren, R., Wai, C. M., Pardo, J., Giarola, V., Ambrosini, S., Song, X., & Bartels, D. (2018). Desiccation tolerance evolved through gene duplication and network rewiring in *Lindernia*. *Plant Cell*, <https://doi.org/10.1105/tpc.18.00517>.
- VanBuren, R., Wai, J., Zhang, Q., Song, X., Edger, P. P., Bryant, D., ... Bartels, D. (2017). Seed desiccation mechanisms co-opted for vegetative desiccation in the resurrection grass *Oropetium thomaeum*. *Plant, Cell & Environment*, 40, 2292–2306. <https://doi.org/10.1111/pce.13027>
- Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., ... Young, S. K. (2014). Pilon: An integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS ONE*, 9, e112963. <https://doi.org/10.1371/journal.pone.0112963>
- Wang, Y., Tang, H., DeBarry, J. D., Tan, X., Li, J., Wang, X., ... Guo, H. (2012). MScanX: A toolkit for detection and evolutionary analysis of gene synteny and collinearity. *Nucleic Acids Research*, 40, e49. <https://doi.org/10.1093/nar/gkr1293>
- Wick, R. R., Schultz, M. B., Zobel, J., & Holt, K. E. (2015). Bandage: Interactive visualization of de novo genome assemblies. *Bioinformatics*, 31, 3350–3352. <https://doi.org/10.1093/bioinformatics/btv383>
- Wicker, T., Buchmann, J. P., & Keller, B. (2010). Patching gaps in plant genomes results in gene movement and erosion of colinearity. *Genome Research*, 20, 1229–1237. <https://doi.org/10.1101/gr.107284.110>
- Xiao, L., Yang, G., Zhang, L., Yang, X., Zhao, S., Ji, Z., ... Chen, M. (2015). The resurrection genome of *Boea hygrometrica*: A blueprint for survival of dehydration. *Proceedings of the National Academy of Sciences of the United States of America*, 112, 5833–5837. <https://doi.org/10.1073/pnas.1505811112>
- Xu, Z., & Wang, H. (2007). LTR_FINDER: An efficient tool for the prediction of full-length LTR retrotransposons. *Nucleic Acids Research*, 35, W265–W268. <https://doi.org/10.1093/nar/gkm286>
- Xu, Z., Xin, T., Bartels, D., Li, Y., Gu, W., Yao, H., ... Zhou, J. (2018). Genome analysis of the ancient tracheophyte *Selaginella tamariscina* reveals evolutionary features relevant to the acquisition of desiccation tolerance. *Molecular Plant*, 11, 983–994. <https://doi.org/10.1016/j.molp.2018.05.003>



Zhang, Q., & Bartels, D. (2018). Molecular responses to dehydration and desiccation in desiccation-tolerant angiosperm plants. *Journal of Experimental Botany*, 69, 3211–3222. <https://doi.org/10.1093/jxb/erx489>

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

How to cite this article: VanBuren R, Wai CM, Keilwagen J, Pardo J. A chromosome-scale assembly of the model desiccation tolerant grass *Oropetium thomaeum*. *Plant Direct*. 2018;2:1–9. <https://doi.org/10.1002/pld3.96>