

1-9-2020

## The *Gossypium longicalyx* genome as a resource for cotton breeding and evolution

Corrinne E. Grover  
*Iowa State University, corrinne@iastate.edu*

Mengqiao Pan  
*Nanjing Agricultural University*

Daojun Yuan  
*Huazhong Agricultural University*

Mark A. Arick II  
*Mississippi State University*

Guanjing Hu  
*Iowa State University, hugj2006@iastate.edu*

*See next page for additional authors*

Follow this and additional works at: [https://lib.dr.iastate.edu/eeob\\_ag\\_pubs](https://lib.dr.iastate.edu/eeob_ag_pubs)



Part of the [Agriculture Commons](#), [Ecology and Evolutionary Biology Commons](#), [Entomology Commons](#), and the [Genetics and Genomics Commons](#)

The complete bibliographic information for this item can be found at [https://lib.dr.iastate.edu/eeob\\_ag\\_pubs/389](https://lib.dr.iastate.edu/eeob_ag_pubs/389). For information on how to cite this item, please visit <http://lib.dr.iastate.edu/howtocite.html>.

---

This Article is brought to you for free and open access by the Ecology, Evolution and Organismal Biology at Iowa State University Digital Repository. It has been accepted for inclusion in Ecology, Evolution and Organismal Biology Publications by an authorized administrator of Iowa State University Digital Repository. For more information, please contact [digirep@iastate.edu](mailto:digirep@iastate.edu).

---

# The *Gossypium longicalyx* genome as a resource for cotton breeding and evolution

## Abstract

Cotton is an important crop that has made significant gains in production over the last century. Emerging pests such as the reniform nematode have threatened cotton production. The rare African diploid species *Gossypium longicalyx* is a wild species that has been used as an important source of reniform nematode immunity. While mapping and breeding efforts have made some strides in transferring this immunity to the cultivated polyploid species, the complexities of interploidal transfer combined with substantial linkage drag have inhibited progress in this area. Moreover, this species shares its most recent common ancestor with the cultivated A-genome diploid cottons, thereby providing insight into the evolution of long, spinnable fiber. Here we report a newly generated *de novo* genome assembly of *G. longicalyx*. This high-quality genome leveraged a combination of PacBio long-read technology, Hi-C chromatin conformation capture, and BioNano optical mapping to achieve a chromosome level assembly. The utility of the *G. longicalyx* genome for understanding reniform immunity and fiber evolution is discussed.

## Keywords

*Gossypium longicalyx*, nematode resistance, cotton fiber, genome sequence, PacBio

## Disciplines

Agriculture | Ecology and Evolutionary Biology | Entomology | Genetics and Genomics

## Comments

This preprint is made available through bioRxiv, doi: [10.1101/2020.01.08.898908](https://doi.org/10.1101/2020.01.08.898908).

## Authors

Corrinne E. Grover, Mengqiao Pan, Daojun Yuan, Mark A. Arick II, Guanqing Hu, Logan Brase, David M. Stelly, Zefu Lu, Robert J. Schmitz, Daniel G. Peterson, Jonathan F. Wendel, and Joshua A. Udall

1           **The *Gossypium longicalyx* genome as a resource for cotton breeding and evolution**

2           Corrinne E. Grover<sup>1</sup>, Mengqiao Pan<sup>2</sup>, Daojun Yuan<sup>3</sup>, Mark A. Arick II<sup>4</sup>, Guanqing Hu<sup>1</sup>, Logan  
3           Brase<sup>5</sup>, David M. Stelly<sup>6</sup>, Zefu Lu<sup>7</sup>, Robert J. Schmitz<sup>7</sup>, Daniel G. Peterson<sup>4</sup>, Jonathan F.  
4           Wendel<sup>1</sup>, and Joshua A. Udall<sup>8\*</sup>

5  
6           <sup>1</sup> Ecology, Evolution, and Organismal Biology Dept., Iowa State University, Ames, IA, 50010

7           <sup>2</sup> State Key Laboratory of Crop Genetics and Germplasm Enhancement, Cotton Hybrid R & D  
8           Engineering Center, Nanjing Agricultural University, Nanjing, 210095, China

9           <sup>3</sup> College of Plant Science and Technology, Huazhong Agricultural University, Wuhan, Hubei,  
10           430070, China

11           <sup>4</sup> Institute for Genomics, Biocomputing & Biotechnology, Mississippi State University, United  
12           States

13           <sup>5</sup> Division of Biology and Biomedical Sciences, Washington University in St. Louis, St. Louis,  
14           MO 63110

15           <sup>6</sup> Department of Soil and Crop Sciences, Texas A&M University, College Station, 77843, USA

16           <sup>7</sup> Department of Genetics, University of Georgia, Athens, GA 30602

17           <sup>8</sup> USDA/Agricultural Research Service, Crop Germplasm Research Unit, College Station, TX  
18           77845

19  
20           ORCID (email):

21           CEG: 0000-0003-3878-5459 ([corrinne@iastate.edu](mailto:corrinne@iastate.edu))

22           MP: ([mengqiaopan@live.com](mailto:mengqiaopan@live.com))

23           DY: 0000-0001-6007-5571([robert@mail.hzau.edu.cn](mailto:robert@mail.hzau.edu.cn))

24           MAA: 0000-0002-7207-3052 ([maa146@IGBB.MsState.Edu](mailto:maa146@IGBB.MsState.Edu))

25           GH: 0000-0001-8552-7394 ([hugj2006@iastate.edu](mailto:hugj2006@iastate.edu))

26           LB: 0000-0002-7175-3208 ([braselogan@gmail.com](mailto:braselogan@gmail.com))

27           DS: 0000-0002-3468-4119 ([stelly@tamu.edu](mailto:stelly@tamu.edu))

28           ZL: ([zefulu@uga.edu](mailto:zefulu@uga.edu))

29           RJS: 0000-0001-7538-6663 ([schmitz@uga.edu](mailto:schmitz@uga.edu))

30           DGP 0000-0002-0274-5968 ([peterperson@IGBB.MsState.Edu](mailto:peterperson@IGBB.MsState.Edu))

31           JFW 0000-0003-2258-5081 ([jfw@iastate.edu](mailto:jfw@iastate.edu))

32           JAU 0000-0003-0978-4764 ([Joshua.Udall@usda.gov](mailto:Joshua.Udall@usda.gov))

33  
34  
35           \*corresponding author: [Joshua.Udall@usda.gov](mailto:Joshua.Udall@usda.gov)

36  
37  
38           **Keywords:**

39           *Gossypium longicalyx*, nematode resistance, cotton fiber, genome sequence, PacBio

## 40 **Abstract**

41 Cotton is an important crop that has made significant gains in production over the last century.  
42 Emerging pests such as the reniform nematode have threatened cotton production. The rare  
43 African diploid species *Gossypium longicalyx* is a wild species that has been used as an  
44 important source of reniform nematode immunity. While mapping and breeding efforts have  
45 made some strides in transferring this immunity to the cultivated polyploid species, the  
46 complexities of interploidal transfer combined with substantial linkage drag have inhibited  
47 progress in this area. Moreover, this species shares its most recent common ancestor with the  
48 cultivated A-genome diploid cottons, thereby providing insight into the evolution of long,  
49 spinnable fiber. Here we report a newly generated *de novo* genome assembly of *G. longicalyx*.  
50 This high-quality genome leveraged a combination of PacBio long-read technology, Hi-C  
51 chromatin conformation capture, and BioNano optical mapping to achieve a chromosome level  
52 assembly. The utility of the *G. longicalyx* genome for understanding reniform immunity and  
53 fiber evolution is discussed.

54

55

## 56 **Introduction**

57

58 Cotton (genus *Gossypium*) is an important crop which provides the largest natural source of  
59 fiber. Colloquially, the term cotton refers to one of four domesticated species, primarily the  
60 tetraploid *G. hirsutum*, which is responsible for over 98% of cotton production worldwide  
61 (Kranthi 2018). *Gossypium* contains over 50 additional wild species related to the domesticated  
62 cottons that serve as potential sources of disease and pest resistance. Among these, *Gossypium*  
63 *longicalyx* J.B. Hutch. & B.J.S. Lee is the only representative of the diploid “F-genome”  
64 (Wendel and Grover 2015) and the only species with immunity to reniform nematode infection  
65 (Yik and Birchfield 1984). Discovered only 60 years ago (Hutchinson and B. J. S. Lee 1958), it  
66 is both cytogenetically differentiated from members of the other genome groups (Phillips 1966)  
67 and morphologically isolated (Fryxell 1971, 1992). Importantly, *G. longicalyx* is sister to the A-  
68 genome cottons (Wendel and Albert 1992; Wendel and Grover 2015; Chen *et al.* 2016), i.e., *G.*  
69 *arboreum* and *G. herbaceum*, the only diploids with long, spinnable fiber.

70

71 Interest in the genome of *G. longicalyx* is two-fold. First, broad-scale screening of the cotton  
72 germplasm collection indicates that domesticated cotton lacks appreciable natural resistance to  
73 reniform nematode (Birchfield *et al.* 1963; Yik and Birchfield 1984), and while several other  
74 species exhibit degrees of resistance, only *G. longicalyx* exhibits immunity to infection (Yik and  
75 Birchfield 1984). This is significant as reniform nematode has emerged as a major source of  
76 cotton crop damage, reducing cotton production by over 205 million bales per year (Lawrence *et al.*  
77 *al.* 2015) and accounting for ~11% of the loss attributable to pests (Khanal *et al.* 2018). Current  
78 reniform resistant lines are derived from complex breeding schemes which are required to  
79 introgress reniform immunity from the diploid *G. longicalyx* into polyploid *G. hirsutum* (Bell  
80 and Robinson 2004; Dighe *et al.* 2009; Khanal *et al.* 2018); however, undesirable traits have

81 accompanied this introgression (Nichols *et al.* 2010) extreme stunting of seedlings and plants  
82 exposed to dense nematode populations, prohibiting commercial deployment (Zheng *et al.* 2016).

83  
84 The genome of *G. longicalyx* is also valuable because it is phylogenetically sister to the only  
85 diploid clade with spinnable fiber (Wendel and Albert, 1992; Wendel and Grover, 2015; Chen *et*  
86 *al.*, 2016), the A-genome species, which contributed the maternal ancestor to polyploid cotton.  
87 Consequently, there has been interest in this species as the ancestor to spinnable fiber (Hovav *et*  
88 *al.* 2008; Paterson *et al.* 2012), although progress has been limited due to lack of genomic  
89 resources in *G. longicalyx*. Comparisons between the *G. longicalyx* genome and other cotton  
90 genomes, including the domesticated diploids (Du *et al.* 2018), may provide clues into the  
91 evolutionary origin of “long” fiber.

92  
93 Here we describe a high-quality, *de novo* genome sequence for *G. longicalyx*, a valuable  
94 resource for understanding nematode immunity in cotton and possibly other species. This  
95 genome also provides a foundation to understand the evolutionary origin of spinnable fiber in  
96 *Gossypium*.

97  
98

## 99 **Methods & Materials**

100

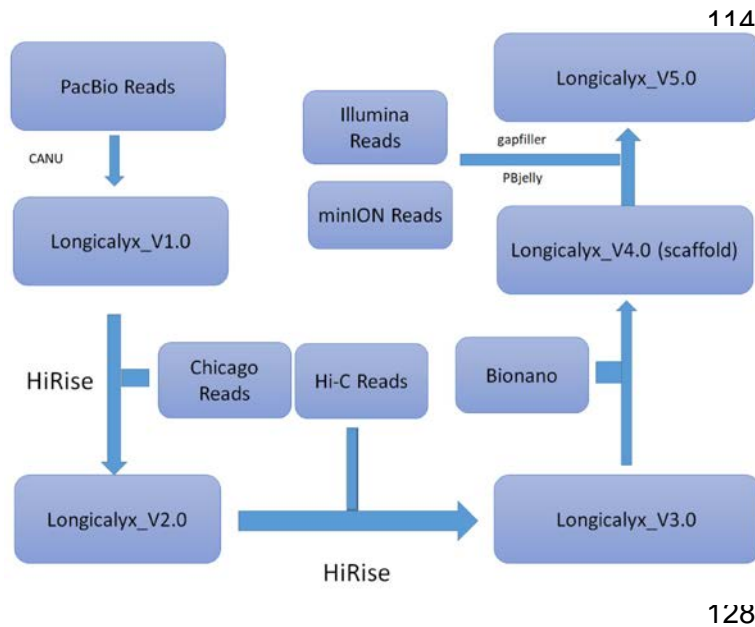
### 101 **Plant material and sequencing methods**

102 Leaf tissue of mature *G. longicalyx* (F1-1) was collected from a Brigham Young University  
103 (BYU) greenhouse. DNA was extracted using CTAB techniques (Kidwell and Osborn 1992),  
104 and the amount recovered was measured via Qubit Fluorometer (ThermoFisher, Inc.). The  
105 sequencing library was constructed by the BYU DNA Sequencing Center (DNASC) using only  
106 fragments >18 kb, which were size selected on the BluePippen (Sage Science, LLC) and verified  
107 in size using a Fragment Analyzer (Advanced Analytical Technologies, Inc). Twenty-six PacBio  
108 cells were sequenced from a single library on the Pacific Biosciences Sequel system. Resulting  
109 reads were assembled using Canu V1.6 using default parameters (Koren *et al.* 2017) to create a  
110 sequence assembly called Longicalyx\_V1.0, composed of 229 large contigs (Figure 1).

111

112

113



**Figure 1.**

Chicago Highrise reads (Dovetail Genomics) provide DNA-DNA proximity information used to improve the Canu sequence assembly (Longicalyx\_V2.0; statistics not calculated), as previously demonstrated for *de novo* human and alligator genomes (Putnam *et al.* 2016). Simultaneously, HiC libraries were constructed from *G. longicalyx* leaf tissue at PhaseGenomics LLC. A second

129 round of HighRise was used to include the HiC data for additional genome scaffolding (Koch  
130 2016; Putnam *et al.* 2016), reducing the contig number to 135 (Longicalyx\_V3.0).

131  
132 High-molecular weight DNA was extracted from young *G. longicalyx* leaves and subsequently  
133 purified, nicked, labeled, and repaired according to Bionano Plant protocol and standard  
134 operating procedures for the Irys platform. BssSI was used in conjunction with the IrysSolve  
135 pipeline to assemble an optical map on the BYU Fulton SuperComputing cluster. The resulting  
136 optical map was aligned to the assembly named Lonigcalyx\_V3.0 using an *in silico* labeled  
137 reference sequence. Bionano maps linked large contigs present in this assembly, producing 17  
138 large scaffolds (Lonigcalyx\_V4.0).

139  
140 Minion sequencing libraries were created and sequenced following the standard protocol from  
141 Oxford Nanopore. Scaffolds from Lonigcalyx\_V4.0 were polished (Supplemental File 1) with  
142 existing Illumina (SRR1174179 and SRR1174182 from the NCBI Short Read Archive) and the  
143 newly generated Minion data for *G. longicalyx* using both PBJelly (English *et al.* 2012) and  
144 GapFiller (Boetzer and Pirovano 2012) to produce the final assembly, Lonigcalyx\_V5.0.

145  
146 **Repeat and gene annotation**

147 Repeats were identified using two methods. The first is a homology-based approach, *i.e.*, a  
148 combination of RepeatMasker (Smit *et al.* 2015) and “One code to find them all” (Bailey-Bechet  
149 *et al.* 2014), whereas the second method (*i.e.*, RepeatExplorer; (Novák *et al.* 2010) clusters reads  
150 based on sequence similarity and automatically annotates the most abundant cluster using  
151 RepeatMasker. Each RepeatMasker run used a custom library, which combines Repbase 23.04  
152 repeats (Bao *et al.* 2015) with cotton-specific repeats. Default parameters were run, except the  
153 run was “sensitive” and was set to mask only TEs (no low-complexity). Parameters are available



154 <https://github.com/Wendellab/longicalyx>. “One code to find them all” was used to aggregate  
155 multiple hits from the first method (RepeatMasker) into TE models using default parameters.  
156 The resulting output was aggregated and summarized in R/3.6.0 (R Core Team 2017) using  
157 *dplyr* /0.8.1 (Wickham *et al.* 2015). Cluster results were obtained from (Grover *et al.* 2019) and  
158 [https://github.com/IGBB/D\\_Cottons\\_USDA](https://github.com/IGBB/D_Cottons_USDA), and these were parsed in R/3.6.0 (R Core Team  
159 2017). All code is available at <https://github.com/Wendellab/longicalyx>.  
160  
161 RNA-Seq libraries were generated from *G. longicalyx* leaf (CL), floral (FF), and stem tissues  
162 (FS) to improve genome annotation. RNA-seq libraries were independently constructed by BGI  
163 Americas (Davis, CA) using Illumina TruSeq reagents and subsequently sequenced (single-end,  
164 50 bp). The newly sequenced *G. longicalyx* RNA-seq was combined with existing RNA-seq  
165 from *G. longicalyx* (SRR1174179) as well as two closely related species, i.e., *G. herbaceum*  
166 (developing fibers and seed; PRJNA595350 and SRR959585, respectively) and *G. arboreum* (5  
167 seed libraries and 1 seedling; SRR617075, SRR617073, SRR617068, SRR617067, SRR959590,  
168 and SRR959508). RNA-seq libraries were mapped to the hard-masked *G. longicalyx* genome  
169 using hisat2 [v2.1.0] (Kim *et al.* 2015). BRAKER2 [v2.1.2] (Hoff *et al.* 2019) was used in  
170 conjunction with GeneMark [v4.36] (Borodovsky and Lomsadze 2011) generated annotations to  
171 train Augustus [v3.3.2] (Stanke *et al.* 2006). Mikado [v1.2.4] (Venturini *et al.* 2018) was used to  
172 produce high quality RNA-seq based gene predictions by combining the RNA-seq assemblies  
173 produced by StringTie [v1.3.6] (Pertea *et al.* 2015) and Cufflinks [v2.2.1] (Ghosh and Chan  
174 2016) with a reference-guided assembly from Trinity [v2.8.5] (Grabherr *et al.* 2011) and a splice  
175 junction analysis from Portcullis [v1.2.2] (Mapleson *et al.* 2018). The Trinity assembly was  
176 formatted using GMAP [v2019-05-12] (Wu and Watanabe 2005). MAKER2 [v2.31.10] (Holt  
177 and Yandell 2011; Campbell *et al.* 2014) was used to integrate gene predictions from (1)  
178 BRAKER2 trained Augustus, (2) GeneMark, and (3) Mikado, also using evidence from all  
179 *Gossypium* ESTs available from NCBI (nucleotide database filtered on “txid3633” and “is\_est”)  
180 and a database composed of all curated proteins in Uniprot Swissprot [v2019\_07] (UniProt  
181 Consortium 2008) combined with the annotated proteins from the *G. hirsutum*  
182 ([https://www.cottongen.org/species/Gossypium\\_hirsutum/jgi-AD1\\_genome\\_v1.1](https://www.cottongen.org/species/Gossypium_hirsutum/jgi-AD1_genome_v1.1)) and *G.*  
183 *raimondii* (Paterson *et al.* 2012) genomes. Maker scored each gene model using the annotation  
184 edit distance (AED - (Eilbeck *et al.* 2009; Holt and Yandell 2011; Yandell and Ence 2012)  
185 metric based on EST and protein evidence provided. Gene models with an AED greater than  
186 0.47 were removed from further analyses, and the remaining gene models were functionally  
187 annotated using InterProScan [v5.35-74.0] (Jones *et al.* 2014) and BlastP [v2.9.0+] (Camacho *et*  
188 *al.* 2009) searches against the Uniprot SwissProt database. Orthologs between the *G. longicalyx*  
189 annotations and the existing annotations for *G. arboreum* (Du *et al.* 2018), *G. raimondii*  
190 (Paterson *et al.* 2012), *G. hirsutum* (Hu *et al.* 2019), and *G. barbadense* (Hu *et al.* 2019) were  
191 predicted by OrthoFinder using default settings (Emms and Kelly 2015, 2019). All genomes are  
192 hosted through CottonGen (<https://www.cottongen.org>; (Yu *et al.* 2014)) and running parameters  
193 are available from <https://github.com/Wendellab/longicalyx>.

194

## 195 **ATAC-seq and data analysis**

196 ATAC-seq was performed as described previously (Lu *et al.* 2017). For each replicate,  
197 approximately 200 mg freshly collected leaves or flash frozen leaves were immediately chopped  
198 with a razor blade in ~ 1 ml of pre-chilled lysis buffer (15 mM Tris-HCl pH 7.5, 20 mM NaCl,  
199 80 mM KCl, 0.5 mM spermine, 5 mM 2-mercaptoethanol, 0.2% Triton X-100). The chopped  
200 slurry was filtered twice through miracloth and once through a 40 µm filter. The crude nuclei  
201 were stained with DAPI and loaded into a flow cytometer (Beckman Coulter MoFlo XDP).  
202 Nuclei were purified by flow sorting and washed in accordance with Lu et al (Lu *et al.* 2017).  
203 The sorted nuclei were incubated with 2 µl Tn5 transposomes in 40 µl of tagmentation buffer (10  
204 mM TAPS-NaOH pH 8.0, 5 mM MgCl<sub>2</sub>) at 37°C for 30 minutes without rotation. The integration  
205 products were purified using a Qiagen MinElute PCR Purification Kit or NEB Monarch™ DNA  
206 Cleanup Kit and then amplified using Phusion DNA polymerase for 10-13 cycles. PCR cycles  
207 were determined as described previously (Buenrostro *et al.* 2013). Amplified libraries were  
208 purified with AMPure beads to remove primers. ATAC-seq libraries were sequenced in paired-  
209 end 35 bp at the University of Georgia Genomics & Bioinformatics Core using an Illumina  
210 NextSeq 500 instrument.

211

212 Reads were adapter and quality trimmed, and then filtered using “Trim Galore” [v0.4.5]  
213 (Krueger 2015). Clean reads were subsequently aligned to the *Lonigcalyx\_V5.0* assembly using  
214 Bowtie2 [v2.3.4] (Langmead and Salzberg 2012) with the parameters “--no-mixed --no-  
215 discordant --no-unal --dovetail”. Duplicate reads were removed using Picard [v2.17.0] with  
216 default parameters (<http://broadinstitute.github.io/picard/>). Only uniquely mapped read pairs with  
217 a quality score of at least 20 were kept for peak calling. Phantompeakqualtools [v1.14] (Landt *et*  
218 *al.* 2012) was used to calculate the strand cross-correlation, and deepTools [v2.5.2] (Ramírez *et*  
219 *al.* 2016) was used to calculate correlation between replicates. The peak calling tool from  
220 HOMER [v4.10] (Heinz *et al.* 2010), i.e., *findpeaks*, was run in “region” mode and with the  
221 minimal distance between peaks set to 150 bp. MACS2 [v2.1.1] (Zhang *et al.* 2008) *callpeak*, a  
222 second peak-calling algorithm, was run with the parameter “-f BAMPE” to analyze only properly  
223 paired alignments, and putative peaks were filtered using default settings and false discovery rate  
224 (FDR) < 0.05. Due to the high level of mapping reproducibility (Pearson’s correlation  $r = 0.99$   
225 and Spearman correlation  $r = 0.77$  by deepTools), peaks were combined and merged between  
226 replicates for each tool using BEDTools [v2.27.1] (Quinlan 2014). BEDTools was also used to  
227 intersect HOMER peaks and MACS2 peaks to only retain peak regions identified by both tools  
228 as accessible chromatin regions (ACRs) for subsequent analyses.

229

230 ACRs were annotated in relation to the nearest annotated genes in the R environment [v3.5.0] as  
231 genic (gACRs; overlapping a gene), proximal (pACRs; within 2 Kb of a gene) or distal (dACRs;  
232 >2 Kb from a gene). Using R package ChIPseeker [v1.18.0] (Yu *et al.* 2015), the distribution of  
233 ACRs was calculated around transcription start sites (TSS) and transcription termination sites



234 (TTS), and peak distribution was visualized with aggregated profiles and heatmaps. To compare  
235 GC contents between ACRs and non-accessible genomic region, the BEDTools *shuffle* command  
236 was used to generate the distal (by excluding genic and 2 Kb flanking regions) and  
237 genic/proximal control regions (by including genic and 2 Kb flanking regions), and the *nuc*  
238 command was used to calculate GC content for each ACR and permuted control regions.

239

#### 240 **Identification of the Ren<sup>Lon</sup> region in *G. longicalyx***

241 Previous research (Dighe *et al.* 2009; Zheng *et al.* 2016) identified a marker (BNL1231) that  
242 consistently cosegregates with resistance and that is flanked by the SNP markers GI\_168758 and  
243 GI\_072641, which are all located in the region of *G. longicalyx* chromosome 11 referred to as  
244 “Ren<sup>Lon</sup>”. These three markers were used as queries of gmap (Wu and Watanabe 2005) against  
245 the assembled genome to identify the genomic regions associated with each. The coordinates  
246 identified by gmap were placed in a bed file; this file was used in conjunction with the *G.*  
247 *longicalyx* annotation and BEDtools intersect (Quinlan 2014) to identify predicted *G. longicalyx*  
248 genes contained within Ren<sup>Lon</sup>. Samtools faidx (Li *et al.* 2009) was used to extract the 52  
249 identified genes from the annotation file, which were functionally annotated using blast2go  
250 (blast2go basics; biobam) and including blastx (Altschul *et al.* 1990), gene ontology (The Gene  
251 Ontology Consortium 2019), and InterPro (Jones *et al.* 2014). Orthogroups containing each of  
252 the 52 Ren<sup>Lon</sup> genes were identified from the Orthofinder results (see above).

253

#### 254 **Comparison between *G. arboreum* and *G. longicalyx* for fiber evolution**

255 Whole-genome alignments were generated between *G. longicalyx* and either *G. arboreum*, *G.*  
256 *raimondii*, *G. turneri*, *G. hirsutum* (A-chromosomes), and *G. barbadense* (A-chromosomes)  
257 using Mummer (Marçais *et al.* 2018) and visualized using dotPlotly  
258 (<https://github.com/tpoorten/dotPlotly>) in R (version 3.6.0) (R Development Core Team and  
259 Others 2011). Divergence between *G. longicalyx* and *G. arboreum* or *G. raimondii* was  
260 calculated using orthogroups that contain a single *G. longicalyx* gene with a single *G. arboreum*  
261 and/or single *G. raimondii* gene. Pairwise alignments between *G. longicalyx* and *G. arboreum* or  
262 *G. raimondii* were generated using the *linsi* from MAFFT (Kato and Standley 2013). Pairwise  
263 distances between *G. longicalyx* and *G. arboreum* and/or *G. raimondii* were calculated in R  
264 (version 3.6.0) using phangorn (Schliep 2011) and visualized using ggplot2 (Wickham 2016). To  
265 identify genes unique to species with spinnable fiber (i.e., *G. arboreum* and the polyploid  
266 species), we extracted any *G. arboreum* gene contained within orthogroups composed solely of  
267 *G. arboreum* or polyploid A-genome gene annotations, and subjected these to blast2go (as  
268 above). Syntenic conservation of genes contained within the Ren<sup>Lon</sup> region, as compared to *G.*  
269 *arboreum*, was evaluated using GEvo as implemented in SynMap via COGE (Lyons and  
270 Freeling 2008; Haug-Baltzell *et al.* 2017).

271

272

273

## 274 **Data availability**

275 The assembled genome sequence of *G. longicalyx* is available at NCBI SUB6483233 and  
276 CottonGen (<https://www.cottongen.org/>). The raw data for *G. longicalyx* are also available at  
277 NCBI PRJNA420071 for PacBio and Minion, and PRJNA420070 for RNA-Seq. Supplemental  
278 files are available from figshare.

279

## 280 **Results and Discussion**

### 281 **Genome assembly and annotation**

282 We report a *de novo* genome sequence for *G. longicalyx*. This genome was first assembled from  
283 ~144x coverage (raw) of PacBio reads, which alone produced an assembly consisting of 229  
284 contigs with an N50 of 28.8MB (Table 1). The contigs were scaffolded using a combination of  
285 Chicago Highrise, Hi-C, and BioNano to produce a chromosome level assembly consisting of 17  
286 contigs with an average length of 70.4 Mb (containing only 8.4kb of gap sequence). Thirteen of  
287 the chromosomes were assembled into single contigs. Exact placement of the three unscaffolded  
288 contigs (~100 kb) was not determined, but these remaining sequences were included in NCBI  
289 with the assembled chromosomes. The final genome assembly size was 1190.7 MB, representing  
290 over 90% of the estimated genome size (Hendrix and Stewart 2005).

291

**Table 1. Statistics for assembly versions**

Method	<i>G. longicalyx</i> assemblies*			
	Longicalyx_V1.0 PacBio/Canu	Longicalyx_V3.0 +Chicago HighRise+HiC	Longicalyx_V4.0 +BioNano	Longicalyx_V5.0 +Illumina+Minion
Coverage	79.45			
Total Contig Number	229	135	17	17
Assembly Length**	1196.17 Mb	1196.19 Mb	1190.66 Mb	1190.67 Mb
Average Contig Length	5.22 Mb	8.86 Mb	70.04 Mb	70.04 Mb
Total Length of Ns	0	18200	18000	8488
N50 value is	28.88 Mb	95.88 Mb	95.88 Mb	95.88 Mb
N90 value is	7.58 Mb	76.48 Mb	76.48 Mb	76.29 Mb

\* Statistics for Longicalyx\_V2.0 not calculated

\*\* Genome size for *G. longicalyx* is 1311 (Hendrix and Stewart, 2005)

292

293

294 BUSCO analysis of the completed genome (Waterhouse *et al.* 2017) recovered 95.8% complete  
295 BUSCOs (from the total of 2121 BUSCO groups searched; Table 2). Most BUSCOs (86.5%)  
296 were both complete and single copy, with only 9.3% BUSCOs complete and duplicated. Less  
297 than 5% of BUSCOs were either fragmented (1.4%) or missing (2.8%), indicating a general  
298 completeness of the genome. Genome contiguity was independently verified using the LTR  
299 Assembly Index (LAI) (Ou *et al.* 2018), which is a reference-free method to assess genome  
300 contiguity by evaluating the completeness of LTR-retrotransposon assembly within the genome.

301 This method, applied to over 100 genomes in Phytozome, suggests that an LAI between 10 and  
 302 20 should be considered “reference-quality”; the *G. longicalyx* genome reported here received an  
 303 LAI score of 10.74. Comparison of the *G. longicalyx* genome to published cotton genomes  
 304 (Table 2) suggests that the quality of this assembly is similar or superior to other currently  
 305 available cotton genomes.

306  
 307 **Table 2.** BUSCO and LAI scores for the *G. longicalyx* genome compared to existing cotton genomes.  
 308

	Complete BUSCO			Incomplete BUSCO		LAI score	Reference
	Total	Single	Duplicated	Fragmented	Missing		
<i>G. longicalyx</i>	95.80%	86.50%	9.30%	1.40%	2.80%	10.74	
<i>G. turneri</i>	95.80%	86.00%	9.80%	1.00%	3.20%	8.51	(Udall <i>et al.</i> 2019)
<i>G. raimondii</i> (BYU)	92.80%	85.10%	7.70%	2.70%	4.50%	10.57	(Udall <i>et al.</i> 2019)
<i>G. raimondii</i> (JGI)	98.00%	87.30%	10.70%	0.70%	1.30%	8.51	(Paterson <i>et al.</i> 2012)
<i>G. arboreum</i> (CRI)	94.70%	85.20%	9.50%	1.00%	4.30%	12.59	(Du <i>et al.</i> 2018)
<i>G. barbadense</i> 3-79 (HAU v2)	96.30%	12.20%	84.10%	0.80%	2.90%	10.38	(Wang <i>et al.</i> 2019)
<i>G. hirsutum</i> TM1 (HAU v1)	97.70%	14.50%	83.20%	0.50%	1.80%	10.61	(Wang <i>et al.</i> 2019)

309  
 310 Genome annotation produced 40,181 transcripts representing 38,378 unique genes.  
 311 Comparatively, the reference sequences for the related diploids *G. raimondii* (Paterson *et al.*  
 312 2012) and *G. arboreum* (Du *et al.* 2018) recovered 37,223 and 40,960 genes, respectively.  
 313 Ortholog analysis between *G. longicalyx* and both diploids suggests a simple 1:1 relationship  
 314 between a single *G. longicalyx* gene and a single *G. raimondii* or *G. arboreum* gene for 67-68%  
 315 of the *G. longicalyx* genes (25,637 and 26,249 genes, respectively; Table 3). Approximately 3-  
 316 4% of the *G. longicalyx* genome (i.e., 1,153-1,438 genes) are in “one/many” (Table 3)  
 317 relationships whereby one or more *G. longicalyx* gene model(s) matches one or more *G.*  
 318 *raimondii* or *G. arboreum* gene model(s). The remaining 5,009 genes were not placed in  
 319 orthogroups with any other cotton genome, slightly higher than the 2,016 - 2,556 unplaced genes  
 320 in the other diploid species used here. While this could be partly due to genome annotation  
 321 differences in annotation pipelines, it is also likely due to differences in the amount of RNA-seq  
 322 available for each genome.

323

324 **Table 3.** Orthogroups between *G. longicalyx* and two related diploid species. Numbers of genes are listed  
 325 and percentages are in parentheses. Relationships listed in the last four lines of the table represent  
 326 one/many *G. longicalyx* genes relative to one or many genes from *G. arboreum* or *G. raimondii*.

	<i>G. longicalyx</i>	<i>G. arboreum</i>	<i>G. raimondii</i>
Number of genes	38,378	40,960	37,223
Genes in orthogroups	33,369 (86.9%)	38,404 (93.8%)	35,207 (94.6%)
Unassigned genes	5,009 (13.1%)	2,556 (6.2%)	2,016 (5.4%)
Orthogroups containing species	26,591 (78.5%)	29,763 (87.8%)	29,153 (86.0%)
Genes in species-specific orthogroups	74 (0.2%)	0	8 (0.0%)
1-to-1 relationship		26,249 (70.5%)	25,637 (68.9%)
1-to-many relationship		1,207 (3.2%)	1,153 (3.1%)
many-to-1 relationship		1,438 (3.9%)	1,172 (3.1%)
many-to-many relationship		513 (1.4%)	290 (0.8%)

327  
 328 **Repeats**  
 329 Transposable element (TE) content was predicted for the genome, both by *de novo* TE prediction  
 330 (Bailly-Bechet *et al.* 2014; Smit *et al.* 2015) and repeat clustering (Novák *et al.* 2010). Between  
 331 44 - 50% of the *G. longicalyx* genome is inferred to be repetitive by RepeatMasker and  
 332 RepeatExplorer, respectively. While estimates for TE categories (e.g., DNA, Ty3/*gypsy*,  
 333 Ty1/*copia*, etc.) were reasonably consistent between the two methods (Table 4), RepeatExplorer  
 334 recovered nearly 100 additional megabases of putative repetitive sequences, mostly in the  
 335 categories of Ty3/*gypsy*, unspecified LTR elements, and unknown repetitive elements.  
 336 Interestingly, RepeatMasker recovered a greater amount of sequence attributable to Ty1/*copia*  
 337 and DNA elements (Table 4); however, this only accounted for 22 Mb (less than 20% of the total  
 338 differences over all categories). The difference between methods with respect to each category  
 339 and the total TE annotation is relatively small and may be attributable to a combination of  
 340 methods (homology-based TE identification method versus similarity clustering), the under-  
 341 exploration of the cotton TE population, and sensitivity differences in each method with respect  
 342 to TE age/abundance.

343  
 344 **Table 4.** Comparison between repeat quantification methods for the *G. longicalyx* genome. Amounts are  
 345 given in megabases (Mb).

	<i>RepeatExplorer</i>	<i>RepeatMasker/OneCode</i>
LTR/Gypsy (Ty3)	557	513
LTR/Copia (Ty1)	39	48
LTR, unspecified	44	0
DNA (all element types)	2.3	15
unknown	18	0
Total repetitive clustered	660	575

346  
 347 Because the RepeatExplorer pipeline allows simultaneous analysis of multiple samples (i.e., co-  
 348 clustering), we used that repeat profile for both description and comparison to the closely related  
 349 sister species, *G. herbaceum* and *G. arboreum* (from subgenus *Gossypium*). Relative to other  
 350 cotton species, *G. longicalyx* has an intermediate amount of TEs, as expected from its  
 351 intermediate genome size (1311 Mb; genome size range for *Gossypium* diploids = 841 - 2778  
 352 Mb). Approximately half of the genome (660 Mb) is composed of repetitive sequences,  
 353 somewhat less than the closely related sister (A-genome) clade, which are slightly bigger in total  
 354 size and have slightly more repetitive sequence (~60% repetitive; Table 5). Over 80% of the *G.*  
 355 *longicalyx* repetitive fraction is composed of Ty3/gypsy elements, a similar proportion to the  
 356 proportion of Ty3/gypsy in subgenus *Gossypium* genomes. Most other element categories were  
 357 roughly similar in total amount and proportion between *G. longicalyx* and the two species from  
 358 subgenus *Gossypium* (Figure 2).

359

360 Table 5. Transposable element content in *G. longicalyx* versus the sister clade (section *Gossypium*)

	<i>Subgenus Longiloba</i>		<i>Subgenus Gossypium</i>	
	<i>G. longicalyx</i>		<i>G. herbaceum</i>	<i>G. arboreum</i>
Genome Size	1311		1667	1711
LTR/Gypsy (Ty3)	557		876	943
LTR/Copia (Ty1)	39		43	41
LTR, unspecified	44		62	57
DNA (all element types)	2.3		2.7	2.4
unknown	18		27	25
Total repetitive clustered	660		1011	1067
% genome is repetitive	50%		61%	62%
% genome is gypsy	42%		53%	55%
% repetitive is gypsy	84%		87%	88%

361

362

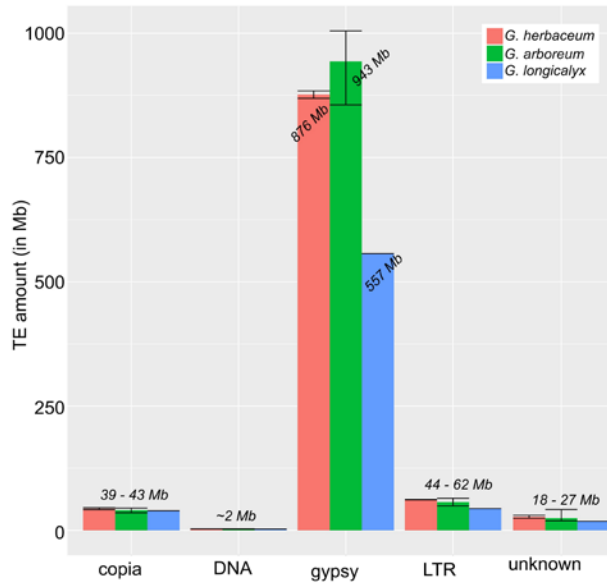
363

364

365

366

**Figure 2.** Repetitive content in *G. longicalyx* relative to the related diploid species *G. herbaceum* and *G. arboreum*.



381

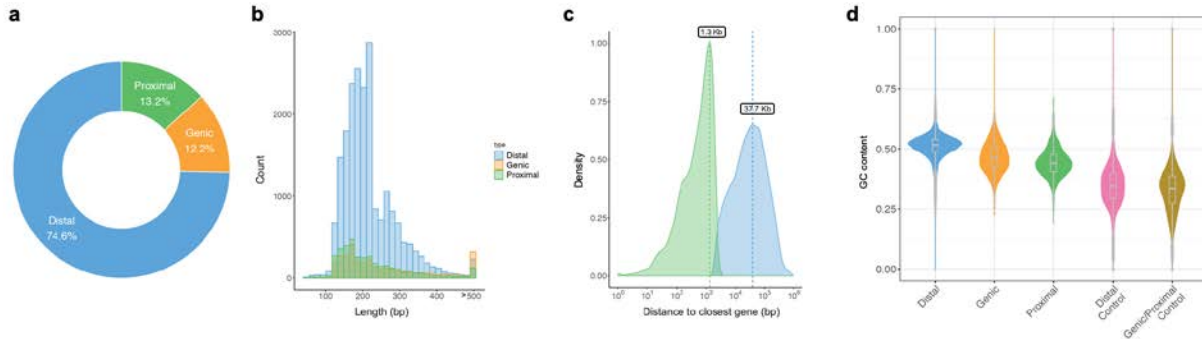
### 382 **Chromatin accessibility in *G. longicalyx***

383 We performed ATAC-seq to map accessible chromatin regions (ACRs) in leaves. Two replicated  
384 ATAC-seq libraries were sequenced to ~25.7 and ~45.0 million reads per sample. The strand  
385 cross-correlation statistics supported the high quality of the ATAC-seq data, and the correlation  
386 of mapping read coverages (Pearson  $r = 0.99$  and Spearman  $r = 0.77$ ) suggested a high level of  
387 reproducibility between replicates (Supplemental Table 1). A total of 28,030 ACRs (6.4 Mb)  
388 were identified ranging mostly from 130 bp to 400 bp in length, which corresponds to ~0.5% of  
389 the assembled genome size (Supplemental Table 2). The enrichment of ACRs around gene  
390 transcription start sites (Supplemental Figure 1) suggested that these regions were functionally  
391 important and likely enriched with *cis*-regulatory elements. Based on proximity to their nearest  
392 annotated genes, these ACRs were categorized as genic (gACRs; overlapping a gene), proximal  
393 (pACRs; within 2 Kb of a gene) or distal (dACRs; >2 Kb from a gene). The gACRs and pACRs  
394 represented 12.2% and 13.2% of the total number of ACRs (952 Kb and 854 Kb in size,  
395 respectively), while approximately 75% (4.6 Mb) were categorized as dACRs, a majority of  
396 which were located over 30 Kb from the nearest gene (Figure 3). This high percentage of dACRs  
397 is greater than expected (~40% of 1 GB genome) given previous ATAC-seq studies in plants (Lu  
398 *et al.* 2019; Ricci *et al.* 2019) and may reflect challenges in annotating rare transcripts. While  
399 more thorough, species-specific RNA-seq will improve later annotation versions and refine our  
400 understanding of ACR proximity to genes, we do note that our observation of abundant dACRs  
401 and potentially long-range *cis*-regulatory elements is consistent with previous results (Lu *et al.*  
402 2019; Ricci *et al.* 2019) The dACRs discovered here were the most GC-rich, followed by gACRs  
403 and pACRs (52%, 46%, and 44%, respectively), all of whom had GC contents significantly  
404 higher than randomly selected control regions with the same length distribution (Figure 3d).  
405 Because high GC content is associated with several distinct features that can affect the *cis*-



406 regulatory potential of a sequence (Landolin *et al.* 2010; Wang *et al.* 2012), these results support  
407 the putative regulatory functions of ACRs.

408  
409



410  
411 **Figure 3.** Accessible chromatin regions (ACRs) in the *G. longicalyx* genome. **a.** Categorization  
412 of ACRs in relation to nearest gene annotations - distal dACRs, proximal pACRs, and genic  
413 gACRs. **b.** Length distribution of ACRs that were identified by both HOMER and MACS2  
414 contained within various genomic regions. **c.** Distance of gACRs and pACRs to nearest  
415 annotated genes. **d.** Boxplot of GC content in ACRs and control regions.

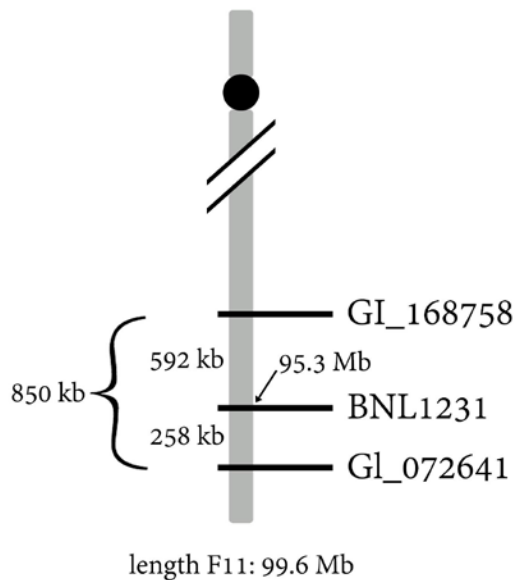
416  
417

### 418 **Genomics of *G. longicalyx* reniform nematode resistance**

419 Reniform nematode is an important cotton parasite that results in stunted growth, delayed  
420 flowering and/or fruiting, and a reduction in both yield quantity and quality (Robinson 2007;  
421 Khanal *et al.* 2018). While domesticated cotton varieties are largely vulnerable to reniform  
422 nematode (Robinson *et al.* 1997), nematode resistance is found in some wild relatives of  
423 domesticated cotton, including *G. longicalyx*, which is nearly immune (Yik and Birchfield 1984).  
424 Recent efforts to elucidate the genetic underpinnings of this resistance in *G. longicalyx* (i.e.,  
425 Ren<sup>Lon</sup>) identified a marker (BNL1231) that consistently cosegregates with resistance and is  
426 flanked by the SNP markers GL\_168758 and GL\_072641 (Dighe *et al.* 2009; Zheng *et al.* 2016).  
427 Located in chromosome 11, this region contains one or more closely-linked nearly dominant  
428 gene(s) (Dighe *et al.* 2009) that confer hypersensitivity to reniform infection (Khanal *et al.*  
429 2018), resulting in the “stunting” phenotype; however, the possible effects of co-inherited R-  
430 genes has not been eliminated. Because the introgressed segment recombines at reduced rates in  
431 interspecific crosses, it has been difficult to fine-map the gene(s) of interest. Additionally,  
432 progress from marker-assisted selection has been lacking, as no recombinants have possessed the  
433 desired combination of reniform resistance and “non-stunting” (Zheng *et al.* 2016). Therefore,  
434 more refined knowledge of the position, identity of the resistance gene(s), mode(s) of immunity  
435 and possible causes of “stunting” will likely catalyze progress on nematode resistance.

436  
437  
438

439  
440



441 **Figure 4.** Diagram of the Ren<sup>Lon</sup> region in *G. longicalyx*. Marker BNL1231, which co-segregates with nematode resistance, is located at approximately 95.3 Mb on chromosome F11.

BLAST analysis of the three Ren<sup>Lon</sup>-associated markers (above) to the assembled *G. longicalyx* genome identifies an 850 kb region on chromosome F11 (positions 94747040..95596585; Figure 4) containing 52 predicted genes (Supplemental Table 3). Functional annotation reveals that over half of the genes (29, or 56%) are annotated as “TMV resistance protein N-like” or similar. In tobacco,

456 TMV resistance protein N confers a hypersensitive

457 response to the presence of the tobacco mosaic virus (TMV; (Erickson *et al.* 1999). Homologs of  
458 this gene in different species can confer resistance to myriad other parasites and pathogens,  
459 including aphid and nematode resistance in tomato (Rossi *et al.* 1998); fungal resistance in  
460 potato (Hehl *et al.* 1999) and flax (Ellis *et al.* 2007); and viral resistance in pepper (Guo *et al.*  
461 2017). Also included in this region are 6 genes annotated as strictosidine synthase-like (SSL),  
462 which may also function in immunity and defense (Sohani *et al.* 2009). While the six SSL-like  
463 genes are tandemly arrayed without disruption, several other genes are intercalated within the  
464 array of TMV resistance-like genes, including the 6 SSL-like genes (Supplemental Table 3).  
465

466 Because there is agronomic interest in transferring nematode resistance from *G. longicalyx* to  
467 other species, we generated orthogroups between *G. longicalyx*, the two domesticated polyploid  
468 species (i.e., *G. hirsutum* and *G. barbadense*), and their model diploid progenitors (*G. raimondii*  
469 and *G. arboreum*; Supplemental Table 4; Supplemental File 2). Interestingly, many of the  
470 defense-relevant *G. longicalyx* genes in the Ren<sup>Lon</sup> region did not cluster into orthogroups with  
471 any other species (15 out of 38; Table 6), including 11 of the 29 TMV resistance-related genes in  
472 the Ren<sup>Lon</sup> region, and fewer were found in syntenic positions in *G. arboreum*. Most of the TMV  
473 resistance-related genes that cluster between *G. longicalyx* and other *Gossypium* species are  
474 present in a single, large orthogroup (OG0000022; Table 4), whereas the remaining TMV-  
475 resistance like genes from *G. longicalyx* are commonly in single gene orthogroups. Since disease  
476 resistance (R) proteins operate by detecting specific molecules elicited by the pathogen during  
477 infection (Martin *et al.* 2003), the increased copy number and variability among the *G.*  
478 *longicalyx* TMV-resistance-like genes may suggest specialization among copies.

479  
480  
481

Table 6: Orthogroup identity (by Orthofinder) for defense-related genes in the Ren<sup>Lon</sup> region and the copy number per species. In *G. longicalyx*, this number includes genes found outside of the Ren<sup>Lon</sup> region. *G. hirsutum* and *G. barbadense* copy numbers are split genes found on the A or D chromosomes, or on scaffolds/contigs not placed on a chromosome.

<b>Description</b>	<b>Orthogroup</b>	<b><i>G. longicalyx</i> gene in Ren<sup>Lon</sup> region</b>	<b><i>G. longicalyx</i></b>	<b><i>G. arboreum</i></b>	<b><i>G. raimondii</i></b>	<b><i>G. hirsutum</i></b>	<b><i>G. barbadense</i></b>
adenyl-sulfate kinase 3-like	OG0053444	Golon.011G359300*	1	---	---	---	---
L-type lectin-domain containing receptor kinase IV.2-like	OG0053450	Golon.011G361200	1	---	---	---	---
T-complex protein 1 subunit theta-like	OG0053447	Golon.011G360400	1	---	---	---	---
protein STRICTOSIDINE SYNTHASE-LIKE 10-like	OG0000242	Golon.011G363400	6	4	2	6 A	9 A, 5 scaffold
		Golon.011G363500					
		Golon.011G363600**					
		Golon.011G363700					
		Golon.011G363800					
OG0053454	Golon.011G363300	1	---	---	---	---	
TMV resistance protein N-like	OG0000022	Golon.011G360100	25	22	5	10 A, 22 D	12 A, 21 D, 1 scaffold
		Golon.011G360300					
		Golon.011G360500					
		Golon.011G360700					
		Golon.011G360800					
		Golon.011G361000					
		Golon.011G361100					
		Golon.011G361400					
		Golon.011G361900					
		Golon.011G362000					
		Golon.011G362400					
		Golon.011G362700					
		Golon.011G362800					
		Golon.011G362900*					
		Golon.011G364000					

		Golon.011G359900					
	OG0028874*	Golon.011G362600	4	---	---	---	---
	OG0028544	Golon.011G363200	3	---	---	1 A	---
	OG0030067	Golon.011G360200	1	---	---	2 A	---
	OG0030069	Golon.011G362500	1	---	---	1 A	1 A
	OG0053445	Golon.011G359800	1	---	---	---	---
	OG0053446	Golon.011G360000	1	---	---	---	---
	OG0053448	Golon.011G360600	1	---	---	---	---
	OG0053451	Golon.011G361700	1	---	---	---	---
	OG0053452	Golon.011G361800	1	---	---	---	---
	OG0053453	Golon.011G362100	1	---	---	---	---
TMV resistance protein N-like isoform X1	OG0053449	Golon.011G360900	1	---	---	---	---
	OG0028874*	Golon.011G362300	4	---	---	---	---
TMV resistance protein N-like isoform X2	OG0033549	Golon.011G363900	1	---	---	---	1 A

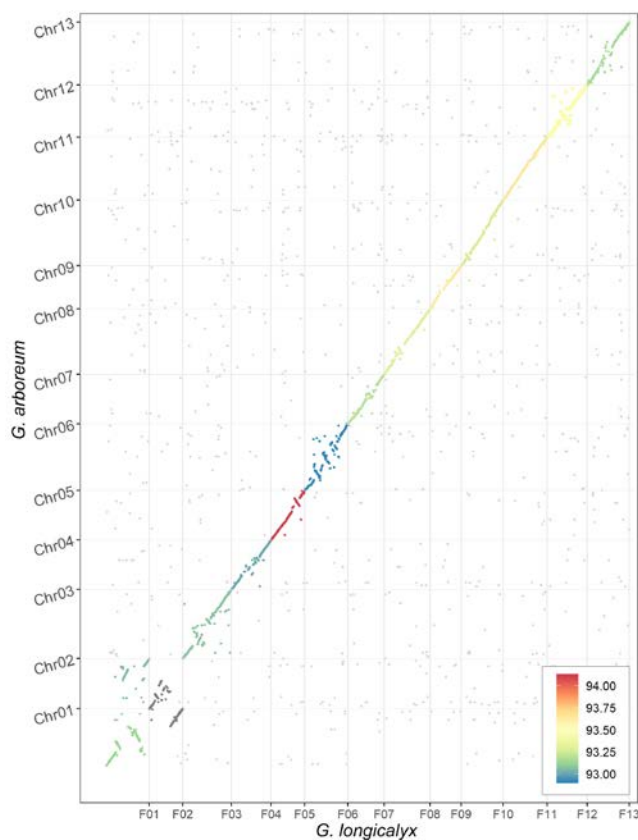
\* This gene is syntenically conserved with *G. arboreum* in the COGE-GEVO analysis.

\*\* This orthogroup is split between two related, but separately named, annotations.

## 483 Comparative genomics and the evolution of spinnable fiber

484 Cotton fiber morphology changed dramatically between *G. longicalyx* and its sister clade,  
485 composed of the A-genome cottons *G. arboreum* and *G. herbaceum*. Whereas *G. longicalyx*  
486 fibers are short and tightly adherent to the seed, A-genome fibers are longer and suitable for  
487 spinning. Accordingly, there has been interest in the changes in the A-genome lineage that have  
488 led to spinnable fiber (Hovav *et al.* 2008; Paterson *et al.* 2012). Progress here has been limited by  
489 the available resources for *G. longicalyx*, relying on introgressive breeding (Nacoulima *et al.*  
490 2012), microarray expression characterization (Hovav *et al.* 2008), and SNP-based surveys  
491 (Paterson *et al.* 2012) of *G. longicalyx* genes relative to *G. herbaceum*. As genomic resources  
492 and surveys for selection are becoming broadly available for the A-genome cottons, our

493 understanding of the evolution of spinnable  
fiber becomes more tangible by the inclusion  
of *G. longicalyx*.



**Figure 5:** Synteny between *G. longicalyx* and domesticated *G. arboreum*. Mean percent identity is illustrated by the color (93-94% identity from blue to red), including intergenic regions.

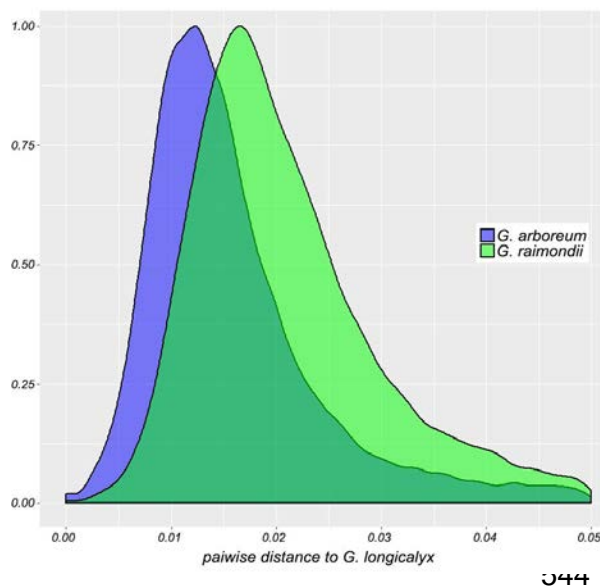
Whole-genome alignment between *G. longicalyx* and the closely related *G. arboreum* (domesticated for long fiber) shows high levels of synteny and overall sequence identity (Figure 5). In general, these two genomes are largely collinear, save for scattered rearrangements and several involving chromosomes 1 and 2; these latter may represent a combination of chromosomal

515 evolution and/or misassembly in one or both genomes. Notably, comparison of *G. longicalyx* to  
516 other recently published genomes (Supplemental Figures 2-5) suggests that an inversion in the  
517 middle of *G. longicalyx* Chr01 exists relative to representatives of the rest of the genus; however,  
518 the other structural rearrangements are restricted to *G. arboreum* and its derived A subgenome in  
519 *G. hirsutum* and *G. barbadense*, suggesting that these differences are limited to comparisons  
520 between *G. longicalyx* and A-(sub)genomes.

521

522 Genic comparisons between *G. longicalyx* and *G. arboreum* suggests a high level of  
523 conservation. Orthogroup analysis finds a one-to-one relationship between these two species for  
524 over 70% of genes. Most of these putative orthologs exhibit <5% divergence (p-distance) in the  
525 coding regions, with over 50% of all putative orthologs exhibiting less than 1.5% divergence.  
526 Comparatively, the median divergence for putative orthologs between *G. longicalyx* and the  
527 more distantly related *G. raimondii* is approximately 2%, with ortholog divergence generally  
528 being higher in the *G. raimondii* comparison (Figure 6).  
529

530



**Figure 6:** Distribution of pairwise p-distances between coding regions of predicted orthologs (i.e., exons only, start to stop) between *G. longicalyx* and either *G. arboreum* (blue) or *G. raimondii* (green). Only orthologs with <5% divergence are shown, which comprises most orthologs in each comparison.

Because *G. longicalyx* represents the ancestor to spinnable fiber, orthogroups containing only *G. arboreum* or polyploid A-genome gene annotations may represent genes important in fiber evolution. Accordingly, we extracted 705 *G.*

545 *arboreum* genes from orthogroups composed solely of *G. arboreum* or polyploid (i.e., *G.*  
546 *hirsutum* or *G. barbadense*) A-genome gene annotations for BLAST and functional annotation.  
547 Of these 705 genes, only 20 represent genes known to influence fiber, i.e., ethylene responsive  
548 genes (10), auxin responsive genes (5), and peroxidase-related genes (5 genes; Supplemental  
549 Table 4). While other genes on this list may also influence the evolution of spinnable fiber,  
550 identifying other candidates will require further study involving comparative coexpression  
551 network analysis or explicit functional studies.

## 552 553 **Conclusion**

554 While several high-quality genome sequences are available for both wild and domesticated  
555 cotton species, each new species provides additional resources to improve both our  
556 understanding of evolution and our ability to manipulate traits within various species. In this  
557 report, we present the first *de novo* genome sequence for *G. longicalyx*, a relative of cultivated  
558 cotton. This genome not only represents the ancestor to spinnable fiber, but also contains the  
559 agronomically desirable trait of reniform nematode immunity. This resource forms a new  
560 foundation for understanding the source and mode of action that provides *G. longicalyx* with this



561 valuable trait, and will facilitate efforts in understanding and exploiting it in modern crop  
562 species.

563

### 564 **Acknowledgements**

565 We thank Emma Miller and Evan Long for technical assistance. We thank the National Science  
566 Foundation Plant Genome Research Program (Grant #1339412) and Cotton Inc. for their  
567 financial support. This research was funded, in part, through USDA ARS Agreements 58-6066-  
568 6-046 and 58-6066-6-059. Support for R.J.S and Z.L. was provided by NSF IOS-1856627 and  
569 the Pew Charitable Trusts. We thank BYU Fulton SuperComputer lab for their resources and  
570 generous support. We also thank ResearchIT for computational support at Iowa State University.  
571 We thank Rise Services for office accommodations in Orem, UT.

572

### 573 **References**

574 Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, 1990 Basic local alignment  
575 search tool. *J. Mol. Biol.* 215: 403–410.

576 Bailly-Bechet, M., A. Haudry, and E. Lerat, 2014 “One code to find them all”: a perl tool to  
577 conveniently parse RepeatMasker output files. *Mob. DNA* 5: 13.

578 Bao, W., K. K. Kojima, and O. Kohany, 2015 Repbase Update, a database of repetitive elements  
579 in eukaryotic genomes. *Mob. DNA* 6: 11.

580 Bell, A., and A. F. Robinson, 2004 Development and characteristics of triple species hybrids  
581 used to transfer reniform nematode resistance from *Gossypium longicalyx* to *Gossypium*  
582 *hirsutum*, pp. 422–426 in *Proceedings of the Beltwide Cotton Conferences*,  
583 [naldc.nal.usda.gov](http://naldc.nal.usda.gov).

584 Birchfield, W., L. R. Brister, and Others, 1963 Susceptibility of cotton and relatives to reniform  
585 nematode in Louisiana. *Plant Disease Reporter* 47: 990–992.

586 Boetzer, M., and W. Pirovano, 2012 Toward almost closed genomes with GapFiller. *Genome*  
587 *Biol.* 13: R56.

588 Borodovsky, M., and A. Lomsadze, 2011 Eukaryotic gene prediction using GeneMark.hmm-E  
589 and GeneMark-ES. *Curr. Protoc. Bioinformatics* Chapter 4: Unit 4.6.1–10.

590 Buenrostro, J. D., P. G. Giresi, L. C. Zaba, H. Y. Chang, and W. J. Greenleaf, 2013  
591 Transposition of native chromatin for fast and sensitive epigenomic profiling of open  
592 chromatin, DNA-binding proteins and nucleosome position. *Nat. Methods* 10: 1213–1218.

593 Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos *et al.*, 2009 BLAST+:  
594 architecture and applications. *BMC Bioinformatics* 10: 421.

595 Campbell, M. S., C. Holt, B. Moore, and M. Yandell, 2014 Genome Annotation and Curation  
596 Using MAKER and MAKER-P. *Curr. Protoc. Bioinformatics* 48: 4.11.1–39.

597 Chen, Z., K. Feng, C. E. Grover, P. Li, F. Liu *et al.*, 2016 Chloroplast DNA Structural Variation,  
598 Phylogeny, and Age of Divergence among Diploid Cotton Species. *PLoS One* 11:  
599 e0157183.

600 Dighe, N. D., A. F. Robinson, A. A. Bell, M. A. Menz, R. G. Cantrell *et al.*, 2009 Linkage

- 601 Mapping of Resistance to Reniform Nematode in Cotton following Introgression from  
602 *Gossypium longicalyx* (Hutch. & Lee). *Crop Sci.* 49: 1151–1164.
- 603 Du, X., G. Huang, S. He, Z. Yang, G. Sun *et al.*, 2018 Resequencing of 243 diploid cotton  
604 accessions based on an updated A genome identifies the genetic basis of key agronomic  
605 traits. *Nat. Genet.* 50: 796–802.
- 606 Eilbeck, K., B. Moore, C. Holt, and M. Yandell, 2009 Quantitative measures for the management  
607 and comparison of annotated genomes. *BMC Bioinformatics* 10: 67.
- 608 Ellis, J. G., P. N. Dodds, and G. J. Lawrence, 2007 Flax rust resistance gene specificity is based  
609 on direct resistance-avirulence protein interactions. *Annu. Rev. Phytopathol.* 45: 289–306.
- 610 Emms, D. M., and S. Kelly, 2019 OrthoFinder: phylogenetic orthology inference for comparative  
611 genomics. *bioRxiv* 466201.
- 612 Emms, D. M., and S. Kelly, 2015 OrthoFinder: solving fundamental biases in whole genome  
613 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol.* 16: 157.
- 614 English, A. C., S. Richards, Y. Han, M. Wang, V. Vee *et al.*, 2012 Mind the gap: upgrading  
615 genomes with Pacific Biosciences RS long-read sequencing technology. *PLoS One* 7:  
616 e47768.
- 617 Erickson, F. L., S. Holzberg, A. Calderon-Urrea, V. Handley, M. Axtell *et al.*, 1999 The helicase  
618 domain of the TMV replicase proteins induces the N-mediated defence response in tobacco.  
619 *Plant J.* 18: 67–75.
- 620 Fryxell, P. A., 1992 A revised taxonomic interpretation of *Gossypium* L (Malvaceae). *Rheeda* 2:  
621 108–165.
- 622 Fryxell, P. A., 1971 PHENETIC ANALYSIS AND THE PHYLOGENY OF THE DIPLOID  
623 SPECIES OF *GOSSYPIUM* L. (MALVACEAE). *Evolution* 25: 554–562.
- 624 Ghosh, S., and C.-K. K. Chan, 2016 Analysis of RNA-Seq Data Using TopHat and Cufflinks.  
625 *Methods Mol. Biol.* 1374: 339–361.
- 626 Grabherr, M. G., B. J. Haas, M. Yassour, J. Z. Levin, D. A. Thompson *et al.*, 2011 Full-length  
627 transcriptome assembly from RNA-Seq data without a reference genome. *Nat. Biotechnol.*  
628 29: 644–652.
- 629 Grover, C. E., M. A. Arick 2nd, A. Thrash, J. L. Conover, W. S. Sanders *et al.*, 2019 Insights  
630 into the Evolution of the New World Diploid Cottons (*Gossypium*, Subgenus *Houzingenia*)  
631 Based on Genome Sequencing. *Genome Biol. Evol.* 11: 53–71.
- 632 Guo, G., S. Wang, J. Liu, B. Pan, W. Diao *et al.*, 2017 Rapid identification of QTLs underlying  
633 resistance to Cucumber mosaic virus in pepper (*Capsicum frutescens*). *Theor. Appl. Genet.*  
634 130: 41–52.
- 635 Haug-Baltzell, A., S. A. Stephens, S. Davey, C. E. Scheidegger, and E. Lyons, 2017 SynMap2  
636 and SynMap3D: web-based whole-genome synteny browsers. *Bioinformatics* 33: 2197–  
637 2198.
- 638 Hehl, R., E. Faurie, J. Hesselbach, F. Salamini, S. Whitham *et al.*, 1999 TMV resistance gene N  
639 homologues are linked to *Synchytrium endobioticum* resistance in potato. *Theor. Appl.*  
640 *Genet.* 98: 379–386.

- 641 Heinz, S., C. Benner, N. Spann, E. Bertolino, Y. C. Lin *et al.*, 2010 Simple Combinations of  
642 Lineage-Determining Transcription Factors Prime cis-Regulatory Elements Required for  
643 Macrophage and B Cell Identities. *Molecular Cell* 38: 576–589.
- 644 Hendrix, B., and J. M. Stewart, 2005 Estimation of the nuclear DNA content of gossypium  
645 species. *Ann. Bot.* 95: 789–797.
- 646 Hoff, K. J., A. Lomsadze, M. Borodovsky, and M. Stanke, 2019 Whole-Genome Annotation  
647 with BRAKER. *Methods Mol. Biol.* 1962: 65–95.
- 648 Holt, C., and M. Yandell, 2011 MAKER2: an annotation pipeline and genome-database  
649 management tool for second-generation genome projects. *BMC Bioinformatics* 12: 491.
- 650 Hovav, R., J. A. Udall, B. Chaudhary, E. Hovav, L. Flagel *et al.*, 2008 The evolution of  
651 spinnable cotton fiber entailed prolonged development and a novel metabolism. *PLoS*  
652 *Genet.* 4: e25.
- 653 Hu, Y., J. Chen, L. Fang, Z. Zhang, W. Ma *et al.*, 2019 *Gossypium barbadense* and *Gossypium*  
654 *hirsutum* genomes provide insights into the origin and evolution of allotetraploid cotton.  
655 *Nat. Genet.* 51: 739–748.
- 656 Hutchinson, J. B., and B. J. S. Lee, 1958 Notes from the East African Herbarium: IX: A New  
657 Species of *Gossypium* from Central Tanganyika. *Kew Bull.* 13: 221–223.
- 658 Jones, P., D. Binns, H.-Y. Chang, M. Fraser, W. Li *et al.*, 2014 InterProScan 5: genome-scale  
659 protein function classification. *Bioinformatics* 30: 1236–1240.
- 660 Katoh, K., and D. M. Standley, 2013 MAFFT multiple sequence alignment software version 7:  
661 improvements in performance and usability. *Mol. Biol. Evol.* 30: 772–780.
- 662 Khanal, C., E. C. McGawley, C. Overstreet, and S. R. Stetina, 2018 The Elusive Search for  
663 Reniform Nematode Resistance in Cotton. *Phytopathology* 108: 532–541.
- 664 Kidwell, K. K., and T. C. Osborn, 1992 Simple plant DNA isolation procedures, pp. 1–13 in  
665 *Plant Genomes: Methods for Genetic and Physical Mapping*, edited by J. S. Beckmann and  
666 T. C. Osborn. Springer Netherlands, Dordrecht.
- 667 Kim, D., B. Langmead, and S. L. Salzberg, 2015 HISAT: a fast spliced aligner with low memory  
668 requirements. *Nat. Methods* 12: 357–360.
- 669 Koch, L., 2016 Chicago HighRise for genome scaffolding. *Nat. Rev. Genet.* 17: 194.
- 670 Kranthi, K. R., 2018 Cotton production practices: snippets from global data 2017. *The ICAC*  
671 *Recorder XXXVI*: 4–14.
- 672 Krueger, F., 2015 Trim galore. A wrapper tool around Cutadapt and FastQC to consistently  
673 apply quality and adapter trimming to FastQ files.
- 674 Landolin, J. M., D. S. Johnson, N. D. Trinklein, S. F. Aldred, C. Medina *et al.*, 2010 Sequence  
675 features that drive human promoter function and tissue specificity. *Genome Res.* 20: 890–  
676 898.
- 677 Landt, S. G., G. K. Marinov, A. Kundaje, P. Kheradpour, F. Pauli *et al.*, 2012 ChIP-seq  
678 guidelines and practices of the ENCODE and modENCODE consortia. *Genome Res.* 22:  
679 1813–1831.
- 680 Langmead, B., and S. L. Salzberg, 2012 Fast gapped-read alignment with Bowtie 2. *Nat.*

- 681           Methods 9: 357–359.
- 682 Lawrence, K., M. Olsen, T. Faske, R. Hutmacher, J. Muller *et al.*, 2015 Cotton disease loss  
683           estimate committee report, 2014, pp. 188–190 in *Proceedings of the 2015 Beltwide Cotton*  
684           *Conferences, San Antonio, TX. Cordova: National Cotton Council.*,
- 685 Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence  
686           Alignment/Map format and SAMtools. *Bioinformatics* 25: 2078–2079.
- 687 Lu, Z., B. T. Hofmeister, C. Vollmers, R. M. DuBois, and R. J. Schmitz, 2017 Combining  
688           ATAC-seq with nuclei sorting for discovery of cis-regulatory regions in plant genomes.  
689           *Nucleic Acids Res.* 45: e41.
- 690 Lu, Z., A. P. Marand, W. A. Ricci, C. L. Ethridge, X. Zhang *et al.*, 2019 The prevalence,  
691           evolution and chromatin signatures of plant regulatory elements. *Nat Plants*.
- 692 Lyons, E., and M. Freeling, 2008 How to usefully compare homologous plant genes and  
693           chromosomes as DNA sequences: How to usefully compare plant genomes. *Plant J.* 53:  
694           661–673.
- 695 Mapleson, D., L. Venturini, G. Kaithakottil, and D. Swarbreck, 2018 Efficient and accurate  
696           detection of splice junctions from RNA-seq with Portcullis. *Gigascience* 7.:
- 697 Marçais, G., A. L. Delcher, A. M. Phillippy, R. Coston, S. L. Salzberg *et al.*, 2018 MUMmer4: A  
698           fast and versatile genome alignment system. *PLoS Comput. Biol.* 14: e1005944.
- 699 Martin, G. B., A. J. Bogdanove, and G. Sessa, 2003 Understanding the functions of plant disease  
700           resistance proteins. *Annu. Rev. Plant Biol.* 54: 23–61.
- 701 Nacoulima, N., J. P. Baudoin, and G. Mergeai, 2012 Introgression of improved fiber fineness  
702           trait in *G. hirsutum* L. from *G. longicalyx* Hutch. & Lee. *Commun. Agric. Appl. Biol. Sci.*  
703           77: 207–211.
- 704 Nichols, R. L., A. Bell, D. Stelly, N. Dighe, F. Robinson *et al.*, 2010 Phenotypic and genetic  
705           evaluation of LONREN germplasm, pp. 798–799 in *Proc. Beltwide Cotton Conf. New*  
706           *Orleans, LA.*,
- 707 Novák, P., P. Neumann, and J. Macas, 2010 Graph-based clustering and characterization of  
708           repetitive sequences in next-generation sequencing data. *BMC Bioinformatics* 11: 378.
- 709 Ou, S., J. Chen, and N. Jiang, 2018 Assessing genome assembly quality using the LTR Assembly  
710           Index (LAI). *Nucleic Acids Res.* 46: e126.
- 711 Paterson, A. H., J. F. Wendel, H. Gundlach, H. Guo, J. Jenkins *et al.*, 2012 Repeated  
712           polyploidization of *Gossypium* genomes and the evolution of spinnable cotton fibres.  
713           *Nature* 492: 423–427.
- 714 Pertea, M., G. M. Pertea, C. M. Antonescu, T.-C. Chang, J. T. Mendell *et al.*, 2015 StringTie  
715           enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat. Biotechnol.*  
716           33: 290–295.
- 717 Phillips, L. L., 1966 the cytology and phylogenetics of the diploid species of GOSSYPIMUM. *Am.*  
718           *J. Bot.* 53: 328–335.
- 719 Putnam, N. H., B. L. O’Connell, J. C. Stites, B. J. Rice, M. Blanchette *et al.*, 2016 Chromosome-  
720           scale shotgun assembly using an in vitro method for long-range linkage. *Genome Res.* 26:

- 721 342–350.
- 722 Quinlan, A. R., 2014 BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Curr.*  
723 *Protoc. Bioinformatics* 47: 11.12.1–34.
- 724 Ramírez, F., D. P. Ryan, B. Grüning, V. Bhardwaj, F. Kilpert *et al.*, 2016 deepTools2: a next  
725 generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* 44: W160–5.
- 726 R Core Team, 2017 *R: A language and environment for statistical computing*. R Foundation for  
727 Statistical Computing., Vienna, Austria.
- 728 R Development Core Team, R., and Others, 2011 *R: A language and environment for statistical*  
729 *computing*.
- 730 Ricci, W. A., Z. Lu, L. Ji, A. P. Marand, C. L. Ethridge *et al.*, 2019 Widespread long-range cis-  
731 regulatory elements in the maize genome. *Nat Plants* 5: 1237–1249.
- 732 Robinson, A. F., 2007 Reniform in U.S. cotton: when, where, why, and some remedies. *Annu.*  
733 *Rev. Phytopathol.* 45: 263–288.
- 734 Robinson, A. F., R. N. Inserra, E. P. Caswell-Chen, N. Vovlas, and A. Troccoli, 1997  
735 *Rotylenchulus* Species: Identification, Distribution, Host Ranges, and Crop Plant Resistance  
736 | *Nematropica*. *Nematropica* 27: 127–180.
- 737 Rossi, M., F. L. Goggin, S. B. Milligan, I. Kaloshian, D. E. Ullman *et al.*, 1998 The nematode  
738 resistance gene Mi of tomato confers resistance against the potato aphid. *Proc. Natl. Acad.*  
739 *Sci. U. S. A.* 95: 9750–9754.
- 740 Schliep, K. P., 2011 phangorn: phylogenetic analysis in R. *Bioinformatics* 27: 592–593.
- 741 Smit, A. F. A., R. Hubley, and P. Green, 2015 RepeatMasker Open-4.0. 2013--2015.
- 742 Sohani, M. M., P. M. Schenk, C. J. Schultz, and O. Schmidt, 2009 Phylogenetic and  
743 transcriptional analysis of a strictosidine synthase-like gene family in *Arabidopsis thaliana*  
744 reveals involvement in plant defence responses. *Plant Biol.* 11: 105–117.
- 745 Stanke, M., O. Keller, I. Gunduz, A. Hayes, S. Waack *et al.*, 2006 AUGUSTUS: ab initio  
746 prediction of alternative transcripts. *Nucleic Acids Res.* 34: W435–9.
- 747 The Gene Ontology Consortium, 2019 The Gene Ontology Resource: 20 years and still GOing  
748 strong. *Nucleic Acids Res.* 47: D330–D338.
- 749 Udall, J. A., E. Long, C. Hanson, D. Yuan, T. Ramaraj *et al.*, 2019 De Novo Genome Sequence  
750 Assemblies of *Gossypium raimondii* and *Gossypium turneri*. G3 g3.400392.2019.
- 751 UniProt Consortium, 2008 The universal protein resource (UniProt). *Nucleic Acids Res.* 36:  
752 D190–5.
- 753 Venturini, L., S. Caim, G. G. Kaithakottil, D. L. Mapleson, and D. Swarbreck, 2018 Leveraging  
754 multiple transcriptome assembly methods for improved gene structure annotation.  
755 *Gigascience* 7.:
- 756 Wang, J., J. Zhuang, S. Iyer, X. Lin, T. W. Whitfield *et al.*, 2012 Sequence features and  
757 chromatin structure around the genomic regions bound by 119 human transcription factors.  
758 *Genome Research* 22: 1798–1812.
- 759 Wang, Maojun, Lili Tu, Daojun Yuan, De Zhu, Chao Shen, Jianying Li, Fuyan Liu, *et al.* 2019.  
760 “Reference Genome Sequences of Two Cultivated Allotetraploid Cottons, *Gossypium*



- 761       *hirsutum* and *Gossypium barbadense*.” *Nature Genetics* 51 (2): 224–29.
- 762 Waterhouse, R. M., M. Seppey, F. A. Simão, M. Manni, P. Ioannidis *et al.*, 2017 BUSCO  
763 applications from quality assessments to gene prediction and phylogenomics. *Mol. Biol.*  
764 *Evol.*
- 765 Wendel, J. F., and V. A. Albert, 1992 Phylogenetics of the Cotton Genus (*Gossypium*):  
766 Character-State Weighted Parsimony Analysis of Chloroplast-DNA Restriction Site Data  
767 and Its Systematic and Biogeographic Implications. *Syst. Bot.* 17: 115–143.
- 768 Wendel, J. F., and C. E. Grover, 2015 Taxonomy and Evolution of the Cotton Genus,  
769 *Gossypium*, pp. 25–44 in *Cotton*, Agronomy Monograph, American Society of Agronomy,  
770 Inc., Crop Science Society of America, Inc., and Soil Science Society of America, Inc.,  
771 Madison, WI.
- 772 Wickham, H., 2016 *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.
- 773 Wickham, H., R. Francois, L. Henry, K. Müller, and Others, 2015 dplyr: A grammar of data  
774 manipulation. R package version 0. 4 3.:
- 775 Wu, T. D., and C. K. Watanabe, 2005 GMAP: a genomic mapping and alignment program for  
776 mRNA and EST sequences. *Bioinformatics* 21: 1859–1875.
- 777 Yandell, M., and D. Ence, 2012 A beginner’s guide to eukaryotic genome annotation. *Nat. Rev.*  
778 *Genet.* 13: 329–342.
- 779 Yik, C. P., and W. Birchfield, 1984 Resistant Germplasm in *Gossypium* Species and Related  
780 Plants to *Rotylenchulus reniformis*. *J. Nematol.* 16: 146–153.
- 781 Yu, J., S. Jung, C.-H. Cheng, S. P. Ficklin, T. Lee *et al.*, 2014 CottonGen: a genomics, genetics  
782 and breeding database for cotton research. *Nucleic Acids Res.* 42: D1229–36.
- 783 Yu, G., L.-G. Wang, and Q.-Y. He, 2015 ChIPseeker: an R/Bioconductor package for ChIP peak  
784 annotation, comparison and visualization. *Bioinformatics* 31: 2382–2383.
- 785 Zhang, Y., T. Liu, C. A. Meyer, J. Eeckhoute, D. S. Johnson *et al.*, 2008 Model-based analysis of  
786 ChIP-Seq (MACS). *Genome Biol.* 9: R137.
- 787 Zheng, X., K. A. Hoegenauer, J. Quintana, A. A. Bell, A. M. Hulse-Kemp *et al.*, 2016 SNP-  
788 Based MAS in Cotton under Depressed-Recombination for <sup>RenLon</sup>-Flanking Recombinants:  
789 Results and Inferences on Wide-Cross Breeding Strategies. *Crop Sci.* 56: 1526–1539.