

1 **Long non-coding RNAs and their proposed functions in fibre**
2 **development of cotton (*Gossypium* spp.)**

3

4 Maojun Wang¹, Daojun Yuan¹, Lili Tu¹, Wenhui Gao¹, Yonghui He¹, Haiyan Hu¹,
5 Pengcheng Wang¹, Nian Liu¹, Keith Lindsey² and Xianlong Zhang^{1*}

6

7 ¹National Key Laboratory of Crop Genetic Improvement, Huazhong Agricultural
8 University, Wuhan 430070, Hubei, China

9 ²Integrative Cell Biology Laboratory, School of Biological and Biomedical Sciences,
10 Durham University, South Road, Durham DH1 3LE, United Kingdom

11 *Corresponding author: Xianlong Zhang

12 E-mail: xlzhang@mail.hzau.edu.cn

13 Tel: +86-27-87280510

14 Fax: +86-27-87280196

15

16 **Summary**

17 Long non-coding RNAs (lncRNAs) are transcripts of at least 200 bp in length, that
18 possess no apparent coding capacity and are involved in various biological regulatory
19 processes. Until now, no systematic identification of lncRNAs has been reported in
20 cotton (*Gossypium* spp.).

21 Here, we describe the identification of 30,550 long intergenic non-coding RNA
22 (lincRNA) loci (50,566 transcripts) and 4,718 long non-coding natural antisense
23 transcript (lncNAT) loci (5,826 transcripts). LncRNAs are rich in repetitive sequences
24 and preferentially expressed in a tissue-specific manner. The detection of abundant
25 genome-specific and/or lineage-specific lncRNAs indicated their weak evolutionary
26 conservation. Approximately 76% of homoeologous lncRNAs exhibit biased
27 expression patterns towards the At or Dt subgenomes. Compared with protein-coding
28 genes, lncRNAs showed overall higher methylation levels and their expression was
29 less affected by gene body methylation.

30 The expression validation in different cotton accessions and co-expression network
31 construction helped identify several functional lncRNA candidates involved in cotton
32 fibre initiation and elongation. Analysis of integrated expression from the
33 subgenomes of lncRNAs generating miR397 and its targets due to genome
34 polyploidization indicated their pivotal functions in regulating lignin metabolism in
35 domesticated tetraploid cotton fibres.

36 This study provides a first comprehensive resource of lncRNAs in *Gossypium*.

37 **Keywords:** cotton lncRNAs methylation polyploidization fibre
38 development

39 **Introduction**

40 Generally, long non-coding RNAs (lncRNAs) are transcripts of at least 200 bp in
41 length, possess no apparent coding capacity but are involved in various biological
42 regulatory processes (Rinn and Chang, 2012). On the basis of their genomic
43 localization with respect to protein-coding genes, lncRNAs can be classified as long
44 intergenic non-coding RNAs (lincRNAs), long non-coding natural antisense
45 transcripts (lncNATs), long intronic non-coding RNAs and overlapping lncRNAs that
46 partially overlap with protein-coding genes (Derrien *et al.*, 2012). Compared to
47 protein-coding genes and even small non-coding RNAs, most lncRNAs lack strong
48 sequence conservation between species (Marques and Ponting, 2009; Necsulea *et al.*,
49 2014). lncRNAs are usually expressed at low levels and often exhibit tissue-specific
50 patterns (Cabili *et al.*, 2011), raising the possibility that lncRNAs regulate tissue
51 development. In animals, lncRNAs have been demonstrated to be involved in
52 chromatin modification, transcriptional regulation and post-transcriptional regulation
53 (Geisler and Coller, 2013; Cech and Steitz, 2014). A recent study shows that lncRNAs
54 may play an important role in *de novo* protein evolution (Ruiz-Orera *et al.*, 2014).

55 With the rapid advances in sequencing technology and transcriptomic analysis,
56 thousands of lncRNAs have been now identified in several plant species. In
57 *Arabidopsis*, more than 6,000 lincRNAs have been identified using Tiling Array and
58 RNA-seq (Liu *et al.*, 2012). More recently, 37,238 lncNATs were identified and their
59 responses to light were characterized (Wang *et al.*, 2014). In a study of the origins of
60 small RNAs, Zhou *et al.* (2009) identified more than 7000 lncNATs in rice. In maize,
61 20,163 lincRNAs were identified by integrating public EST databases and RNA-seq
62 data (Li *et al.*, 2014). The public databases PLncDB and PlantNATsDB store
63 lincRNAs from *Arabidopsis* and lncNATs from 69 plant species, respectively (Chen
64 *et al.*, 2012; Jin *et al.*, 2013).

65 While many sequences have been identified, a detailed functional analysis of
66 plant lncRNAs is still in its infancy. For example, lncNAT COOLAIR and intronic
67 lncRNA COLDAIR have been demonstrated to be vital for vernalization in
68 *Arabidopsis* (Swiezewski *et al.*, 2009; Wang *et al.*, 2014). Viroids, a class of sub-viral
69 plant-pathogenic lncRNAs, can regulate gene expression through a small RNA-guided
70 pathway after their degradation (Navarro *et al.*, 2012). LDMAR in rice was found to

71 be required for normal pollen development under long-day conditions (Ding *et al.*,
72 2012). In addition, the DNA-dependent RNA Polymerase V (Pol V)-dependent
73 lncRNAs are involved in RNA-directed DNA methylation (RdDM) by acting as
74 scaffold RNAs (He *et al.*, 2014; Matzke and Mosher, 2014).

75 Cotton (*Gossypium* spp.) is widely cultivated and utilized for its single-celled
76 fibre in the textile industry and is also an important oilseed crop. *Gossypium* belongs
77 to the Malvaceae and diverged from a common ancestor with *Theobroma cacao*
78 (Paterson *et al.*, 2012; Wang *et al.*, 2012). Generally, the genus *Gossypium* is
79 categorized into 45 diploid species (A-G,K; $2n = 2x = 26$) and 5 tetraploid species
80 (AADD, $2n = 4x = 52$), with genome sizes varying about 3-fold, from ~880 Mb to
81 ~2.5 Gb (Hawkins *et al.*, 2006; Wendel *et al.*, 2010). The tetraploid species were
82 formed approximately 1-2 million years ago by the reunification of two divergent
83 diploid species *Gossypium arboretum* (A2) and *Gossypium raimondii* (D5) (Senchina
84 *et al.*, 2003). Human domestication has produced the high-yielding tetraploid
85 *Gossypium hirsutum* (Upland cotton, AADD, AD1 genome), whereas *Gossypium*
86 *barbadense* (Sea-Island cotton, AADD, AD2 genome) is exploited for the superior
87 length, strength, and fineness of the fibres (Kim and Triplett, 2001). Because of its
88 excellent genetic and genomic resources, cotton is regarded as a good model to study
89 genome polyploidization (Paterson *et al.*, 2012), and the cotton fibre is an excellent
90 experimental system for studying cell fate determination, cell elongation and cell wall
91 formation (Guan and Chen, 2013).

92 Studies on non-coding RNAs in cotton have been largely limited to small RNAs
93 until now, and RNA sequencing has helped identify hundreds of small non-coding
94 RNAs. For example, Wei *et al.* (2013) identified miRNAs expressed during anther
95 development in genetic male sterile and wild type cottons and Yang *et al.* (2013)
96 identified miRNAs in cotton somatic embryogenesis. Gong *et al.* (2013) identified 33
97 miRNA families that were conserved between the A and D genomes. Xue *et al.* (2013)
98 confirmed the expression of 79 miRNA families and identified 257 novel miRNAs
99 related to cotton fibre elongation. Functional analysis of miR828 and miR858
100 identified roles in the regulation of homoeologous MYB2 in allotetraploid *G.*
101 *hirsutum* fibre development (Guan *et al.*, 2014). Recent transgenic analysis of
102 miRNA156/157 indicated a fundamental role in fibre elongation (Liu *et al.*, 2014).

103 We aimed to identify lncRNAs in the allotetraploid cotton species *G. babardense*,
104 following genomic and RNA sequencing. We integrated 162 public unstranded
105 transcriptomic sequencing datasets and generated 9 stranded transcriptomic sequences
106 representing the main tissues of cotton to identify lncRNAs. In total, we identified
107 50,566 lincRNAs and 5,826 lncNATs in *G. babardense*. To assign these lncRNAs to
108 subgenomes, we studied their homoeologous expression bias, and characterized the
109 methylation profiles of lncRNAs and compared them with protein-coding genes. We
110 went on to identify functional lncRNA candidates by differential expression analysis
111 and co-expression network construction during cotton fibre development.

112

113 **Materials and Methods**

114 **Plant material, library construction and sequencing**

115 Plant seeds of cotton accession 3-79 (*Gossypium barbadense*) were sown in the
116 glasshouse. When two fully expanded leaves appeared, root, hypocotyl and leaf were
117 excised separately, frozen immediately in liquid nitrogen and stored at -70°C until
118 use. To collect cotton fibre samples, plants were grown in the field in Wuhan, China.
119 Flowers were tagged at the day of blooming (0 day post anthesis, 0 DPA), and bolls
120 were collected at 10 DPA and 20 DPA (Table S2). Samples from different plants were
121 pooled. Total RNA was isolated from these samples using the Spectrum Plant Total
122 RNA Kit (Sigma-Aldrich). Libraries were constructed using the Illumina TruSeq
123 Stranded RNA Kit following the kit's recommendation. Strand-specific sequencing
124 was performed on the Illumina HiSeq 2000 system (paired end 100 bp reads).

125

126 **Publicly available datasets used in this study**

127 We downloaded 154 RNA datasets of *Gossypium* species from the NCBI Sequence
128 Read Archive collection sequenced on the Illumina platform, which include
129 Zebularine-treated RNA and control datasets released by the Plant Industry of
130 Commonwealth Scientific and Industrial Research Organisation (CSIRO) (Table S1).
131 We downloaded 13 *Gossypium* 454 long reads sequencing datasets from the NCBI
132 Sequence Read Archive and integrated all the public ESTs of cotton (Table S3). We
133 also obtained 4 whole genome DNA methylation sequencing datasets released by

134 Joshua A. Udall laboratory (SRX331701). The 7 small RNA and 3 degradation
135 sequencing datasets of cotton fibre tissues were from our laboratory (Liu *et al.*, 2014).

136

137 **lncRNA identification**

138 All the RNA datasets were processed by removing adaptors and trimming low-quality
139 bases ($Q > 20$). The clean sequencing reads were mapped independently to the
140 *Gossypium barbadense* genome using the spliced read aligner Tophat (Trapnell *et al.*,
141 2009). We then applied two iterations of Tophat alignments proposed by Cabili *et al.*
142 (2011) to maximize the splice junction site information from all samples. We
143 separately assembled the transcriptomes using Cufflinks (Trapnell *et al.*, 2010). The
144 Cuffcompare procedure was applied to compare all the assemblies to the genome
145 annotation of *G. barbadense*.

146 We then adopted 6 steps to identify *bona fide* lncRNAs from the novel and
147 antisense transcripts of transcriptome assemblies: 1) transcripts were removed that
148 were detected in fewer than two experiments; 2) transcripts with mapping coverage
149 less than half of transcript length were removed; 3) transcripts were removed that
150 derived from rRNA and tRNA (cutoff E-value 0.001); 4) transcripts with length less
151 than 200 bp were removed; 5) transcripts were searched against the Swiss-Prot and
152 Pfam databases to eliminate transcripts encoding proteins and protein-coding domains
153 (cutoff E-value 0.001); 6) transcripts were removed that did not pass
154 protein-coding-score test by the Coding Potential Calculator (CPC) and
155 Coding-Non-Coding Index (CNCI) softwares (Sun *et al.*, 2013). The optimized
156 parameters of Coding-Non-Coding Index were trained using a lncRNA dataset from
157 *Arabidopsis* (Liu *et al.*, 2012). To verify the lncRNA identification, the public
158 datasets and ESTs were mapped to the lncRNA transcripts by blastn (E-value cutoff
159 $1e-10$, coverage > 0.8).

160

161 **Expression analysis**

162 We employed the Tophat software (with -G parameter) to map all clean RNA-seq
163 reads to the *G. barbadense* genome. The normalized expression of lncRNA and
164 protein-coding transcripts were estimated using all mapped reads by Cufflinks. The
165 multi-read and fragment bias correction methods embedded in Cufflinks were adopted
166 to improve the accuracy of expression level estimation. The differentially expressed

167 genes were identified using DESeq package (adjusted p value 0.01 and at least
168 two-fold change) (Anders and Huber, 2010).

169

170 **Nearest neighbour analysis**

171 Based on the genome location of the lincRNAs and protein-coding genes, the nearest
172 protein-coding genes around each lincRNA at upstream and downstream positions
173 within 5 kb were identified. For lincNATs, we identified the protein-coding genes on
174 the antisense strand. Pearson correlation was employed to explore the expression
175 relationship between these lincRNA/protein-coding gene and lincNAT/protein-coding
176 gene pairs. The GO terms of nearest protein-coding genes with highly similar
177 expression patterns were mapped to lincRNAs for enrichment analysis, similar to the
178 method described by Pauli *et al.* (2012).

179

180 **Tissue specificity analysis**

181 To determine the tissue specificity of lincRNAs and protein-coding genes, we
182 followed the entropy-based measure suggested by Cabili *et al.* (2011). Expression
183 values of genes in samples were firstly normalized to density vectors. Then, the
184 distance between two tissue expression patterns was defined by JS divergence. Finally,
185 we defined the tissue specificity score per transcript using the maximal tissue
186 specificity score of all tissues.

187

188 **Genome synteny of lincRNA**

189 The scaffolds of the At and Dt subgenomes were aligned to *G. arboreum* and *G.*
190 *raimondii* diploid genomes using LASTZ respectively (Harris, 2007). The best
191 mapping results allowing at least 60% coverage were sorted along the diploid
192 chromosomes to construct pseudochromosomes. The syntenic blocks with at least five
193 genes between At and Dt subgenomes were identified using MCScanX software
194 (Wang *et al.*, 2012). We referred the homoeologous lincRNA pairs based on the
195 overlapping of these transcript loci to syntenic blocks and also evidenced by blastn
196 reciprocal best hits with coverage of at least 90%.

197

198 **Methylation data analysis**

199 After clipping adapters and trimming low quality reads, the clean bisulphate-treated
200 DNA sequencing reads were aligned to the *G. barbadense* genome using Bismark
201 software (-N 1, -L 30) (Krueger and Andrews, 2011). Only unique mapping reads
202 were retained for further analysis. Methylated cytosines covered by at least three
203 reads were identified using binomial distribution (p value cutoff 1e-5). Customized
204 Perl scripts were programmed to calculate the CG, CHG and CHH ratio per transcript.

205

206 **miRNA prediction**

207 The clean data of small RNA sequencing (miRNAs and small RNAs, smRNA) were
208 mapped to *G. barbadense* using Bowtie, which allowed 200 multiple mapping
209 positions and zero mismatch for each read. We adopted structure-based annotation
210 and probability-based annotation to predict miRNA loci as suggested by Paterson *et al.*
211 (2012). For the structure-based annotation, RNAfold was employed to predict
212 secondary structures and miRcheck was used to evaluate secondary structures
213 (Jones-Rhoades and Bartel, 2004). We then utilized miRDP to filter the putative
214 precursors of the structure-based annotation (Yang and Li, 2011). All the annotated
215 mature miRNAs were searched against the miRBase (Release 20) to categorize them
216 into cotton conserved and non-conserved miRNA gene families (Kozomara *et al.*,
217 2013). We also employed the CleaveLand pipeline to predict putative miRNA targets
218 based on the degradation data (Addo-Quaye *et al.*, 2009). The *bona fide* miRNA
219 targets were detected based on the criteria suggested by Addo-Quaye *et al.* (2008).

220

221 **Network construction**

222 Weighted gene co-expression analysis (WGCNA) was employed to construct the
223 network (Langfelder *et al.*, 2008). The framework for network construction can be
224 summarized as: 1) defining a gene co-expression similarity by the pearson correlation;
225 2) applying an adjacency function to transform the co-expression similarities to
226 connection strengths with a soft thresholding power of 10; 3) identifying network
227 modules consisting of the highly correlated gene expression patterns using the
228 hierarchical clustering with topological overlap matrix. Non-module genes were
229 categorized by a 'grey' colour. All the steps for network analysis were completed
230 using language R. The software VisANT was used to graphically visualize networks
231 (Hu *et al.*, 2013).

232

233 **Quantitative Real-Time PCR**

234 RNA samples from ovules at -1, 0, 4 and 5 DPA and fibres at 10 DPA and 20 DPA
235 were collected, and quantitative real-time PCR was performed as described previously
236 and the expression levels were normalized using UB7 (Tan *et al.*, 2013). The PCR
237 products at 10 DPA and 20 DPA fibres were cloned into the pGEM-T vector and the
238 randomly selected 100 clones were each sequenced.

239

240 **RLM-RACE**

241 The RLM-RACE was performed to validate the splicing site of miRNA target genes
242 using GeneRacer kit (Invitrogen, <https://www.lifetechnologies.com>). Total RNA (5
243 µg) from 10 DPA and 20 DPA fibres were ligated to RNA adapter without calf
244 intestinal phosphatase treatment. Further PCR reactions using 5' adaptor primers and
245 3' gene-specific primers were guided by the manufacturer's instructions.

246

247 **Data access**

248 The stranded RNA-seq data have been submitted to the NCBI Sequence Read Archive
249 under the Bioproject ID PRJNA266265. The lncRNA sequences and genome
250 coordinate files can be accessed from our genome website at
251 <http://cotton.cropdb.org/cotton/download/data.php>.

252

253 **Results**

254 **Identification and characterization of cotton lncRNAs**

255 In order to develop a comprehensive catalogue of lncRNAs in *Gossypium*, a
256 prerequisite is to integrate a high-quality and high-depth RNA-seq dataset. We
257 collected 154 public and 8 in-house Illumina transcriptomes (Table S1). To determine
258 the orientation of transcripts accurately, we also generated 9 transcriptomes covering
259 different developmental stages of *G. babardense* using the stranded sequencing
260 method (Table S2). In total, this collection represents more than 5 billion clean reads
261 for lncRNA identification.

262 We mapped RNA-seq data from diploids and tetraploids to the subgenomes and
263 the whole genome of *G. babardense* independently (data from our unpublished *G.*

264 *barbadense* genome sequence, 29,751 scaffolds, N50 260.06 kb, encoding 80,876
265 protein-coding genes) in order to perform *de novo* transcript assembly using the
266 Tophat-Cufflinks pipeline. Some filtering steps were conducted to retain *bona fide*
267 lncRNAs (Fig. 1a). This pipeline provided 30,550 lncRNA loci (50,566 transcripts)
268 and 4,718 lncNAT loci (5,826 transcripts).

269 To verify the reliability of prediction, we aligned all the lncRNAs to 425,526
270 public cotton ESTs. A total of 2,929 lncRNAs (5.8%) were supported by at least one
271 EST. We also aligned all the lncRNAs to the collected 454 sequencing reads (Table
272 S3) and observed that 12,029 lncRNAs (23.8%) were supported by at least one read.
273 Attributing lncRNAs to subgenomes showed that the number of lncRNAs in the At
274 subgenome was approximately 2,900 larger than that in the Dt subgenome (Table S4).
275 The exon number distribution of lncRNAs showed that the *G. barbadense* genome
276 encoded 63% single-exonic lncRNAs and 77% single-exonic lncNATs, which are
277 significantly higher proportions than those of protein-coding transcripts (15%; Fig.
278 1b). The mean transcript length of lncRNAs was typically shorter than protein-coding
279 genes (average length: 504 bp for lncRNAs, 713 bp for lncNATs and 1,621 bp for
280 protein-coding transcripts; Fig. 1c).

281 GC content is believed to be related to the biased intergenomic nonreciprocal
282 DNA exchanges in the tetraploid cotton genomes (Guo *et al.*, 2014). In this study, we
283 observed that both the distributions of GC content amongst lncRNAs (lncRNAs and
284 lncNATs) and protein-coding genes exhibit no apparent differences between the At
285 and Dt subgenomes (Kolmogorov-Smirnov test, lncRNA p-value 0.1486,
286 protein-coding genes p-value 0.1803; Fig. 1d). However, lncRNAs show the lowest
287 GC content (median 37.1%), followed by lncNATs (median 40.6%), and
288 protein-coding genes (median 41.8%) the highest in each subgenome.

289 The *G. barbadense* genome is highly enriched for repetitive sequences (70%),
290 with the At subgenome at 74% and the Dt subgenome at 63%. Overlapping
291 coordinates of lncRNAs with transposable elements (TE), we found that 55.8% of
292 lncRNAs contained TE, corresponding to the At subgenome with 58.1%, the Dt
293 subgenome with 54.8% and ungrouped scaffolds with 48.8% (Fig. 1e). The fraction of
294 TE-containing lncNATs was less than half relative to lncRNAs, with At subgenome
295 at 23.2%, Dt subgenome at 21.7% and ungrouped scaffolds at 23.7%. This result is
296 comparable to the studies in animals, such as mouse, zebrafish and human (Kapusta *et*

297 *al.*, 2013). The LTR retrotransposons of the Gypsy family occupied a dominant
298 proportion of repetitive sequences in lincRNAs, which was the same as its distribution
299 at the genome level (Fig. 1f). Long interspersed nuclear elements (LINE) only
300 occupied 6% of the genome, but showed an increased abundance to 14% in lincRNAs,
301 and up to 37% in lincNATs.

302

303 **Expression of cotton lincRNAs among tissues**

304 The stranded RNA-seq data were adopted to systematically explore lincRNA
305 expression among 9 different tissues/samples. The results showed that the highly
306 differentiated tissues anther and cotton fibres at 20 DPA expressed fewer genes than
307 others (Fig. 2a). The overall expression levels of both lincRNAs and lincNATs were
308 lower than of protein-coding transcripts (Fig. 2b), consistent with a previous study
309 (Cabili *et al.*, 2011). Given that lincRNAs may function in regulating adjacent
310 protein-coding genes and thus possess similar expression patterns, we examined this
311 possibility by computing the Pearson correlation coefficients (r_p) between lincRNAs
312 and the nearest protein-coding genes (within 5 kb) (lincRNA-PCgene); lincNATs and
313 the corresponding protein-coding genes on the opposite strand (lincNAT-PCgene); and
314 the nearest protein-coding pairs lacking an intervening gene (PCgene-PCgene). In
315 total, we identified 10,749 lincRNA-PCgene pairs, 5,826 lincNAT-PCgene pairs and
316 25,449 PCgene-PCgene pairs. Compared with randomly sampled transcript pairs, we
317 observed high ratios of extremely positive correlations between lincRNA-PCgene (16%
318 vs. 6%, $r_p > 0.8$), lincNAT-PCgene (35% vs. 4%, $r_p > 0.8$) and PCgene-PCgene (24% vs.
319 6%, $r_p > 0.8$) pairs (Fig. 2c). The expression relationships between these pairs provide
320 candidates to be tested in further functional studies.

321 To evaluate the tissue specificity of expression, the JS scores (an entropy-based
322 measure) of transcripts were calculated (Cabili *et al.*, 2011). The density distributions
323 of lincRNAs and lincNATs were significantly different from protein-coding transcripts
324 (Kolmogorov-Smirnov test, p value $< 2.2e-16$; Fig. 2d). Using a JS score of 0.5 as a
325 cutoff, we found that 42% of lincRNA and 51% of lincNAT transcripts were
326 tissue-preferentially expressed, dramatically higher than the percentage of
327 protein-coding transcripts (18%) across the 9 tissues/samples. Further quantitative
328 analysis showed that anther expressed the largest number of tissue-preferential genes
329 (3,140 protein-coding transcripts, 3,925 lincRNAs and 787 lincNATs) though the total

330 number of expressed transcripts was smaller than for other samples (Fig. 2e). In
331 contrast, fibres at 20 DPA expressed a relatively small number of specific genes (973
332 protein-coding transcripts, 852 lincRNAs and 230 lincNATs), slightly higher than for
333 stigma. Randomly selected tissue-preferentially expressed lincRNAs were verified by
334 RT-PCR (Fig. 2f). These results indicate that a large number of lincRNAs were
335 expressed preferentially in particular tissues.

336

337 **Evolution history and subgenome expression partition**

338 It is believed that the sequences of lincRNAs are less conserved than protein-coding
339 transcripts (Marques and Ponting, 2009; Necsulea *et al.*, 2014), and we were
340 interested to know how many cotton lincRNAs are inherited from closely related
341 species.

342 We firstly aligned the lincRNAs of the At and Dt subgenomes to each reciprocally,
343 then to the diploid A and D genomes, and also to the closely related species
344 *Theobroma cacao* and the more distant dicot *Vitis vinifera* (Jaillon *et al.*, 2007;
345 Argout *et al.*, 2011). Using all the lincRNA transcripts in the At subgenome as queries,
346 we found that 99.5% had homologous copies in the diploid A genome, 76.7 % in the
347 Dt subgenome and 75.6 % in the diploid D genome (Fig. 3a). However, only 6.8% of
348 the lincRNAs in the At subgenome were found to match homologous regions in the *T.*
349 *cacao* genome and 2.6 % in the *V. vinifera* genome. Similar results were observed
350 when lincRNAs in the Dt-subgenome were used as query sequences (Fig. S1a). These
351 results suggest that the vast majority of lincRNAs were species-specific or limited to
352 closely related species.

353 As relatively highly expressed neighbour protein-coding genes may have
354 functional relationships with lincRNAs, we mapped the GO terms of such
355 protein-coding genes ($r_p > 0.9$) to lincRNAs in order to predict their possible functions.
356 The results showed that the At subgenome-specific lincRNAs were enriched in
357 ribosome assembly, spermine biosynthesis process and microtubule cytoskeleton
358 organization (Fig. 3b). Dt subgenome-specific lincRNAs were enriched in lignin
359 catabolic process, response to biotic stimulus and carbon utilization (Fig. 3c). The
360 conserved lincRNAs in *T. cacao* and *V. vinifera* were enriched in fundamental
361 biological processes, such as translation elongation, peroxisome organization and
362 L-phenylalanine catabolism (Fig. S1b).

363 Despite rapid gene fractionation, the majority of lncRNAs were conserved
364 between the At and Dt subgenomes. Using data from the recently released *G.*
365 *arboreum* and *G. raimondii* genomes, we ordered the scaffolds of At and Dt
366 subgenomes to pseudochromosomes based on whole genome alignment (Fig. S2).
367 Through genome-wide synteny analysis, we identified 377 syntenic blocks between
368 the At and Dt subgenomes representing 9,262 protein-coding gene pairs (Fig. 3d).
369 Overlapping lncRNAs with these syntenic blocks and using a reciprocal best hit
370 alignment (coverage cutoff 0.9), we identified 1,090 homoeologous lincRNA pairs
371 between the At and Dt subgenomes, of which 900 pairs were anchored on
372 pseudochromosomes. Genomic landscape analysis showed that both of lncRNAs and
373 protein-coding genes were preferentially located in regions with poor repetitive
374 sequences assuming as a negative correlation (Fig. S3), especially for the
375 protein-coding genes (Fig. 3d).

376 As highlighted in recent studies, the non-additivity of gene expression, also
377 referred as 'transcriptomic shock', appears to be widespread in newly formed
378 allopolyploids (Yoo *et al.*, 2013). Hierarchical clustering of homoeologous lincRNAs
379 showed that those from a total of 8 tissues/samples were clustered in a
380 subgenome-specific manner with the exception of those derived from anther (Fig.
381 S4a), contrasted with the result by clustering protein-coding genes (Fig. S4b). The
382 averaged expressions of lincRNA pairs across tissues were compared (Fig. 3d). This
383 led to the identification of 196 pairs expressed dominantly in At-subgenome and 188
384 pairs expressed dominantly in Dt subgenome. However, the overall comparison
385 ignored the detailed bias in patterns in different tissues, and so we categorized the
386 expression patterns into four types.

387 Based on these analyses, the expression of 305 pairs were At-biased, 315 pairs
388 were Dt-biased and 67 pairs were chimeric-biased. Therefore, we conclude that
389 expression bias of lincRNAs was extensive in tetraploid cottons in a
390 subgenome-specific manner, and the numbers of bias-expressed pairs in each
391 subgenome were comparable.

392

393 **Methylation of lncRNAs**

394 DNA methylation is widespread as a means of regulating protein-coding gene
395 transcription in diverse organisms. To characterize the methylation patterns of

396 lincRNAs, we obtained 4 bisulphate-converted DNA sequencing datasets of petal in
397 cotton species, including diploid *G. arboreum*, *G. raimondii*, an F1-hybrid between *G.*
398 *arboreum* and *G. raimondii*, and the natural tetraploid *G. hirsutum*. The clean reads
399 were uniquely mapped to the *G. barbadense* genome to dissect cytosine methylation
400 (Table S5). The numbers of methylated sites in the At and Dt subgenomes were
401 summarized using each dataset and the percentages of DNA methylation in CG, CHG
402 and CHH contexts were compared (Table S6).

403 At the chromosomal level, highly methylated regions showed preferentially a
404 particular abundance of TEs, seen as a broadly positive correlation. However,
405 protein-coding genes in these regions were expressed at generally low levels (Fig. 4a).
406 This phenomenon was observed in all the four datasets used to analyse diploids and
407 tetraploids. Compared with protein-coding genes, lincRNAs showed higher
408 methylation levels in CG and CHG contexts, but comparable methylation levels in a
409 CHH context (Fig. 4b; Fig. S5). Specifically, the CG methylation levels in exon
410 regions of protein-coding genes rapidly increased when departing from the
411 transcription starting sites and termination sites. However, no such obvious
412 methylation patterns were seen for lincRNAs. For CHG and CHH methylation, the
413 upstream, exon and downstream regions of lincRNAs showed no obvious differences.

414 Many studies have found that the methylation levels of upstream and genic
415 sequences are negatively correlated with the expression levels of protein-coding genes.
416 However, few studies have focused on the relationship between DNA methylation and
417 lincRNA expression. To investigate this, we used RNA-seq data from the same sample
418 as bisulphite-converted DNA sequencing to quantify expression levels of lincRNAs in
419 petals. It was found that in all the three methylation contexts, genes with very high
420 expression levels displayed low methylation levels while highly methylated genes
421 displayed low expression levels, indicating a negative correlation between DNA
422 methylation and gene expression for both of lincRNAs and protein-coding genes (Fig.
423 4c). Specifically, in upstream regions, the scatter-plots of protein-coding genes tended
424 to cover lincRNAs in all three methylation contexts. Interestingly, for gene body
425 methylation, protein-coding genes showed a tighter distribution of methylation levels
426 in each of the three contexts than did lincRNAs. Analysis of accumulated frequency
427 distribution of methylation levels to the relative gene number demonstrated that gene
428 body methylation of lincRNAs in each methylation context was significantly different

429 from that for protein-coding genes, whereas upstream methylation showed no
430 significant differences. These studies suggest that gene body methylation has a
431 generally stronger effect on protein coding gene expression than for lincRNAs.

432 To reveal the direct effects of methylation on lincRNA expression, we collected
433 RNA-seq data from Upland cotton ovules at 0 DPA treated with zebularine, a DNA
434 methylation inhibitor forming a covalent complex with DNA methyltransferases
435 (Zhou *et al.*, 2002). After analysing the quality of RNA-seq (Fig. S6a), we observed
436 that the expression levels of lincRNAs were quite variable and up-regulated
437 expression was clearly consistent along each chromosome after zebularine treatment,
438 while the expression levels of protein-coding genes varied less (Fig. S7).

439 We then conducted a differential gene expression analysis (Fig. 4d). The results
440 showed that a total of 9,917 lincRNA transcripts were differentially expressed, among
441 which the majority (94.4%) were highly expressed in treated ovule samples. In
442 contrast, only 52.2% of differentially expressed protein-coding transcripts were highly
443 expressed in treated samples. Intriguingly, the 86% of up-regulated lincRNAs in the
444 At subgenome and 87% in the Dt subgenome contained repetitive sequences (Fig. 4e),
445 which was a value much higher than for the down-regulated lincRNAs (32% of the At
446 subgenome and 36% of the Dt subgenome) and also higher than ratios of all the
447 lincRNAs in the At and Dt subgenomes (58% of the At subgenome and 55% of the Dt
448 subgenome). Further functional enrichment of the differentially expressed transcripts
449 revealed that up-regulated lincRNAs in treated samples were enriched in DNA
450 integration, cytoskeleton organization, regulation of pH and cell death, while
451 down-regulated lincRNAs were enriched in respiratory gaseous exchange, protein
452 ubiquitination and nucleoside metabolic process (Fig. S6b).

453

454 **Small RNAs generated by lincRNAs**

455 lincRNAs can be small RNA precursors and can also negatively regulate miRNA
456 maturation (Plosky, 2014). We collected 7 sets of small RNA sequencing data for *G.*
457 *barbadense* fibres, representing three important developmental stages (-3 DPA, 0
458 DPA, 3 DPA for fibre initiation stage, 7 DPA and 12 DPA for fibre elongation stage,
459 20 DPA and 25 DPA for fibre secondary cell wall synthesis stage) to identify putative
460 small RNA precursors. The miRNA prediction resulted in a total of 318 conserved
461 miRNAs and 227 non-conserved miRNAs (Table S7, S8). All the lincRNAs were

462 then overlapped to precursors of miRNAs from genome-wide miRNA predictions.
463 We found 128 lincRNAs as possible precursors of conserved miRNAs related to 25
464 families and 101 lincRNAs as possible precursors of non-conserved miRNAs (Table
465 S9). Three well-known miRNAs were covered in this study and presented as
466 examples (Fig. S8). In addition to functioning as miRNA precursors, abundant
467 lincRNA transcripts may be degraded to form smRNAs. The mapping of smRNA
468 reads showed that 4,707 lincRNA transcripts (9.3%) were mapped sense and 4,131
469 (8.2%) were mapped antisense to endo-smRNA reads (Table S9). Future experimental
470 studies are necessary to demonstrate the function of these lincRNAs, but are beyond
471 the scope of the current work.

472

473 **Functional lincRNA candidates in cotton fibre development**

474 Cotton fibre initiation is a fundamental stage determining the fate of the fibre cell.
475 Lint fibres are believed to appear on the day of anthesis (0 DPA) and fuzz fibres
476 develop on the fourth day post anthesis (4 DPA) (Zhang *et al.*, 2007). To identify
477 putative functional lincRNAs contributing to the initiation of lint and fuzz fibres, the
478 expression of 20 randomly selected lincRNAs that were highly expressed in ovules of
479 *G. barbadense* 3-79 was determined in 8 different genotypes of Upland cotton (*G.*
480 *hirsutum*). These cotton accessions include 3 lint-fuzz (TM-1, Xuzhou-142 and YZ1)
481 wild types, 2 lintless-fuzzless mutants (Xuzhou-142 lintless-fuzzless (XZ142WX) and
482 Xinxiangxiaoji lintless-fuzzless (XinWX)) and 3 linted-fuzzless mutants (n2, GZnn
483 and GZNn) (Fig. 5a).

484 Hierarchical clustering analysis showed that most lincRNAs were preferentially
485 expressed in lint-fuzz cotton ovules at -1 and 0 DPA or 4 and 5 DPA (Fig. 5b, c).
486 Specifically, the expression of one lincRNA (LINC02) was highlighted, the expression
487 of which might in part underlie the development of lint and fuzz fibres. This lincRNA
488 produced significantly higher transcription levels in lint-fuzz/linted-fuzzless cottons
489 than that in lintless-fuzzless cottons (p-value < 0.05), but no different transcription
490 levels were seen between lint-fuzz and linted-fuzzless cotton ovules at -1 or 0 DPA
491 ovules (Fig. 5d). We also observed the higher transcription levels in lint-fuzz cottons
492 than that in lintless-fuzzless/linted-fuzzless cottons at 4 DPA or 5 DPA (p-value <
493 0.05) (Fig. 5e).

494 To predict the functional roles of lincRNAs in the 'fibre elongation' and 'secondary
495 cell wall synthesis' stages of fibre development, we applied a weighted gene
496 co-expression network analysis (WGCNA) using published cotton fibre
497 transcriptomes at 10 DPA and 20 DPA (Fig. S9). After removing the low-expressed
498 transcript pairs, 720 lincRNA pairs and 6,858 protein-coding gene pairs were retained
499 for network construction. The network was partitioned into 17 modules (Fig. 6a).
500 Hierarchical clustering and functional enrichment of these modules showed they
501 displayed different characteristics (Fig. S10, S11).

502 The module M12 is highlighted here (Fig. 6c). Transcripts in this module were
503 At-biased in their expression and significantly enriched in heterocyclic metabolic and
504 cofactor metabolic processes (Fig. S10). Hub genes often play founder roles and can
505 define the functional foci in networks (Langfelder *et al.*, 2008). The
506 phosphoenolpyruvate carboxylase-related kinase 2, involved in protein
507 phosphorylation, and a ubiquitin-specific protease were regarded as two hub genes.
508 Interestingly, one lincRNA pair, designated as P1, was highlighted as a hub gene,
509 suggesting a vital functional role in this module (Fig. 6c).

510 Another module, M16, was highlighted as a representative of a Dt-bias
511 expression module (Fig. 6b). This module involved 18 lincRNA pairs and was
512 enriched in oxidation-reduction and small molecule metabolic processes. Previous
513 studies have showed that regulation of reactive oxygen species levels plays a pivotal
514 role in the formation of spinnable cotton fibre (Hovav *et al.*, 2008). Consistent with
515 this, we found that key genes related to reactive oxygen species metabolism, such as
516 2-oxoglutarate (2OG) and Fe (II)-dependent oxygenase, flavin-binding
517 monooxygenase and alpha-helical ferredoxin, were involved in this module. The
518 RabGAP domain-containing protein related to small GTPase mediated signal
519 transduction, categorized as 'small molecule metabolic process', was also involved
520 (Fig. 6d).

521

522 **Integrated expression of lincRNAs generating miR397 and their targets in cotton** 523 **fibre development**

524 Comparative analysis of lincRNAs with small RNA sequencing data helped identify
525 one pair of lincRNAs preferentially expressed in fibres, that were precursors of
526 miR397 from the At and Dt subgenomes (Fig. S12). The Dt-derived lincRNA was

527 highly expressed, and suppressed its At subgenome homoeologue at 10 DPA (Fig. 7a).
528 Conversely, at 20 DPA, the expression of At-subgenome copy reached a very high
529 level, while the expression of Dt-subgenome copy was reduced to a quite low level.
530 This observation was confirmed by the sequencing of 100 randomly picked PCR
531 clones (Fig. 7b). Moreover, the expression level of At-subgenome copy at 20 DPA
532 was significantly higher (~10 fold) than that of the highly expressed Dt-subgenome
533 copy at 10 DPA, which was verified by qRT-PCR detecting the total expression at 10
534 DPA and 20DPA (Fig. 7c).

535 The expression of these two lncRNAs was further analysed in two diploid
536 progenitors and in domesticated and wild tetraploid cottons, using public RNA-seq
537 data. In both diploids, we found the At-subgenome and Dt-subgenome copies were
538 highly expressed in 20 DPA fibres (Fig. S13). We also found that the expression
539 pattern of the At-subgenome copy in all the domesticated and wild Upland and
540 Sea-Island cotton accessions was consistent with the observation in Sea-Island cotton
541 3-79 (Fig. 7d). For the Dt-subgenome copy, we observed the same expression pattern
542 in domesticated Upland and the other 1 Sea-Island cottons, but a reverse expression
543 pattern between 10 DPA and 20 DPA fibres in wild cottons. These results showed that
544 strong directional human selection for enhancing fibre yield has prioritized the
545 expression of the Dt-subgenome copy of lncRNA generating miR397 at 10 DPA, but
546 retained the expression pattern of the At subgenome copy as the same as the diploid A
547 genome and wild tetraploid cottons.

548 MiR397 was validated to target laccase (LAC) transcripts which are important
549 regulators in lignin metabolism (Wang *et al.*, 2012). We detected two types of such
550 LAC genes (*LAC4a* and *LAC4b*; one gene locus in the At subgenome and one locus in
551 the Dt subgenome for each type) in tetraploid cotton genomes (Fig. 7e). RNA-seq
552 data showed that *LAC4a* in the At and Dt subgenomes retained the same expression
553 pattern as diploid progenitors. Nevertheless, the Dt subgenome copy of *LAC4b*
554 underwent an expression transition event the same as the Dt subgenome lncRNA (Fig.
555 7e). *LAC4b* was highly expressed at 20 DPA in *G. raimondii* (proposed Dt
556 subgenome progenitor), which suppressed the expression level at 10 DPA. However,
557 in tetraploid cotton, the Dt subgenome *LAC4b* (Gbscaffold30529.8.0) was highly
558 expressed at 10 DPA and reduced at 20 DPA. These results were validated by
559 qRT-PCR and random clone sequencing analysis (Fig. S14). Degradome sequencing

560 data showed an obvious cleavage activity of miR397 in *LAC4a* (Fig. 7f), indicating
561 that miR397 could repress *LAC4a* by guiding mRNA degradation. In contrast, no
562 cleavage signal was detected in *LAC4b*. Sequence alignment showed a SNP at the
563 tenth site, which was crucial for miRNA-guided mRNA cleavage (Zheng *et al.*, 2012),
564 of miRNA binding region between *LAC4a* and *LAC4b*. The RLM-RACE results
565 confirmed this finding (Fig. 7f).

566 To study the putative mechanisms of expression transition of the Dt subgenome
567 *LAC4b*, we aligned its promoter and downstream regions with the diploid *G.*
568 *raimondii* genome. Intriguingly, little evolutionary variations were observed at the
569 upstream region (3 kb; Fig. 7g). However, an approximate 500 bp transposon inserted
570 into the region downstream of *LAC4b* in the Dt subgenome, which was coming from a
571 region downstream of the At subgenome *LAC4b* and might induce the expression
572 transition (Fig. 7g, h). We confirmed this observation by directly sequencing these
573 two regions from the At and Dt subgenomes.

574

575 **Discussion**

576 Increasing numbers of functional studies on protein-coding genes and small
577 non-coding RNAs are revealing the high level of complexity of eukaryotic
578 transcriptomes, especially when we consider the extensive abundance of long
579 non-coding RNAs (Kapusta and Feschotte, 2014). However, limited data are available
580 for plants. One of the reasons is the poor availability of complete reference genomes
581 and high-depth transcriptome datasets. In cotton, several studies have identified small
582 non-coding RNAs through small RNA sequencing but there are no data presented for
583 lncRNAs. The recent publication of genome sequences and the accumulation of
584 RNA-seq data make it feasible for genome-wide identification of lncRNAs.

585 In this study, we integrated high-quality RNA-seq data with high depth stranded
586 RNA sequencing to explore lncRNAs. We obtained 50,566 lincRNA and 5,826
587 lncNAT transcripts. Due to the tetraploid genomic characteristics and large genome
588 size of cotton, the number of lncRNAs is larger than previous identifications in
589 *Arabidopsis* and maize (Liu *et al.*, 2012; Li *et al.*, 2014). We also believe that more
590 lncRNAs may be identified using stressed plants, as reported for *Arabidopsis* (Liu *et*
591 *al.*, 2012). After attributing these lncRNAs to the At and Dt subgenomes, we observed

592 the number encoded by the At subgenome was 2,900 larger than in the Dt subgenome.
593 Further homoeologous sequence alignments showed that the At subgenome encoded
594 nearly 23% specific lncRNAs (Dt subgenome 17%), which is higher than the ratio of
595 protein-coding genes between these two genomes (Li *et al.*, 2014). When compared
596 with data for the *T. cacao* and *V. vinifera* genomes, we found that lncRNAs diverged
597 quickly among closely related species and even in different genomes of *Gossypium*.
598 Further studies should be conducted to elucidate the functional roles of specific
599 lncRNAs in the At and Dt subgenomes, and those of other species.

600 Genome-wide methylation characterization of protein-coding genes has been
601 explored widely in animals and plants, but few systematic analyses of lncRNAs have
602 been carried out (Zemach *et al.*, 2010). Therefore, we characterized the methylation of
603 lncRNAs using bisulphite-converted DNA sequencing data. It was found that the
604 methylation levels of lncRNAs were higher overall than for protein-coding genes. A
605 large proportion of differentially expressed lincRNAs were up-regulated in ovule
606 samples when treated with methyltransferase inhibitor, and the majority of these
607 lincRNAs overlapped with transposable elements.

608 Furthermore, the genome landscape of averaged gene expression levels in 500 kb
609 windows showed that the expression levels of lncRNAs were more obviously changed
610 compared to protein-coding genes. These results are consistent with the fact that more
611 than half of lncRNAs originated from transposable elements, which are generally
612 heavily methylated (Fedoroff, 2012), indicating that a large number of lincRNAs are
613 silenced in developing cotton ovules due to DNA methylation. These results suggest a
614 functional relationship between transposable elements, lncRNAs and DNA
615 methylation.

616 Functional characterization of lncRNAs is still in its infancy. High-throughput
617 methods, such as Chromatin isolation by RNA purification (ChIRP) and RNA
618 immunoprecipitation (RIP), have proved to be useful and have been utilized in many
619 studies (Chu *et al.*, 2011; Quinn *et al.*, 2014). In this study, we identified several
620 differentially expressed lncRNAs in cotton fibre initiation stage in different cotton
621 accessions, which might be in part associated with the development of lint and fuzz
622 fibres. These lncRNAs represent functional candidates for future experimental studies.
623 We then used a co-expression network strategy to predict function in cotton fibre

624 elongation and secondary cell wall synthesis stages by combining the expression of
625 homoeologous protein-coding genes and lncRNAs across the At and Dt subgenomes.

626 We systematically explored the expression of one lncRNA pair generating
627 miR397. The function of miR397 has been well studied in rice by down-regulating its
628 target laccase-like gene transcripts (Zhang *et al.*, 2013). The target of miR397, *LAC4*,
629 can promote constitutive lignification in *Arabidopsis* (Berthet *et al.*, 2011). In cotton
630 fibres, accumulation of lignin will reinforce the fibre cell walls (Han *et al.*, 2013).
631 Therefore, we focused on the expression of lncRNAs and their target *LAC4* in
632 developing cotton fibres.

633 The expression of two lncRNAs were biased in their subgenomes at different
634 stages, and analysis in diploids and several domesticated and wild tetraploid cottons
635 suggested that human domestication changed the expression pattern of the Dt
636 subgenome lncRNA. Intriguingly, the expression pattern of the Dt subgenome *LAC4b*
637 was also changed in the same manner as for the lncRNA. We speculate that the
638 expression transition of the Dt subgenome *LAC4b* was induced by a TE insertion from
639 the At subgenome. The finding of a SNP in the miRNA binding region between
640 *LAC4a* and *LAC4b* suggests that *LAC4b* might be regulated by miR397 via
641 translational inhibition (Li *et al.*, 2013). Our study provides a framework to explore
642 gene expression bias in tetraploid cotton and the molecular basis of miR397-guided
643 lignin metabolism in fibre development.

644 In summary, our study is the first to characterize lncRNAs in *Gossypium* using
645 high-depth RNA-seq data, although we were able to verify only part of lncRNAs by
646 expression analysis. Future work will aim to dissect their biological functions in
647 relation to cotton development and the genetics underpinning improved agronomic
648 traits. In allopolyploid organisms, such as cotton, wheat and rapeseed, gene
649 expression is to a significant level likely to be regulated by diverse epigenetic
650 modifications (Chen, 2007), and therefore studies on lncRNAs are imperative, as
651 some are most likely involved in epigenetic regulation, such as through chromatin
652 modification and RNA-directed DNA methylation (RdDM). Our study provides new
653 information that underpins the functional characterization of lncRNAs in
654 allopolyploid plants.

655

656 **Acknowledgements**

657 We are very grateful to the laboratory of Dr Joshua A Udall for releasing the
658 bisulphite converted DNA and transcriptome sequencing data in petals. We are also
659 very grateful to the laboratory of Dr Elizabeth S. Dennis for releasing cotton ovule
660 RNA-seq data treated with DNA methyltransferases inhibitor. This work was
661 financially supported by National Natural Science Foundation of China (NO.
662 31230056 and NO. 31201251) and Huazhong Agricultural University Independent
663 Scientific & Technological Innovation Foundation (NO. 2014bs03).

664

665 **References**

- 666 **Addo-Quaye C, Eshoo TW, Bartel DP, Axtell MJ. 2008.** Endogenous siRNA and
667 miRNA targets identified by sequencing of the *Arabidopsis* degradome.
668 *Current Biology* **18**, 758-762.
- 669 **Addo-Quaye C, Miller W, Axtell MJ. 2009.** CleaveLand: a pipeline for using
670 degradome data to find cleaved small RNA targets. *Bioinformatics* **25**,
671 130-131.
- 672 **Anders S, Huber W. 2010.** Differential expression analysis for sequence count data.
673 *Genome Biology* **11**, R106.
- 674 **Argout X, Salse J, Aury J-M, Guiltinan MJ, Droc G, Gouzy J, Allegre M,
675 Chaparro C, Legavre T, Maximova SN et al. 2011.** The genome of
676 *Theobroma cacao*. *Nature Genetics* **43**, 101-108.
- 677 **Berthet S, Demont-Caulet N, Pollet B, Bidzinski P, C ézard L, Le Bris P, Borrega
678 N, Herv é J, Blondet E, Balzergue S, et al. 2011.** Disruption of LACCASE4
679 and 17 results in tissue-specific alterations to lignification of *Arabidopsis*
680 *thaliana* stems. *Plant Cell* **23**, 1124-1137.
- 681 **Cabili MN, Trapnell C, Goff L, Koziol M, Tazon-Vega B, Regev A, Rinn JL.
682 2011.** Integrative annotation of human large intergenic noncoding RNAs
683 reveals global properties and specific subclasses. *Genes Development* **25**,
684 1915-1927.
- 685 **Cech TR, Steitz JA. 2014.** The noncoding RNA revolution-trashing old rules to forge
686 new ones. *Cell* **157**, 77-94.
- 687 **Chen DJ, Yuan CH, Zhang J. et al. 2012.** PlantNATsDB: a comprehensive database
688 of plant natural antisense transcripts. *Nucleic Acids Research* **40**, 1187-1193.
- 689 **Chen ZJ. 2007.** Genetic and epigenetic mechanisms for gene expression and
690 phenotypic variation in plant polyploids. *Annual Review of Plant Biology* **58**,
691 377-406.
- 692 **Chu C, Qu K, Zhong FL, Artandi SE, Chang HY. 2011.** Genomic maps of long
693 noncoding RNA occupancy reveal principles of RNA-chromatin interactions.
694 *Molecular Cell* **44**, 667-678.
- 695 **Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, Guernec G,
696 Martin D, Merkel A, Knowles DG. et al. 2012.** The GENCODE v7 catalog

697 of human long noncoding RNAs: Analysis of their gene structure, evolution,
698 and expression. *Genome Research* **22**, 1775-1789.

699 **Ding JH, Lu Q, Ouyang YD, Mao HL, Zhang PB, Yao JL, Xu CG, Li XH, Xiao**
700 **JH, Zhang QF. 2012.** A long noncoding RNA regulates photoperiod-sensitive
701 male sterility, an essential component of hybrid rice. *Proceedings of the*
702 *National Academy of Sciences, USA* **109**, 2654-2659.

703 **Fedoroff NV. 2012.** Transposable elements, epigenetics, and genome evolution.
704 *Science* **338**, 758-767.

705 **Geisler S, Collier J. 2013.** RNA in unexpected places: long non-coding RNA
706 functions in diverse cellular contexts. *Nature Review Molecular Cell Biology*
707 **14**, 699-712.

708 **Gong L, Kakrana A, Arikiti S, Meyers BC, Wendel JF. 2013.** Composition and
709 expression of conserved microRNA genes in diploid cotton (*Gossypium*)
710 species. *Genome Biology and Evolution* **5**, 2449-2459.

711 **Guan X, Chen ZJ. 2013.** Cotton Fiber Genomics. In Seed Genomics, pp. 203-216.
712 Wiley-Blackwell.

713 **Guan X, Pang M, Nah G, Shi X, Ye W, Stelly DM, Chen ZJ. 2014.** miR828 and
714 miR858 regulate homoeologous MYB2 gene functions in *Arabidopsis*
715 trichome and cotton fibre development. *Nature Communication* **5**, 3050.

716 **Guo H, Wang X, Gundlach H, Mayer KFX, Peterson DG, Scheffler BE, Chee**
717 **PW, Paterson AH. 2014.** Extensive and biased intergenomic nonreciprocal
718 DNA exchanges shaped a nascent polyploid genome, *Gossypium* (Cotton).
719 *Genetics* **197**, 1153-1163.

720 **Han LB, Li YB, Wang HY, Wu XM, Li CL, Luo M, Wu SJ, Kong ZS, Pei Y, Jiao**
721 **GL, et al. 2013.** The dual functions of WLIM1a in cell elongation and
722 secondary wall formation in developing cotton fibers. *Plant Cell* **25**,
723 4421-4438.

724 **Harris RS. 2007.** Improved pairwise alignment of genomic DNA. Ph.D. thesis,
725 Pennsylvania State University.

726 **Hawkins JS, Kim H, Nason JD, Wing RA, Wendel JF. 2006.** Differential
727 lineage-specific amplification of transposable elements is responsible for
728 genome size variation in *Gossypium*. *Genome Research* **16**, 1252-1261.

729 **He XJ, Ma ZY, Liu ZW. 2014.** Non-coding RNA transcription and RNA-directed
730 DNA methylation in *Arabidopsis*. *Molecular Plant* **7**, 1406-1414.

731 **Hovav R, Udall JA, Chaudhary B, Hovav E, Flagel L, Hu G, Wendel JF. 2008.**
732 The evolution of spinnable cotton fiber entailed prolonged development and a
733 novel metabolism. *PLoS Genetics*. **4**, e25.

734 **Hu Z, Chang YC, Wang Y, Huang CL, Liu Y, Tian F, Granger B, DeLisi C. 2013.**
735 VisANT 4.0: Integrative network platform to connect genes, drugs, diseases
736 and therapies. *Nucleic Acids Research* **41**, 225-231.

737 **Jaillon O. et al. 2007.** The grapevine genome sequence suggests ancestral
738 hexaploidization in major angiosperm phyla. *Nature* **449**, 463-467.

739 **Jin J, Liu J, Wang H, Wong L, Chua NH. 2013.** PLncDB: plant long non-coding
740 RNA database. *Bioinformatics* **29**, 1068-1071.

741 **Jones-Rhoades MW, Bartel DP. 2004.** Computational identification of plant
742 microRNAs and their targets, including a stress-induced miRNA. *Molecular*
743 *Cell* **14**, 787-799.

744 **Kapusta A, Kronenberg Z, Lynch VJ, Zhuo X, Ramsay L, Bourque G, Yandell**
745 **M, Feschotte C. 2013.** Transposable elements are major contributors to the
746 origin, diversification, and regulation of Vertebrate long noncoding RNAs.
747 *PLoS Genetics* **9**, e1003470.

748 **Kapusta A, Feschotte C. 2014.** Volatile evolution of long noncoding RNA
749 repertoires: mechanisms and biological implications. *Trends in Genetics* **30**,
750 439-452.

751 **Kim HJ, Triplett BA. 2001.** Cotton fiber growth in planta and in vitro. Models for
752 plant cell elongation and cell wall biogenesis. *Plant Physiology* **127**,
753 1361-1366.

754 **Kozomara A, Griffiths-Jones S. 2013.** miRBase: annotating high confidence
755 microRNAs using deep sequencing data. *Nucleic Acids Research* **42**, 68-73.

756 **Krueger F, Andrews SR. 2011.** Bismark: a flexible aligner and methylation caller
757 for Bisulfite-Seq applications. *Bioinformatics* **27**, 1571-1572.

758 **Langfelder P, Horvath S. 2008.** WGCNA: an R package for weighted correlation
759 network analysis. *BMC bioinformatics* **9**, 559.

760 **Li F, Fan G, Wang K, Sun F, Yuan Y, Song G, Li Q, Ma Z, Lu C, Zou C. et al.**
761 **2014.** Genome sequence of the cultivated cotton *Gossypium arboreum*. *Nature*
762 *Genetics* **46**, 567-572.

763 **Li L, Eichten SR, Shimizu R, Petsch K, Yeh CT, Wu W, Chettoor AM, Givan SA,**
764 **Cole RA, Fowler JE. et al. 2014.** Genome-wide discovery and
765 characterization of maize long non-coding RNAs. *Genome Biology* **15**, R40.

766 **Li S, Liu L, Zhuang X, Yu Y, Liu X, Cui X, Ji L, Pan Z, Cao X, Mo B et al. 2013.**
767 MicroRNAs inhibit the translation of target mRNAs on the Endoplasmic
768 Reticulum in *Arabidopsis*. *Cell* **153**, 562-574.

769 **Liu J, Jung C, Xu J, Wang H, Deng S, Bernad L, Arenas-Huertero C, Chua NH.**
770 **2012.** Genome-wide analysis uncovers regulation of long intergenic
771 noncoding RNAs in *Arabidopsis*. *Plant Cell* **24**, 4333-4345.

772 **Liu N, Tu LL, Tang WX, Gao WH, Lindsey K, Zhang XL. 2014.** Small RNA and
773 degradome profiling reveals a role for miRNAs and their targets in the
774 developing fibers of *Gossypium barbadense*. *Plant Journal* **80**, 331-344.

775 **Marques AC, Ponting CP. 2009.** Catalogues of mammalian long noncoding RNAs:
776 modest conservation and incompleteness. *Genome Biology* **10**, R124.

777 **Matzke MA, Mosher RA. 2014.** RNA-directed DNA methylation: an epigenetic
778 pathway of increasing complexity. *Nature Review Genetics* **15**, 394-408.

779 **Navarro B, Gisel A, Rodio ME, Delgado S, Flores R, Di Serio F. 2012.** Small
780 RNAs containing the pathogenic determinant of a chloroplast-replicating
781 viroid guide the degradation of a host mRNA as predicted by RNA silencing.
782 *Plant Journal* **70**, 991-1003.

783 **Necsulea A, Soumillon M, Warnefors M, Liechti A, Daish T, Zeller U, Baker JC,**
784 **Grützner F, Kaessmann H. 2014.** The evolution of lncRNA repertoires and
785 expression patterns in tetrapods. *Nature* **505**, 635-640.

786 **Paterson AH, Wendel JF, Gundlach H, Guo H, Jenkins J, Jin D, Llewellyn D,**
787 **Showmaker KC, Shu S, Udall J. et al. 2012.** Repeated polyploidization of
788 *Gossypium* genomes and the evolution of spinnable cotton fibres. *Nature* **492**,
789 423-427.

790 **Pauli A, Valen E, Lin MF, Garber M, Vastenhouw NL, Levin JZ, Fan L,**
791 **Sandelin A. et al. 2012.** Systematic identification of long noncoding RNAs
792 expressed during zebrafish embryogenesis. *Genome Research* **22**, 577-591.

793 **Plosky BS. 2014.** An ultraconserved lnc to miRNA processing. *Molecular Cell* **55**,
794 3-4.

795 **Quinn JJ, Ilik IA, Qu K, Georgiev P, Chu C, Akhtar A, Chang HY. 2014.**
796 Revealing long noncoding RNA architecture and functions using
797 domain-specific chromatin isolation by RNA purification. *Nature*
798 *Biotechnology* **32**, 933-940.

799 **Rinn JL, Chang HY. 2012.** Genome regulation by long noncoding RNAs. *Annual*
800 *Review of Biochemistry* **81**, 145-166.

801 **Ruiz-Orera J, Messeguer X, Subirana JA, Alba MM. 2014.** Long non-coding
802 RNAs as a source of new peptides. *eLife* **3**, e03523.

803 **Senchina DS, Alvarez I, Cronn RC, Liu B, Rong J, Noyes RD, Paterson AH,**
804 **Wing RA, Wilkins TA, Wendel JF. 2003.** Rate variation among nuclear
805 genes and the age of polyploidy in *Gossypium*. *Molecular biology and*
806 *evolution* **20**, 633-643.

807 **Sturn A, Quackenbush J, Trajanoski Z. 2002.** Genesis: Cluster analysis of
808 microarray data. *Bioinformatics* **18**, 207-208.

809 **Sun L, Luo H, Bu D, Zhao G, Yu K, Zhang C, Liu Y, Chen R, Zhao Y. 2013.**
810 Utilizing sequence intrinsic composition to classify protein-coding and long
811 non-coding transcripts. *Nucleic Acids Research* **41**, e166.

812 **Swiezewski S, Liu F, Magusin A, Dean C. 2009.** Cold-induced silencing by long
813 antisense transcripts of an *Arabidopsis* Polycomb target. *Nature* **462**, 799-802.

814 **Tan JF, Tu LL, Deng FL, Hu HY, Nie YC, Zhang XL. 2013.** A genetic and
815 metabolic analysis revealed that cotton fiber cell development was retarded by
816 flavonoid naringenin. *Plant Physiology* **162**, 86-95.

817 **Trapnell C, Pachter L, Salzberg SL. 2009.** TopHat: discovering splice junctions
818 with RNA-Seq. *Bioinformatics* **25**, 1105-1111.

819 **Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ,**
820 **Salzberg SL, Wold BJ, Pachter L. 2010.** Transcript assembly and
821 quantification by RNA-Seq reveals unannotated transcripts and isoform
822 switching during cell differentiation. *Nature Biotechnology* **28**, 511-515.

823 **Wang H, Chung PJ, Liu J, Jang IC, Kean MJ, Xu J, Chua NH. 2014.**
824 Genome-wide identification of long noncoding natural antisense transcripts
825 and their responses to light in *Arabidopsis*. *Genome Research* **24**, 444-453.

826 **Wang K, Wang Z, Li F, Ye W, Wang J, Song G, Yue Z, Cong L, Shang H, Zhu S.**
827 *et al.* 2012. The draft genome of a diploid cotton *Gossypium raimondii*.
828 *Nature Genetics* **44**, 1098-1103.

829 **Wang Y, Tang H, Debarry JD, Tan X, Li J, Wang X, Lee TH, Jin H, Marler B,**
830 **Guo H. et al.** 2012. MCScanX: a toolkit for detection and evolutionary
831 analysis of gene synteny and collinearity. *Nucleic Acids Research* **40**, e49.

832 **Wang ZM, Xue W, Dong CJ, Jin LG, Bian SM, Wang C, Wu XY, and Liu JY.**
833 **2012.** A comparative miRNAome analysis reveals seven fiber
834 initiation-related and 36 novel miRNAs in developing cotton ovules.
835 *Molecular Plant* **5**, 889-900.

836 **Wang ZW, Wu Z, Raitskin O, Sun Q, Dean C. 2014.** Antisense-mediated FLC
837 transcriptional repression requires the P-TEFb transcription elongation factor.
838 *Proceedings of the National Academy of Sciences, USA* **111**, 7468-7473.

839 **Wei MM, Wei HL, Wu M. et al. 2013.** Comparative expression profiling of miRNA
840 during anther development in genetic male sterile and wild type cotton. *BMC*
841 *Plant Biology* **13**, 66.

842 **Wendel JF., Brubaker CL, Seelanan T. 2010.** The origin and evolution of
843 *Gossypium*. 1-18.

844 **Xue W, Wang Z, Du M, Liu Y, Liu JY. 2013.** Genome-wide analysis of small
845 RNAs reveals eight fiber elongation-related and 257 novel microRNAs in
846 elongating cotton fiber cells. *BMC Genomics* **14**, 629.

847 **Yang X, Li L. 2011.** miRDeep-P: a computational tool for analyzing the microRNA
848 transcriptome in plants. *Bioinformatics* **27**, 2614-2615.

849 **Yang XY, Wang LC, Yuan DJ, Lindsey K, Zhang XL. 2013.** Small RNA and
850 degradome sequencing reveal complex miRNA regulation during cotton
851 somatic embryogenesis. *Journal of experimental botany* **64**, 1521-1536.

852 **Yoo MJ, Szadkowski E, Wendel JF. 2013.** Homoeolog expression bias and
853 expression level dominance in allopolyploid cotton. *Heredity* **110**, 171-180.

854 **Zemach A, McDaniel IE, Silva P, Zilberman D. 2010.** Genome-wide evolutionary
855 analysis of eukaryotic DNA methylation. *Science* **328**, 916-919.

856 **Zhang DY, Zhang TZ, Sang ZQ, Guo WZ. 2007.** Comparative development of lint
857 and fuzz using different cotton fiber-specific developmental mutants in
858 *Gossypium hirsutum*. *Journal of Integrative Plant Biology* **49**, 1038-1046.

859 **Zhang YC, Yu Y, Wang CY, Li ZY, Liu Q, Xu J, Liao JY, Wang XJ, Qu LH,**
860 **Chen F. *et al.* 2013.** Overexpression of microRNA OsmiR397 improves rice
861 yield by increasing grain size and promoting panicle branching. *Nature*
862 *Biotechnology* **31**, 848-852.

863 **Zheng Y, Li YF, Sunkar R, Zhang W. 2012.** SeqTar: an effective method for
864 identifying microRNA guided cleavage sites from degradome of
865 polyadenylated transcripts in plants. *Nucleic Acids Research* **40**, e28.

866 **Zhou L, Cheng X, Connolly BA, Dickman MJ, Hurd PJ, Hornby DP. 2002.**
867 Zebularine: A novel DNA methylation inhibitor that forms a covalent complex
868 with DNA methyltransferases. *Journal of Molecular Biology* **321**, 591-599.

869 **Zhou X, Sunkar R, Jin H, Zhu JK, Zhang W. 2009.** Genome-wide identification
870 and analysis of small RNAs originated from natural antisense transcripts in
871 *Oryza sativa*. *Genome Research* **19**, 70-78.

872

873 **Figure Legends**

874 **Fig. 1 Identification and characterization of lincRNAs in *G. barbardense*.** (a) The
875 pipeline of long non-coding RNAs (lincRNAs) identification in *G. barbardense*. (b)
876 Exon number distribution per transcript of long intergenic non-coding RNAs
877 (lincRNAs), long non-coding natural antisense transcripts (lincNATs) and
878 protein-coding genes (PCgenes). (c) Length density distributions of lincRNAs,
879 lincNATs and protein-coding transcripts. (d) The GC content of lincRNA, lincNAT
880 and protein-coding transcripts in At (GbAt), Dt (GbDt) subgenomes and ungrouped
881 (GbUn) scaffolds of *G. barbadense* genome. (e) The percentages of lincRNA and
882 lincNAT transcripts overlapped with repetitive sequences in At, Dt subgenomes and
883 ungrouped scaffolds. Transcripts with at least 10 bp overlapping regions with repetitive
884 sequences are counted. (f) The percentage of total length of different repetitive
885 sequences in all the lincRNA and lincNAT transcripts, which were compared with At,
886 Dt subgenomes and ungrouped scaffolds.

887

888 **Fig. 2 Expression of lincRNAs across 9 tissues or developmental stages.** (a) The
889 number of expressed lincRNA and protein-coding transcripts in each tissue or stage.
890 The FPKM cutoff for determining expressed transcripts is 0.1 for lincRNAs and 0.5
891 for protein-coding transcripts. (b) Boxplot shows the distribution of maximum FPKM
892 across samples in lincRNAs, lincNATs and protein-coding transcripts. (c) Pearson
893 correlation coefficient distribution for homoeologous transcript pairs. The
894 lincRNA-PCgene pairs and PCgene-PCgene pairs were restricted to adjacent 5 kb
895 regions. (d) The distributions of maximal tissue specificity scores (JS score)
896 calculated for lincRNA and protein-coding transcripts across all tissues. (e) Venn
897 diagram shows the numbers of tissue-preferentially expressed transcripts in each
898 tissues. The cutoff of maximum JS score per transcript is 0.5. (f) RT-PCR validation
899 of tissue-preferentially expressed lincRNAs (LINC1 to LINC9).

900

901 **Fig. 3 Evolution history and genomic landscape of lincRNAs.** The homoeologous
902 chromosomes are in the same colour. The grey lines show syntenic blocks and
903 coloured lines show homoeologous lincRNA pairs between At and Dt subgenomes. (a)
904 Pie chart showing the proportions of homologous lincRNAs in closely related species.

905 All the At subgenome lincRNAs in *G. barbadense* are aligned to Dt subgenome, *G.*
906 *raimondii*, *G. arboretum*, *T. cacao* and *V. vinifera*. (b) GO enrichment of At
907 subgenome specific lincRNAs. (c) GO enrichment of Dt subgenome specific
908 lincRNAs. (d) Features of lincRNAs in At (green track) and Dt (red track) subgenomes
909 of *G. barbadense*, (a) ratio of GC content in 500 kb windows, (b) percentage of
910 repetitive sequences in 500 kb windows, (c) number of protein-coding genes in 500
911 kb windows, (d) number of lincRNA loci in 500 kb windows, (e) log₂ ratio of
912 averaged FPKM values for homoeologous lincRNA pairs ($\log_2(\text{At/Dt}) \geq 1$). The red
913 dots show At-biased expression, green dots show Dt-biased expression and grey dots
914 show equivalent expression. The right panel shows the categories of biased expression
915 of homoeologous lincRNA pairs. The grey dashed lines shows the cutoff
916 ($\log_2(\text{At/Dt}) \geq 1$ or $\log_2(\text{At/Dt}) \leq -1$) for determining biased expression.

917

918 **Fig. 4 Characterization of lincRNA methylation.** (a) The DNA methylation and
919 gene expression levels (lincRNAs and protein-coding genes) in *G. barbadense* (At
920 subgenome green track, Dt subgenome red track). The homoeologous chromosomes
921 are represented by the same color. Each chromosome is divided into 500 kb windows.
922 The four track groups represent *G. arboretum* (a), *G. raimondii* (b), F1-hybrid
923 between *G. arboretum* and *G. raimondii* (A2 x D5) (d) and natural tetraploid (e). For
924 each track group, the CG methylation level, CHG methylation level, CHH
925 methylation level, averaged lincRNA expression and averaged protein-coding gene
926 expression are depicted outside-to-inside. The track c shows the TE density along
927 each chromosomes. (b) DNA methylation in lincRNA and protein-coding gene
928 regions. For each gene, the up-stream 1 kb, gene body and down-stream 1 kb are
929 characterized and divided into 50 bins, respectively. (c) Correlations of the DNA
930 methylation in CG, CHG and CHH contexts with gene expression. For each
931 methylation context, the averaged DNA methylation levels of up-stream 1kb and gene
932 body were plotted against the gene expression level. The accumulated frequency
933 distribution of transcript numbers against DNA methylation level of lincRNAs and
934 protein-coding genes are compared on the upper-right corner. The significant levels (p
935 value) of distribution divergence are indicated. (d) Scatter-plot shows the
936 differentially expressed lincRNAs and protein-coding genes between
937 zebularine-treated ovule and controls. (e) The proportions of TE-contained

938 up-regulated and down-regulated lincRNAs after treated with zebularine are
939 compared to that of all the lincRNAs in At and Dt subgenomes.

940

941 **Fig. 5 Identification of lincRNAs associated with cotton fibre initiation.** (a) The
942 mature fibres or naked seeds of eight Upland cottons used in this study, including
943 three lint-fuzz wild-type genotypes (TM-1, YZ1, XZ142), two lintless-fuzzless mutant
944 genotypes (XZ142WX, XinWX) and three linted-fuzzless mutant genotypes (n2,
945 GZnn, GZNn). (b, c) Heatmaps show the real-time PCR validation of expression of 20
946 lincRNAs at -1 DPA and 0 DPA ovules (b) and 4 DPA and 5 DPA ovules (c). The
947 relative expression levels of each gene in different samples were normalized in the
948 same data interval (-2 to 2) and visualized using Genesis (Sturn *et al.*, 2002). (d)
949 Real-time PCR validation of the differential expression of one lincRNA (LINC02)
950 between lint-fuzz/linted-fuzzless cottons and lintless-fuzzless cottons at -1 and 0 DPA
951 ovules (p-value < 0.05). (e) Real-time PCR validation of the differential expression of
952 one lincRNA (LINC02) between lint-fuzz cottons and lintless-fuzzless/linted-fuzzless
953 cottons at 4 and 5 DPA ovules (p-value < 0.05).

954

955 **Fig. 6 Functional implications of lincRNAs in cotton fibre elongation and**
956 **transition to secondary cell wall synthesis stages.** (a) Clustering dendrogram of
957 homeologous gene duplets between At and Dt subgenomes and assigned modules
958 (labeling M1 to M17). These modules are constructed using gene expression data
959 from 10 DPA and 20 DPA cotton fibre transcriptomes. (b) Heatmaps of gene pairs
960 expression in M12 (left) and M16 (right) combined with the normalized expression of
961 hub genes. (c) Module network of M12. The lincRNA pairs and their involved
962 co-expression relationships with protein-coding genes are colored in red. The
963 protein-coding genes significantly enriched in organic cyclic compound metabolic
964 process are colored in green and orthologs in *Arabidopsis* are annotated. (d) Module
965 network of M16. The lincRNA pairs and their involved co-expression relationships
966 with protein-coding genes are colored as M12. The protein-coding genes significantly
967 enriched in oxidation-reduction process are colored in blue and small molecule
968 metabolic process in cyan.

969

970 **Fig. 7 Expression and functional analysis of lncRNAs generating miR397.** (a)
971 RNA-seq mapping of the lncRNAs pair generating miR397. The mature sequences of
972 miR397 are labeled in red boxes. (b) Ratio of the clone sequences in the At and Dt
973 subgenomes at 10 DPA and 20 DPA. (c) Real-time PCR of the total expression of
974 lncRNA pair in the At and Dt subgenomes. Error bars show three biological replicates.
975 (d) Comparison of the normalized expression of lncRNA pair in domesticated and
976 wild *G. hirsutum* and *G. barbadense* accessions by RNA-seq. (e) Phylogenetic tree of
977 *LAC4* in diploid A and D genomes, and the At and Dt subgenomes of *G. barbadense*.
978 The *Arabidopsis LAC4* is regarded as an outgroup. Light red symbols show genes in
979 diploid A genome (triangle) and the At subgenome (diamond), and light green
980 symbols show genes in diploid D genome (square) and the Dt subgenome (round).
981 The expression of each gene at 10 DPA and 20 DPA in diploid/tetraploid cottons is
982 indicated. (f) Degradome sequencing shows the signature abundance in the position of
983 *LAC4* (left *LAC4a*, right *LAC4b*) targeted by miR397. The red dot shows significant
984 signature as indicated by red arrow. The target cleavage site is identified through
985 RLM-RACE, as shown below the target plot. The numbers indicate the cleavage
986 frequency through clone sequencing. (g) Sequence alignment of the upstream (left)
987 and downstream (right) 3k regions between the Dt subgenome and diploid D genome
988 by LASTZ software. (h) Model of the TE insertion from the At subgenome to the Dt
989 subgenome in *G. barbadense*. TSS, transcription start site; TTS, transcription
990 termination site.
991