

# Identification of a genome-specific repetitive element in the *Gossypium* D genome

Hejun Lu<sup>1,2,\*</sup>, Xinglei Cui<sup>2,\*</sup>, Yanyan Zhao<sup>2</sup>, Richard Odongo Magwanga<sup>2,3</sup>, Pengcheng Li<sup>2</sup>, Xiaoyan Cai<sup>2</sup>, Zhongli Zhou<sup>2</sup>, Xingxing Wang<sup>2</sup>, Yuling Liu<sup>4</sup>, Yanchao Xu<sup>2</sup>, Yuqing Hou<sup>2</sup>, Renhai Peng<sup>4</sup>, Kunbo Wang<sup>2,5</sup> and Fang Liu<sup>2</sup>

<sup>1</sup>Gembloux Agro-Bio Tech, University of Liège, Gembloux, Namur, Belgium

<sup>2</sup>Research Base of Tarium University, State Key Laboratory of Cotton Biology, Institute of Cotton Research of Chinese Academy of Agricultural Science, Anyang, Henan, China

<sup>3</sup>School of Biological and Physical Sciences (SBPS), Jaramogi Oginga Odinga University of Science and Technology (JOOUST), Bondo-Kenya, Bondo, Kenya

<sup>4</sup>Anyang Institute of Technology, Anyang, Henan, China

<sup>5</sup>Tarium University, Alar, Xinjiang, China

\*These authors contributed equally to this work.

## ABSTRACT

The activity of genome-specific repetitive sequences is the main cause of genome variation between *Gossypium* A and D genomes. Through comparative analysis of the two genomes, we retrieved a repetitive element termed *ICRd* motif, which appears frequently in the diploid *Gossypium raimondii* (D<sub>5</sub>) genome but rarely in the diploid *Gossypium arboreum* (A<sub>2</sub>) genome. We further explored the existence of the *ICRd* motif in chromosomes of *G. raimondii*, *G. arboreum*, and two tetraploid (AADD) cotton species, *Gossypium hirsutum* and *Gossypium barbadense*, by fluorescence *in situ* hybridization (FISH), and observed that the *ICRd* motif exists in the D<sub>5</sub> and D-subgenomes but not in the A<sub>2</sub> and A-subgenomes. The *ICRd* motif comprises two components, a variable tandem repeat (TR) region and a conservative sequence (CS). The two constituents each have hundreds of repeats that evenly distribute across 13 chromosomes of the D<sub>5</sub> genome. The *ICRd* motif (and its repeats) was revealed as the common conservative region harbored by ancient Long Terminal Repeat Retrotransposons. Identification and investigation of the *ICRd* motif promotes the study of A and D genome differences, facilitates research on *Gossypium* genome evolution, and provides assistance to subgenome identification and genome assembling.

Submitted 14 June 2019

Accepted 4 December 2019

Published 3 January 2020

Corresponding authors

Kunbo Wang, [wkbcricri@163.com](mailto:wkbcricri@163.com)

Fang Liu, [liufcri@163.com](mailto:liufcri@163.com)

Academic editor

Gerard Lazo

Additional Information and

Declarations can be found on

page 12

DOI 10.7717/peerj.8344

© Copyright

2020 Lu et al.

Distributed under

Creative Commons CC-BY 4.0

OPEN ACCESS

**Subjects** Agricultural Science, Evolutionary Studies, Genetics, Genomics, Plant Science

**Keywords** *Gossypium*, Repetitive element, D genome, Fluorescence in situ hybridization (FISH), Genome-specific, Evolution

## INTRODUCTION

Repetitive DNA sequences are common in eukaryotic genomes, and account for a large fraction of the host genome (*Ibarra-Laclette et al., 2013*). Their amount is highly correlated with the size of the host genome (*Feschotte, 2008*). Repetitive DNA is divided into two major groups based on their structures: tandem repeats and interspersed repeats (*Jurka et*

*al.*, 2005). Tandem repeats are known as simple sequence repeat (SSR), and include micro-satellites, mini-satellites, and satellites (Jeffreys, Neumann & Wilson, 1990). Interspersed repeats are also referred to as transposable elements (TEs) or transposons.

After the first TE of Ac/Ds was reported in maize (McClintock, 1950; Brink & Williams, 1973; Goldschmidt, 2002), further TEs have been identified in many eukaryotic species (Munoz-Lopez & Garcia-Perez, 2010). There are thousands of different TE families in plants, which display extreme diversity (Sanmiguel & Bennetzen, 1998; Bennetzen, 2005; Morgante, 2006). Finnegan first proposed a TE classification system, which includes two classes based on their transposition mechanisms, viz., those mediated by RNA (Retrotransposons) and those by DNA (DNA transposons) (Bowen & Jordan, 2002; Wessler, 2006; Arkhipova, 2018). Wicker unified TEs nomenclature and classification by applying mechanistic and enzymatic criteria (Wicker *et al.*, 2007). TEs play important roles in the genome through diverse ways, such as variation in intron size (Deutsch & Long, 1999; Zhang *et al.*, 2011; Koonin, Csuros & Rogozin, 2013), segmental duplication (Del Pozo & Ramirez-Parra, 2015), transfer of organelle DNA to the nucleus (Adams & Palmer, 2003), expansion/contraction of tandem repeats, and illegitimate recombination (Finnegan, 1989; Koike, Nakai & Takagi, 2002). Long Terminal Repeat Retrotransposons (LTR-TEs), which are usually scattered throughout genomes, are the most abundant TE type and can cause genome expansion over a short evolutionary period particularly in plants (Piegu *et al.*, 2006). The investigation of genome-specific TE is an efficient approach to studying species formation and genome evolution (Dong *et al.*, 2018).

*Gossypium*, a genus of flowering plants from which cotton is harvested, diverged from the common ancestor with *Theobroma cacao* approximately 33.7 million years ago (MYA) (Wang *et al.*, 2012). *Gossypium* comprises eight diploid ( $2n = 2x = 26$ ) genomic groups: A, B, C, D, E, F, G, K, and one allotetraploid ( $2n = 4x = 52$ ) genomic group: AD (Wang, Wendel & Hua, 2018). *Gossypium* species are good subjects for research on polyploidization, genomic organization and genome-size variation because of their high genome diversity: from the smallest New World D genome with an average of 885 Mb to the Australian K-genome with an average of 2,576 Mb (Hendrix & Stewart, 2005). The accumulation of different lineage-specific TEs was thought to be responsible for the variation in genome size in *Gossypium* genomic groups (Hawkins *et al.*, 2006; Lu *et al.*, 2018b). Of the eight genomic groups, the A and D groups are the main ones investigated in cotton genomics research (Du *et al.*, 2018). *Gossypium hirsutum*, the major cultivated cotton species, is known to have originated from the progenitors of *G. arboreum* (A<sub>2</sub>) and *G. raimondii* (D<sub>5</sub>) (Paterson *et al.*, 2012). The key phenotype difference between *G. arboreum* and *G. raimondii* is the production of spinnable fibers in the former but not the latter. In terms of the genomics, the former has a genome size of 1,746 Mb/1C, which is about two times that of the latter (885 Mb/1C) (Hendrix & Stewart, 2005). Genome sequencing showed that the difference in the numbers of protein-coding genes between the A (41,330) and D (37,505) genomes is not obvious, while the lineage-specific TE content is the main reason for the size gap between the A and D genome (Li *et al.*, 2015; Du *et al.*, 2018). Moreover, Wang, Huang & Zhu (2016) suggested that the transposable elements play an important role during cotton genome evolution and fiber cell development. Thus, research on the lineage-specific

**Table 1** Plant materials used in this work, together with ploidy, studied genome, and specimen accession code.

Species	Ploidy	Genome	Accession
<i>G. arboreum</i>	2x	A <sub>2</sub>	Shixiya-1
<i>G. raimondii</i>	2x	D <sub>5</sub>	D5-07
<i>G. hirsutum</i>	4x	(AD) <sub>1</sub>	CCRI-12
<i>G. barbadense</i>	4x	(AD) <sub>2</sub>	Xinhai-7

repetitive sequences between A and D genomes is a meaningful path to investigate speciation dynamics.

Fluorescence *in situ* hybridization (FISH) is a versatile tool to visualize the distribution of certain DNA sequences in chromosomes and plays a vital role in cytogenetic research. In tetraploid cotton, FISH has played a key role in cytological experiments that have contributed to the understanding the allotetraploid event. FISH with DNA segments harboring dispersed repeats has identified genome-specific repeats between the A and D genome, and showed that some A genome repeats appear to have spread to the D genome (*Hanson et al., 1998; Zhao et al., 1998a; Zhao et al., 1998b*). Although the repetitive DNA fragments are more common in the A than in the D genome, one tandem repeat family (B77) has been well-characterized from the D Chromosome (*Zhao et al., 1998a; Zhao et al., 1998b*). Recently, more repetitive sequences were observed with FISH in the cotton genome after construction of a cotton cytogenetic map *Cui et al., 2015; Liu et al. 2016. Lu et al., 2018*) suggested that *CICR* was an important contributor to the size gap between the A and D genome. The identification and localization of these repetitive sequences benefit genome assembly and facilitate understanding of the mechanism of genome evolution.

The D genomic group represents a diverse group of diploids that diverged from a branch of A, B, C, E, F, G, and K genomic groups about 5–10 MYA (*Senchina et al., 2003*). Although the D genome has the smallest size of all *Gossypium* species, this study has revealed the presence of a set of repeat elements with high proliferation, which is absent in the A genome. The discovery and characterization of these novel repetitive elements provides components for a repetitive sequences database and new insight into *Gossypium* evolution.

## MATERIALS AND METHODS

### Plant materials

Cotton plants were obtained from the National Wild Cotton Nursery in Hainan Island, China, sponsored by the Institute of Cotton Research of Chinese Academy of Agricultural Sciences (ICR-CAAS). They were also conserved in the greenhouse at ICR-CAAS' headquarters in Anyang City, Henan Province, China. The DNA and cells came from specimens listed in [Table 1](#), which is based on the latest nomenclature of *Gossypium* species (*Wang, Wendel & Hua, 2018*).

The repeat elements were characterized in the *G. raimondii* genome (*Paterson et al., 2012*), and compared to genomes in other *Gossypium* genomes, viz., *G. arboreum* (*Li et al., 2014*), *G. hirsutum* (AD)<sub>1</sub> (BGI (*Li et al., 2015*), NBI (*Zhang et al., 2015*), HAU (*Wang et*

*al.*, 2019), ZJU (*Hu et al.*, 2019)), *G. barbadense* (AD)<sub>2</sub> (HAUv1 (*Yuan et al.*, 2015), CAS (*Liu et al.*, 2015), HAUv2 (*Wang et al.*, 2019), and ZJU (*Hu et al.*, 2019)). All genome data was downloaded from Cottongen (<https://www.cottongen.org/>), except the (AD)<sub>2</sub>-CAS which was obtained from GenBank under PRJNA251673.

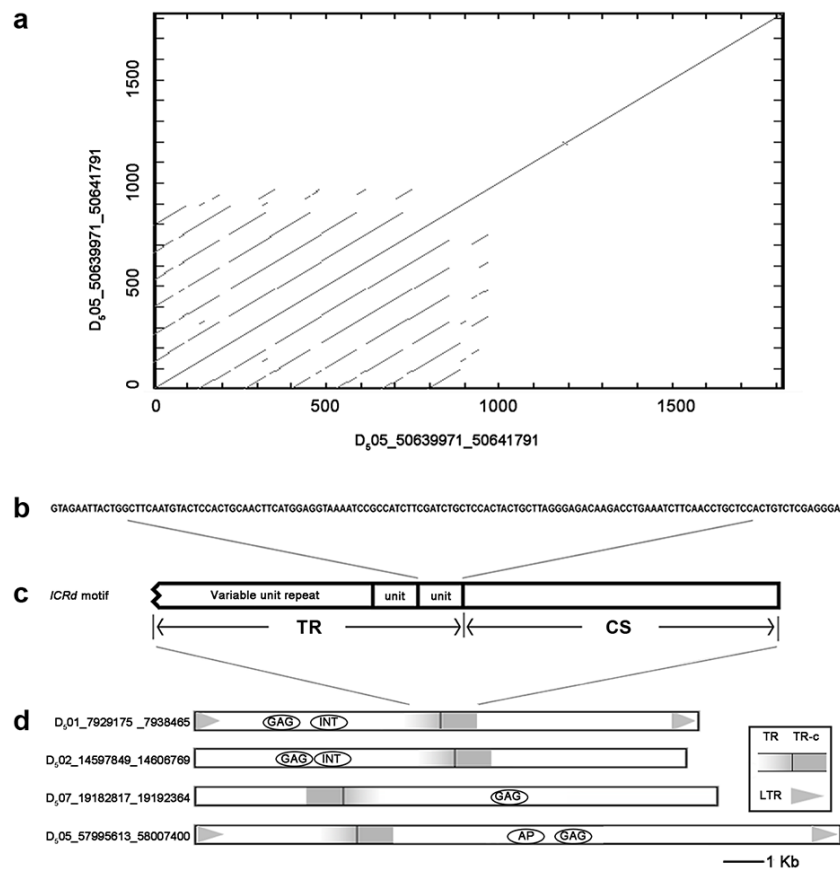
### Characterization of the repetitive element and bioinformatics analysis

BLASTN (v2.6.0) (*Camacho et al.*, 2009) was used to identify repeat elements in the genomes of the plant material, and in the genomes stored in the databases. We used a threshold of greater than or equal to 80% matching ratio and an 80% similarity following the 80–80 rule suggested by *Wicker et al.* (2007). The tandem repeats (TRs) were identified with Tandem Repeats Finder (v4.09) (*Benson*, 1999). We used Perl script for batch extracting sequences from the genome (Doc S1). Sequence alignments were obtained from MUSCLE (v3.81) (*Edgar*, 2004). The Unipro UGENE (v1.31) was used to present the alignments and train consensus sequences (*Okonechnikov et al.*, 2012). The inner enzyme annotation was obtained by online CD-search in NCBI (*Marchler-Bauer et al.*, 2017). GIRI Repbase (*Chen et al.*, 2007) were queried for annotation. RepeatMasker (v4.07) was used to annotate the insertions and estimate the proportion of repetitive sequences in genomes (<http://www.repeatmasker.org>).

Flanking LTRs of LTR-TEs were identified with LTRharvest (v1.5.8) (*Ellinghaus, Kurtz & Willhoeft*, 2008). Subsequently, Vmatch (v2.3.0) was used to cluster the LTRs (*Kurtz*, 2003). The divergence time of the LTR-TEs was estimated using the formula  $T = d/2r$ , where  $r$  represents a substitution rate of  $1.3 \times 10^{-8}$  per site per year (*Ma & Bennetzen*, 2004), and  $d$  represents the distances of paired LTRs, which was calculated based on the Kimura two-parameter (*Kimura*, 1980). The insertions of the repeat elements were obtained based on the BLASTN result, and the LTR-TE and Coding-sequence (CDS) information was obtained from genome annotation (*Paterson et al.*, 2012), which were illustrated by the ggplot2 R package (*Wickham*, 2016) with a sliding 500 kb window for LTR-TE and CDS. The synteny blocks of the homologous segments were shown by a Perl script (Doc S1) based on the BLASTN results.

### Fluorescence in situ hybridization (FISH)

A probe was designed with the PCR product of the *ICRd* motif, which was obtained from the forward primer: TTCTATTTTATCCATCGTTATG, reverse: GGAGATAGGATTTGTTGCT; and followed the amplification procedure: firstly, 95 °C for 5 min of pre-degeneration; then, 30 cycles at 95 °C for 30 s, 52 °C for 30 s, and 72 °C for 2 min. The final extension was done at 72 °C for 6 min. Composition of the reaction mix used the following: gDNA (~5 µg/ml), primers (~0.8 µM), PCR Master Mix (Thermo), and H<sub>2</sub>O. The gDNA was extracted from the leaves of the cotton plants (Table 1). The probe was purified and labeled with digoxigenin-dUTP via nick translation, according to manufacturer's instructions (Roche Diagnostics, USA). Mitotic chromosome preparation and FISH procedures were conducted using a modified protocol (*Wang et al.*, 2001).



**Figure 1** The structure of *ICRd* motif. (A) The self-blast of the *ICRd* motif showed the inner repeats; (B) the structure of *ICRd* motif; (C) the basic TR unit; (D) the examples of the structure illustration of the LTR-TEs inserted with *ICRd* motif.

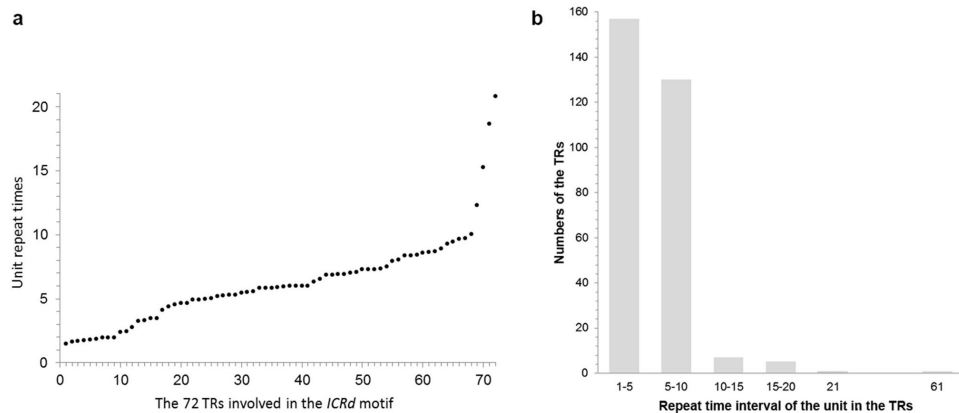
Full-size DOI: 10.7717/peerj.8344/fig-1

## RESULTS

### One specific repetitive sequence in the *Gossypium* D<sub>5</sub> genome

We performed BLAST to query all of the interspersed repetitive sequences of *G. raimondii* (Paterson *et al.*, 2012) with the whole genome of *G. arboreum* (A<sub>2</sub>) (Li *et al.*, 2014). One segment in the *G. raimondii* (D<sub>5</sub>) genome (Chr05: 50639971-50641791) was filtered out and recognized as D<sub>5</sub> genome-specific. This sequence repeats frequently and is spread over 13 chromosomes of the D<sub>5</sub> genome (Table S1), while it is absent from the A<sub>2</sub> genome. Searches in Repbase (Chen *et al.*, 2007) and NCBI found no related annotation and LTRharvest (Ellinghaus, Kurtz & Willhoeft, 2008) and a CD-search (Marchler-Bauer *et al.*, 2017) revealed it is neither LTR nor a coding sequence.

Manual inspection revealed the structure of the genome-specific sequence as having two constituents: a tandem repeats array (referred as TR hereafter) composed of 133 bp basic units, and an unknown conservative sequence (referred as CS hereafter) (Fig. 1). Based on this feature, we identified 72 sequences in total from the D<sub>5</sub> genome with RepeatMasker (Table S2), referred to here as the *ICRd* motif following our previous work (Lu *et al.*,



**Figure 2** The content of the basic unit in the TRs. (A) The basic unit content in the TRs involved in the *ICRd* motifs, displayed from small to large; (B) the number of *ICRd* TRs that harboring different unit content, the *x*-axis adopt the intervals of unit content for convenient exhibition.

Full-size [DOI: 10.7717/peerj.8344/fig-2](https://doi.org/10.7717/peerj.8344/fig-2)

2018b). Among the 72 *ICRd* motifs, the TRs are length-variable having 2–20 times of basic units (Fig. 2A), while the CSs are stable and have an average size  $\sim$ 860 bp.

To verify the genome specificity and chromosome distribution of the *ICRd* motif, we used the PCR product of the *ICRd* motif from *G. raimondii* to design the probe for FISH on the mitotic chromosomes of diploid  $A_2$  and  $D_5$ , and tetraploid *G. hirsutum* ( $(AD)_1$ ) and *G. barbadense* ( $(AD)_2$ ). The probe generated bright signals covering all the chromosomes of the  $D_5$  and  $D$ -subgenome, but no signals on the  $A_2$  and  $A$ -subgenome (Fig. 3). These cytogenetic inspections were in accordance with the genomic comparative analysis and further revealed that the *ICRd* motif is a genome-specific and highly repetitive element in the  $D_5$  genome, as well as in the  $D$ -subgenome of tetraploid cotton.

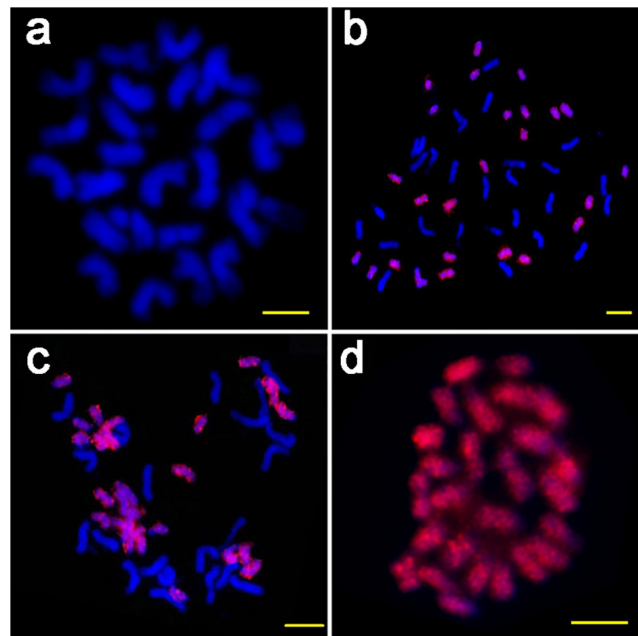
### LTR-TEs inserted with the *ICRd* motif

We compared the insertion loci of 72 *ICRd* motifs with the whole genome repeats annotation (gff file) of the  $D_5$  genome (Paterson et al., 2012) and found that each of the motifs is one-to-one harbored within the 72 LTR-TEs (Table S3), which meant the former is the inner part of the latter.

We extracted the 72 LTR-TEs sequences from the  $D_5$  genome and parsed their structure, which showed all sequences are incomplete, lacking either enzyme or flanking LTRs, the required elements for an intact LTR-TE (Wicker et al., 2007). A consensus accumulation histogram obtained from aligning all of these LTR-TEs together (Fig. S1) showed these TEs to have a vast sequence variation and a single conservative region representing the insertion region of the *ICRd* motif (Fig. 4). The *ICRd* motif appears to be more stable than other parts of the TEs along with degradation and evolution. This stability implies the importance of *ICRd* motif to the TEs.

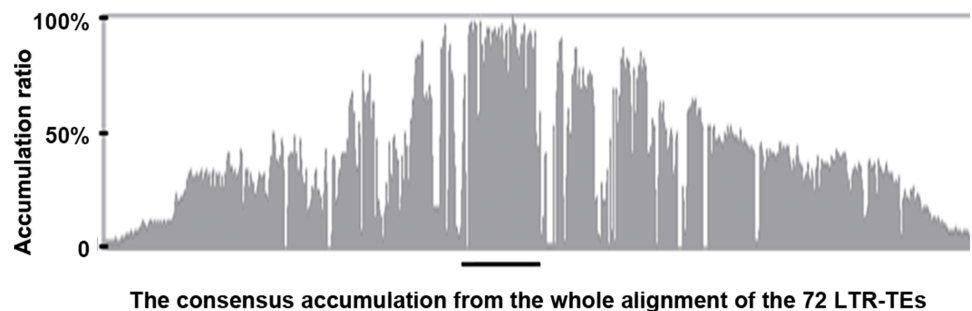
Of the 72 LTR-TEs, 25 were identified as having paired flanking LTRs, and were used to represent the classification and evolution of these TEs. The LTR cluster results showed that, except for two TEs having similar LTR regions, the other 23 TEs are totally different





**Figure 3** The FISH images of *ICRd* motif (red) hybridized to mitotic chromosomes of four species. (A) *G. arboreum* (AA); (B) *G. hirsutum* (AADD); (C) *G. barbadense* (AADD); (D) *G. raimondii* (DD). Bar = 5  $\mu$ m.

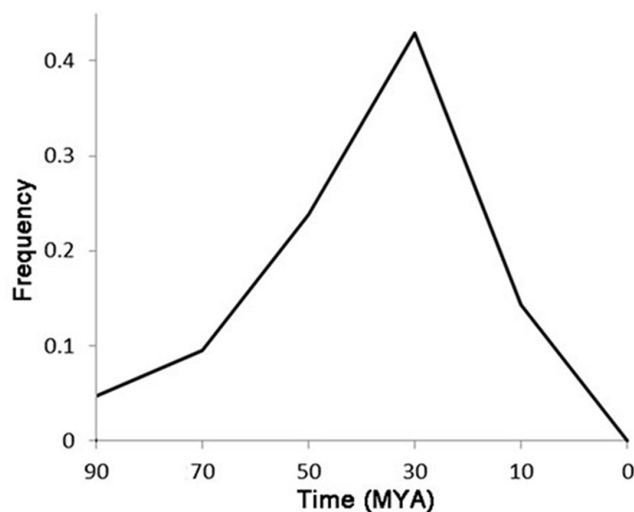
Full-size DOI: 10.7717/peerj.8344/fig-3



**Figure 4** The consensus accumulation histogram from the whole alignment of the 72 LTR-TEs. The region marked with the black line is the *ICRd* motif region.

Full-size DOI: 10.7717/peerj.8344/fig-4

from each other, indicating that they do not belong to the same family based on the LTR similarity rules (Wicker *et al.*, 2007). The estimated active date curve of these TEs—almost all prior to 10 MYA and peaking at  $\sim$ 30 MYA (Fig. 5)—shows the peak is close to the time that *G. raimondii* and *T. cacao* diverged approximately 33.7 MYA (Wang *et al.*, 2012), far earlier than the putative divergence time of the *Gossypium* A and D genomes (Wendel & Cronn, 2001). These results indicate that these LTR-TEs are ancient and potentially contributed to speciation of *Gossypium*.



**Figure 5** The accumulation of putative active date of the LTR-TEs.

Full-size DOI: [10.7717/peerj.8344/fig-5](https://doi.org/10.7717/peerj.8344/fig-5)

### Abundant constituents of the *ICRd* motif in the $D_5$ genome

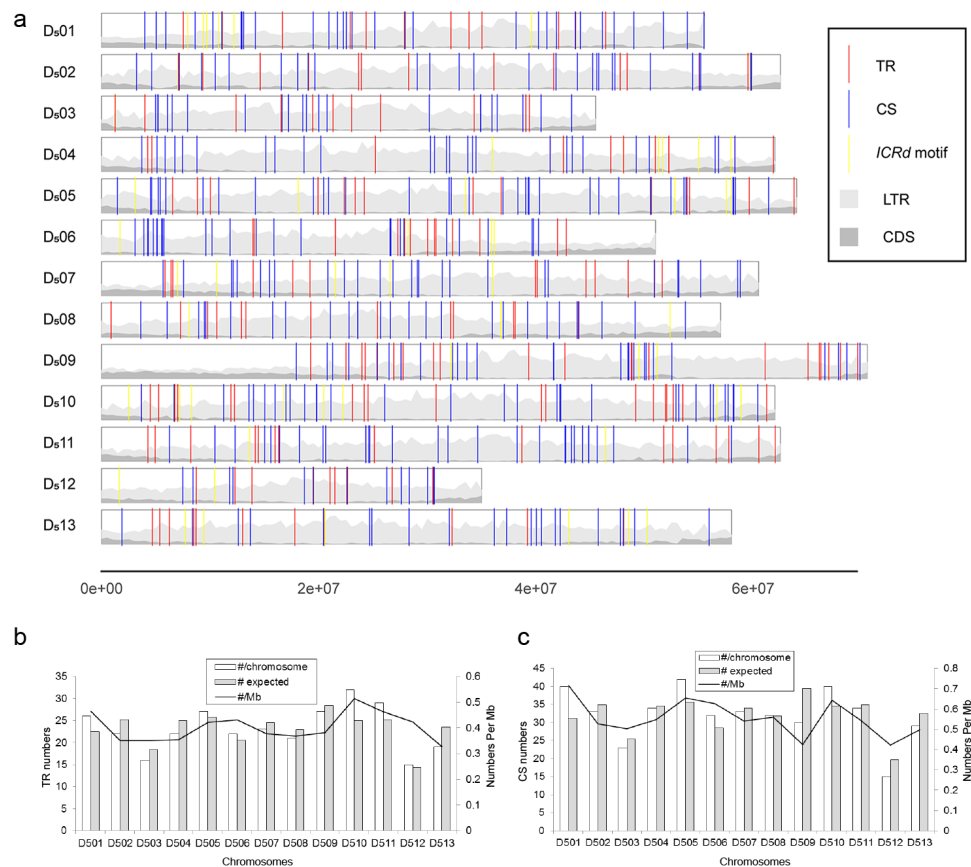
To further analyze the genomic features of the *ICRd* motif, we separately investigated the content and distribution of its two constituents (TR and CS) in the  $D_5$  genome (Fig. 6A). In total 350 TR insertions were detected (Table S2). Insertions varied in length (due to the unit repeating at different times) between 2–21, but mainly 2–10 times the basic unit length (Fig. 2B). The longest TR insertion in  $D_5$  ( $D_503$ : 25689303–25697234) was an extraordinary 61 units up to 8 kb; how it was formed is unknown. On the other hand, a total of 463 CSs were found (Table S2). Combining the analyses of the insertion loci of the two constituents, we found 72 TRs and 72 CSs constituting the *ICRd* motifs (Fig. 1).

Further analysis showed that the TR and CS were evenly distributed on the chromosomes based on an  $\chi^2$  test, with the number of insertions being proportional to the size of the chromosome [TR insertions,  $\chi^2 = 5.56$  ( $df = 12$ ,  $P > 0.9$ ); CS insertions,  $\chi^2 = 9.08$  ( $df = 12$ ,  $P > 0.5$ )]. The even distributions meant that the CS and TR are possible ancient repetitive sequences that have evolved along with the chromosomes. Previous *G. raimondii* genome sequencing work reported that most TEs in *G. raimondii* are deletion derivatives lacking the domains that are typically necessary for transposition and that only 3% of LTR base pairs derived from full-length LTR-TEs (Paterson et al., 2012). We show that hundreds of constituents of the *ICRd* motif in  $D_5$  are potentially the fragments produced from the ancient LTR-TEs.

### Disappearance of the *ICRd* motif from *Gossypium*

Aiming to observe the disappearance of the *ICRd* motif in the newly formed *Gossypium* A genome, we selected two homologous segments from the highly collinear Chromosome 1 of *G. raimondii* ( $D_501$ ) and *G. arboreum* ( $A_201$ ) (Li et al., 2014), respectively. The segment from Chromosome 1 of *G. raimondii* ( $D_501$ ) harbored one *ICRd* motif and its homologous segment from  $A_201$  was obtained from homologous SSR markers (Table S4). The





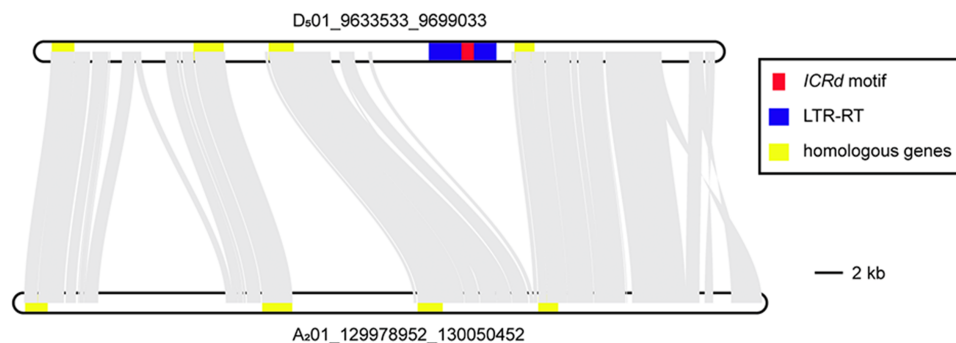
**Figure 6** The distribution of the *ICRd* motif and its constituent in the  $D_5$  genome. (A) Insertions of the *ICRd* motif and its constituents in the  $D_5$  genomes; (B) (C) *ICRd* TR and TR-c chromosomal distribution, the expected (grey) and actual (white) distributions across all chromosomes are illustrated; in addition, the density per megabase is shown for each chromosome.

Full-size [DOI: 10.7717/peerj.8344/fig-6](https://doi.org/10.7717/peerj.8344/fig-6)

illustration of the syntenic block of the two segments showed that the *ICRd* motif together with its host LTR-TE were lost on the  $A_{201}$  segment, while their up- and downstream conservative regions remained (Fig. 7). In the upstream, we observed two insertions rich in repeat sequences especially on the  $A_{201}$  segment (Table S4), which was potentially due to a recent TE expanding event happening in the A genome (Lu et al., 2018a). Thus, we observed that the *ICRd* motifs and host LTR-TEs were lost from the genome with the recent formation of the A genome (Wendel & Cronn, 2001; Wendel, Flagel & Adams, 2012), but remained in the D genome despite mass damage accumulation.

### Distributions of *ICRd* motifs in tetraploid cotton

Tetraploid cotton, *G. hirsutum* and *G. barbadense*, are the major cultivated fiber-producing cotton species. Research on the genome of these two species is an important approach to improving cotton yield and quality. However, due to the large number of homologous segments between A and D-subgenomes, the tetraploid cotton genome assemblage has been a great challenge to molecular geneticists (Bowers et al., 2003; Chen et al., 2007). Through



**Figure 7** The colinearity of the two homologous segments.

Full-size [DOI: 10.7717/peerj.8344/fig-7](https://doi.org/10.7717/peerj.8344/fig-7)

high-throughput sequencing methods, two versions of the *G. hirsutum* genome assembly ((AD)<sub>1</sub>-BGI (Li et al., 2015), (AD)<sub>1</sub>-NBI (Zhang et al., 2015), and two *G. barbadense* versions (AD)<sub>2</sub>-HAU (Yuan et al., 2015) and (AD)<sub>2</sub>-CAS (Liu et al., 2015) were completed in 2015. With the advance of sequencing techniques, the tetraploid genome assemblies were improved in quality (Hu et al., 2019; Wang et al., 2019). However, to benefit research in the post-genome era, such as facilitating molecular breeding of cotton, suitable evaluation is needed to provide accurate reference data. Application of the lineage-specific repetitive element and the *ICRd* motifs are important tools in evaluating the quality of the genome assembly of tetraploid cotton.

To observe the assembling quality of the *ICRd* motif in tetraploid genomes, we queried it with BLAST in all published tetraploid cotton genomes, including four versions of *G. hirsutum* ((AD)<sub>1</sub>) and four versions of *G. barbadense* ((AD)<sub>2</sub>) (Table 2). In the case of (AD)<sub>1</sub>, the two recently published (Hu et al., 2019; Wang et al., 2019) versions and the previous NBI version were in agreement with our FISH inspection results, viz., that the *ICRd* motifs only generated the signals on the D-subgenome chromosomes (Fig. 3). However, the BGI version (Li et al., 2015) is inconsistent with the FISH results in that the *ICRd* motif was misassembled into the A-subgenome. For the (AD)<sub>2</sub> assemblies, the two newly published (Hu et al., 2019; Wang et al., 2019) and CAS versions were better assembled than the HAUv1 version. The HAUv1 showed the number of matches in the chromosome-unassembled scaffolds, while the HAUv2 has improved quality (Table S5). This means that with advances in genome sequencing techniques, tetraploid genomes can be more precisely assembled though the existence of homologous segments from *At* and *Dt*.

## DISCUSSION

### Identification of *ICRd* motif and *Gossypium* evolution

TEs have played an important function in *Gossypium* speciation and the accumulation of different genomic-specific TEs were thought to be responsible for genome-size variation in *Gossypium* (Hawkins et al., 2006). Through FISH inspection, some A genome-specific repetitive elements have been well identified and characterized (Liu et al., 2016), but similar

**Table 2** The distribution of *ICRd* motifs on different genome assemblies of tetraploid cotton.

Tetraploid	Version	Reference	<i>ICRd</i> motif
<i>G. hirsutum</i> (AD) <sub>1</sub>	BGI	<i>Li et al. (2015)</i>	D <sub>h</sub> 01-D <sub>h</sub> 13; A <sub>h</sub> 02, A <sub>h</sub> 05, A <sub>h</sub> 07, A <sub>h</sub> 08
	NBI	<i>Zhang et al. (2015)</i>	D <sub>h</sub> 01-D <sub>h</sub> 13; None in A-sub
	HAU	<i>Wang et al. (2019)</i>	D <sub>h</sub> 01-D <sub>h</sub> 13; None in A-sub
	ZJU	<i>Hu et al. (2019)</i>	D <sub>h</sub> 01-D <sub>h</sub> 13; None in A-sub
<i>G. barbadense</i> (AD) <sub>2</sub>	CAS	<i>Liu et al. (2015)</i>	D <sub>b</sub> 01-D <sub>b</sub> 13; None in A-sub
	HAUv1	<i>Yuan et al. (2015)</i>	D <sub>b</sub> 01, D <sub>b</sub> 02, D <sub>b</sub> 06-D <sub>b</sub> 09, D <sub>b</sub> 12; None in A-sub
	HAUv2	<i>Wang et al. (2019)</i>	D <sub>h</sub> 01-D <sub>h</sub> 13; None in A-sub
	ZJU	<i>Hu et al. (2019)</i>	D <sub>h</sub> 01-D <sub>h</sub> 13; None in A-sub

work in the D genome have been rare; this may be because the genome-specific repetitive sequences in the A genome are much more numerous than in the D genome (*Liu et al., 2018a*). However, in the present study, starting with comparative genomic data, we have screened out one kind of specific sequence in the D genome, and subsequently, we have identified and characterized TEs.

The TEs harboring the *ICRd* motif showed an ancient active date of much earlier than 10 MYA, while the time of divergence of the A and D genomes from the common ancestor is estimated to have occurred 5–10 MYA (*Grover et al., 2004*). Thus the *ICRd* motifs have existed in the ancestor of A and D genome, while disappeared along with the formation of the A genome. Previous researchers have considered that the accumulation of lineage-specific TEs, which is thought to be responsible for the variation of *Gossypium* genomes (*Hawkins et al., 2006*), and the LTR-TE activities after 5 MYA mainly contributed to the two-fold size difference of the A and D genomes (*Zhang et al., 2015*). Based on our analysis, we presumed that as in the activity of new repetitive sequences the extinction of ancient repetitive sequences, such as the disappearance of the *ICRd* motif in the A genome, also contributed significantly to genome evolution. Through FISH, we observed that the *ICRd* motifs were only distributed in D-subgenome chromosomes, and the results were in agreement with a previous study which reported that the TE have proliferated in the progenitor genomes but were retained after allopolyploid formation in the D-subgenome (*Zhang et al., 2015*).

### ***ICRd* motif support cytogenetic markers for tetraploid cotton**

The identification of the *ICRd* motif provides a new subgenome marker for the accurate assembling of tetraploid cotton (*Chen et al., 2007*). Chromosome identification is the foundation of plant genetics, evolution and genomics research (*Saranga, 2007; Xie et al., 2012*). Although many species have been sequenced, the rapid identification of the subgenome is still useful in applied research. FISH has been used as a reliable cytological technique for chromosome identification in many species (*Wang, Guo & Zhang, 2007*), but has only been used recently for the identification of cotton chromosomes (*Gan et al., 2012*). In the present study, the identified *ICRd* motifs can be used as a new cytological marker in *Gossypium*, especially in tetraploids. Further, the repetitive sequence probes are easier and more successfully detected than other probes. Several similar markers have been

reported ([Liu et al., 2016](#)). The addition of these new cytological markers will enrich the marker database for chromosome identification and facilitate cotton genomic studies.

Eukaryotic genomes have a high proportion of TEs and these TEs make eukaryotic genome assembly much more difficult than simple genome assembly ([Treangen & Salzberg, 2012](#)). Many reported genome sequences have gaps because of the high proportion of TEs ([Adams et al., 2000](#)). Allopolyploid genomes are especially difficult to assemble homologous fragments from subgenomes ([Chen et al., 2007](#)). Incorrect assembling of the genomes leads to ambiguity in research which, in turn, produces biases and errors when interpreting results ([Adams et al., 2000](#)). The repetitive sequences analysis in this work were screened out from the whole genome comparison, we characterize the distribution feature on referenced genome assembly, moreover, FISH observation on chromosomes of somatic cell verified the lineage-specific feature. Combining FISH with genome-specific repeat segments is a convenient and practical approach to observe chromosome differences, in addition to assisting polyploid genome assembling, and evaluating assembling accuracy. With the progress of genome sequencing and assembling, genome assembly will become increasingly more precise and convincing, and it is likely that the latter published tetraploid genome will adopt the BioNano and Hi-C approaches ([Hu et al., 2019](#); [Wang et al., 2019](#)) and improve the identification of homologous segments from subgenomes. The improved tetraploid cotton genome assemblies were consistent with FISH, which provides a reference for researchers deciding which genomes to adopt in their research.

## CONCLUSIONS

We identified and characterized a new type of repetitive sequence termed *ICRd* motif in the *Gossypium* D genome. The motifs are interspersed in 13 chromosomes of the D genome, but absent in the A genome, and retained in D-subgenome in tetraploid cotton. We analyzed their structure, genomic distribution, affiliation, and evolution, which revealed a conserved region harbored in ancient LTR-TEs. The identification and characterization of the *ICRd* motif provided new insight into understanding TE evolution along with the formation of cotton genomes as well as providing a convenient and practical tool to distinguish the A and D genome subsets of the tetraploid cotton genome assembly. The *ICRd* motif has a novel structure and affiliation; how the structure was formed and what function the *ICRd* motif plays in LTR-TEs would be valuable areas for future research.

## ACKNOWLEDGEMENTS

We are indebted to Dr Syed Shan-e-Ali Zaidi of the University of Liège, Belgium, for his guidance in analysis and interpretation of the data.

## ADDITIONAL INFORMATION AND DECLARATIONS

### Funding

This research was supported by the National Key Research and Development Plan of China (grants 2016YFD0100306 and 2016YFD0100203) and The Natural Science Foundation of

China (grants 31530053 and 31671745). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

### Grant Disclosures

The following grant information was disclosed by the authors:

The National Key Research and Development Plan of China: 2016YFD0100306, 2016YFD0100203.

The Natural Science Foundation of China: 31530053, 31671745.

### Competing Interests

The authors declare there are no competing interests.

### Author Contributions

- Hejun Lu and Xinglei Cui analyzed the data, conceived and designed the experiments, performed the experiments, prepared figures and/or tables, authored or reviewed drafts of the paper, and approved the final draft.
- Yanyan Zhao performed the experiments, prepared figures and/or tables, authored or reviewed drafts of the paper, and approved the final draft.
- Richard Odongo Magwanga analyzed the data, performed the experiments, prepared figures and/or tables, authored or reviewed drafts of the paper, data analysis and interpretation, and approved the final draft.
- Pengcheng Li performed the experiments, prepared figures and/or tables, data collection, and approved the final draft.
- Xiaoyan Cai performed the experiments, and approved the final draft.
- Zhongli Zhou, Xingxing Wang and Yuling Liu performed the experiments, authored or reviewed drafts of the paper, and approved the final draft.
- Yanchao Xu and Yuqing Hou performed the experiments, prepared figures and/or tables, and approved the final draft.
- Renhai Peng performed the experiments, authored or reviewed drafts of the paper, and approved the final draft.
- Kunbo Wang and Fang Liu conceived and designed the experiments, authored or reviewed drafts of the paper, and approved the final draft.

### Data Availability

The following information was supplied regarding data availability:

The sources of genome assemblies involved in this work: *G. raimondii*, [https://www.cottongen.org/species/Gossypium\\_raimondii/jgi\\_genome\\_221](https://www.cottongen.org/species/Gossypium_raimondii/jgi_genome_221); *G. arboreum*, [https://www.cottongen.org/species/Gossypium\\_arboreum/bgi-A2\\_genome\\_v2.0](https://www.cottongen.org/species/Gossypium_arboreum/bgi-A2_genome_v2.0); *G. hirsutum*, [https://www.cottongen.org/species/Gossypium\\_hirsutum/bgi-AD1\\_genome\\_v1.0](https://www.cottongen.org/species/Gossypium_hirsutum/bgi-AD1_genome_v1.0), [https://www.cottongen.org/species/Gossypium\\_hirsutum/nbi-AD1\\_genome\\_v1.1](https://www.cottongen.org/species/Gossypium_hirsutum/nbi-AD1_genome_v1.1), [https://www.cottongen.org/species/Gossypium\\_hirsutum/HAU-AD1\\_genome\\_v1.0\\_v1.1](https://www.cottongen.org/species/Gossypium_hirsutum/HAU-AD1_genome_v1.0_v1.1), [https://www.cottongen.org/species/Gossypium\\_hirsutum/ZJU-AD1\\_v2.1](https://www.cottongen.org/species/Gossypium_hirsutum/ZJU-AD1_v2.1); *G. barbadense*, [https://www.cottongen.org/species/Gossypium\\_barbadense/nbi-AD2\\_genome\\_v1.0](https://www.cottongen.org/species/Gossypium_barbadense/nbi-AD2_genome_v1.0), [https://www.cottongen.org/species/Gossypium\\_barbadense/ZJU-AD2\\_v1.1](https://www.cottongen.org/species/Gossypium_barbadense/ZJU-AD2_v1.1), [https://www.cottongen.org/species/Gossypium\\_barbadense/HAU-AD2\\_genome\\_v2.0](https://www.cottongen.org/species/Gossypium_barbadense/HAU-AD2_genome_v2.0).

The CAS version of *G. barbadense* is available in GenBank: [PRJNA251673](https://www.ncbi.nlm.nih.gov/nuclot/PRJNA251673).

The Perl scripts are available as [Supplementary File](#).

## Supplemental Information

Supplemental information for this article can be found online at <http://dx.doi.org/10.7717/peerj.8344#supplemental-information>.

## REFERENCES

- Adams MD, Celniker SE, Holt RA, Evans CA, Gocayne JD, Amanatides PG, Scherer SE, Li PW, Hoskins RA, Galle RF, George RA, Lewis SE, Richards S, Ashburner M, Henderson SN, Sutton GG, Wortman JR, Yandell MD, Zhang Q, Chen LX, Brandon RC, Rogers YHC, Blazej RG, Champe M, Pfeiffer BD, Wan KH, Doyle C, Baxter EG, Helt G, Nelson CR, Gabor Miklos GL, Abril JF, Agbayani A, An HJ, Andrews-Pfannkoch C, Baldwin D, Ballew RM, Basu A, Baxendale J, Bayraktaroglu L, Beasley EM, Beeson KY, Benos PV, Berman BP, Bhandari D, Bolshakov S, Borkova D, Botchan MR, Bouck J, Brokstein P, Brottier P, Burtis KC, Busam DA, Butler H, Cadieu E, Center A, Chandra I, Michael Cherry J, Cawley S, Dahlke C, Davenport LB, Davies P, De Pablos B, Delcher A, Deng Z, Deslattes Mays A, Dew I, Dietz SM, Dodson K, Doup LE, Downes M, Dugan-Rocha S, Dunkov BC, Dunn P, Durbin KJ, Evangelista CC, Ferraz C, Ferriera S, Fleischmann W, Fosler C, Gabrielian AE, Garg NS, Gelbart WM, Glasser K, Glodek A, Gong F, Gorrell JHarley, Gu Z, Guan P, Harris M, Harris NL, Harvey D, Heiman TJ, Hernandez JR, Houck J, Hostin D, Houston KA, Howland TJ, Wei MH, Ibegwam C, Jalali M, Kalush F, Karpen GH, Ke Z, Kennison JA, Ketchum KA, Kimmel BE, Kodira CD, Kraft C, Kravitz S, Kulp D, Lai Z, Lasko P, Lei Y, Levitsky AA, Li J, Li Z, Liang Y, Lin X, Liu X, Mattei B, McIntosh TC, McLeod MP, McPherson D, Merkulov G, Milshina NV, Mobarry C, Morris J, Moshrefi A, Mount SM, Moy M, Murphy B, Murphy L, Muzny DM, Nelson DL, Nelson DR, Nelson KA, Nixon K, Nusskern DR, Pacleb JM, Palazzolo M, Pittman GS, Pan S, Pollard J, Puri V, Reese MG, Reinert K, Remington K, Saunders RDC, Scheeler F, Shen H, Christopher Shue B, Siden-Kiamos I, Simpson M, Skupski MP, Smith T, Spier E, Spradling AC, Stapleton M, Strong R, Sun E, Svirskas R, Tector C, Turner R, Venter E, Wang AH, Wang X, Wang ZY, Wassarman DA, Weinstock GM, Weissenbach J, Williams SM, Woodage T, Worley KC, Wu D, Yang S, Alison Yao Q, Ye J, Yeh RF, Zaveri JS, Zhan M, Zhang G, Zhao Q, Zheng L, Zheng XH, Zhong FN, Zhong W, Zhou X, Zhu S, Zhu X, Smith HO, Gibbs RA, Myers EW, Rubin GM, Craig Venter J. 2000. The genome sequence of *Drosophila melanogaster*. *Science* 287:2185–2195 DOI [10.1126/science.287.5461.2185](https://doi.org/10.1126/science.287.5461.2185).
- Adams KL, Palmer JD. 2003. Evolution of mitochondrial gene content: gene loss and transfer to the nucleus. *Molecular Phylogenetics and Evolution* 29:380–395 DOI [10.1016/S1055-7903\(03\)00194-5](https://doi.org/10.1016/S1055-7903(03)00194-5).



- Arkhipova IR. 2018.** Neutral theory, transposable elements, and eukaryotic genome evolution. *Molecular Biology and Evolution* **35**:1332–1337 DOI [10.1093/molbev/msy083](https://doi.org/10.1093/molbev/msy083).
- Bennetzen JL. 2005.** Transposable elements, gene creation and genome rearrangement in flowering plants. *Current Opinion in Genetics and Development* **15**:621–627 DOI [10.1016/j.gde.2005.09.010](https://doi.org/10.1016/j.gde.2005.09.010).
- Benson G. 1999.** Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Research* **27**:573–580 DOI [10.1093/nar/27.2.573](https://doi.org/10.1093/nar/27.2.573).
- Bowen NJ, Jordan IK. 2002.** Transposable elements and the evolution of eukaryotic complexity. *Current Issues in Molecular Biology* **4**:65–76.
- Bowers JE, Chapman BA, Rong J, Paterson AH. 2003.** Unravelling angiosperm genome evolution by phylogenetic analysis of chromosomal duplication events. *Nature* **422**:433–438 DOI [10.1038/nature01521](https://doi.org/10.1038/nature01521).
- Brink RA, Williams E. 1973.** Mutable R-navajo alleles of cyclic origin in maize. *Genetics* **73**:273–296.
- Camacho C, Coulouris V, Avagyan N, Ma J, Papadopoulos K, Bealer TL, Madden G. 2009.** BLAST+: architecture and applications. *BMC Bioinformatics* **10**(1):421 DOI [10.1186/1471-2105-10-421](https://doi.org/10.1186/1471-2105-10-421).
- Chen ZJ, Scheffler BE, Dennis E, Triplett BA, Zhang T, Guo W, Chen X, Stelly DM, Rabinowicz PD, Town CD, Arioli T, Brubaker C, Cantrell RG, Lacape J-M, Ulloa M, Chee P, Gingle AR, Haigler CH, Percy R, Saha S, Wilkins T, Wright RJ, Van Deynze A, Zhu Y, Yu S, Abdurakhmonov I, Katageri I, Kumar PA, Rahman MU, Zafar Y, Yu JZ, Kohel RJ, Wendel JF, Paterson AH. 2007.** Toward sequencing cotton (*Gossypium*) genomes. *Plant Physiology* **145**:1303–1310 DOI [10.1104/pp.107.107672](https://doi.org/10.1104/pp.107.107672).
- Cui X, Liu F, Liu Y, Zhou Z, Zhou Y, Wang C, Wang X, Cai X, Wang Y, Meng F, Peng R, Wang K. 2015.** Construction of cytogenetic map of *Gossypium herbaceum* chromosome 1 and its integration with genetic maps. *Molecular Cytogenetics* **8**(1):2 DOI [10.1186/s13039-015-0106-y](https://doi.org/10.1186/s13039-015-0106-y).
- Del Pozo JC, Ramirez-Parra E. 2015.** Whole genome duplications in plants: an overview from Arabidopsis. *Journal of Experimental Botany* **66**:6991–7003 DOI [10.1093/jxb/erv432](https://doi.org/10.1093/jxb/erv432).
- Deutsch M, Long M. 1999.** Intron-exon structures of eukaryotic model organisms. *Nucleic Acids Research* **27**:3219–3228 DOI [10.1093/nar/27.15.3219](https://doi.org/10.1093/nar/27.15.3219).
- Dong G, Shen J, Zhang Q, Wang J, Yu Q, Ming R, Wang K, Zhang J. 2018.** Development and applications of chromosome-specific cytogenetic BAC-FISH probes in *S. spontaneum*. *Frontiers in Plant Science* **9**:218 DOI [10.3389/fpls.2018.00218](https://doi.org/10.3389/fpls.2018.00218).
- Du X, Huang G, He S, Yang Z, Sun G, Ma X, Li N, Zhang X, Sun J, Liu M, Jia Y, Pan Z, Gong W, Liu Z, Zhu H, Ma L, Liu F, Yang D, Wang F, Fan W, Gong Q, Peng Z, Wang L, Wang X, Xu S, Shang H, Lu C, Zheng H, Huang S, Lin T, Zhu Y, Li F. 2018.** Resequencing of 243 diploid cotton accessions based on an updated A genome identifies the genetic basis of key agronomic traits. *Nature Genetics* **50**:796–802 DOI [10.1038/s41588-018-0116-x](https://doi.org/10.1038/s41588-018-0116-x).

- Edgar RC. 2004.** MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research* **32**:1792–1797 DOI [10.1093/nar/gkh340](https://doi.org/10.1093/nar/gkh340).
- Ellinghaus D, Kurtz S, Willhoeft U. 2008.** LTRharvest, an efficient and flexible software for de novo detection of LTR retrotransposons. *BMC Bioinformatics* **9**:18 DOI [10.1186/1471-2105-9-18](https://doi.org/10.1186/1471-2105-9-18).
- Feschotte C. 2008.** Transposable elements and the evolution of regulatory networks. *Nature Reviews. Genetics* **9**:397–405 DOI [10.1038/nrg2337](https://doi.org/10.1038/nrg2337).
- Finnegan DJ. 1989.** Eukaryotic transposable elements and genome evolution. *Trends in Genetics* **5**:103–107 DOI [10.1016/0168-9525\(89\)90039-5](https://doi.org/10.1016/0168-9525(89)90039-5).
- Gan Y, Liu F, Peng R, Wang C, Li S, Zhang X, Wang Y, Wang K. 2012.** Individual chromosome identification, chromosomal collinearity and genetic-physical integrated map in *Gossypium darwinii* and four D genome cotton species revealed by BAC-FISH. *Genes & Genetic Systems* **87**:233–241 DOI [10.1266/ggs.87.233](https://doi.org/10.1266/ggs.87.233).
- Goldschmidt RB. 2002.** Marginalia to McClintock's Work on Mutable Loci in Maize. *The American Naturalist* **84**:437–455 DOI [10.1086/281640](https://doi.org/10.1086/281640).
- Grover CE, Kim HR, Wing RA, Paterson AH, Wendel JF. 2004.** Incongruent patterns of local and global genome size evolution in cotton. *Genome Research* **14**:1474–1482 DOI [10.1101/gr.2673204](https://doi.org/10.1101/gr.2673204).
- Hanson RE, Zhao X, Islam-Faridi MN, Paterson AH, Zwick MS, Crane CF, McKnight TD, Stelly DM, Price HJAJOB. 1998.** Evolution of interspersed repetitive elements in *Gossypium* (Malvaceae). *American Journal of Botany* **85**(10):1364–1368 DOI [10.2307/2446394](https://doi.org/10.2307/2446394).
- Hawkins JS, Kim HR, Nason JD, Wing RA, Wendel JF. 2006.** Differential lineage-specific amplification of transposable elements is responsible for genome size variation in *Gossypium*. *Genome Research* **16**:1252–1261 DOI [10.1101/gr.5282906](https://doi.org/10.1101/gr.5282906).
- Hendrix B, Stewart JM. 2005.** Estimation of the nuclear DNA content of *Gossypium* species. *Annals of Botany* **95**:789–797 DOI [10.1093/aob/mci078](https://doi.org/10.1093/aob/mci078).
- Hu Y, Chen J, Fang L, Zhang Z, Ma W, Niu Y, Ju L, Deng J, Zhao T, Lian J, Baruch K, Fang D, Liu X, Ruan YL, Rahman MU, Han J, Wang K, Wang Q, Wu H, Mei G, Zang Y, Han Z, Xu C, Shen W, Yang D, Si Z, Dai F, Zou L, Huang F, Bai Y, Zhang Y, Brodt A, Ben-Hamo H, Zhu X, Zhou B, Guan X, Zhu S, Chen X, Zhang T. 2019.** *Gossypium barbadense* and *Gossypium hirsutum* genomes provide insights into the origin and evolution of allotetraploid cotton. *Nature Genetics* **51**(4):739–748 DOI [10.1038/s41588-019-0371-5](https://doi.org/10.1038/s41588-019-0371-5).
- Ibarra-Laclette E, Lyons E, Hernández-Guzmán G, Pérez-Torres CA, Carretero-Paulet L, Chang TH, Lan T, Welch AJ, Juárez MJA, Simpson J, Fernández-Cortés A, Arteaga-Vázquez M, Góngora-Castillo E, Acevedo-Hernández G, Schuster SC, Himmelbauer H, Minoche AE, Xu S, Lynch M, Oropeza-Aburto A, Cervantes-Pérez SA, De Jesús Ortega-Estrada M, Cervantes-Luevano JI, Michael TP, Mockler T, Bryant D, Herrera-Estrella A, Albert VA, Herrera-Estrella L. 2013.** Architecture and evolution of a minute plant genome. *Nature* **498**:94–98 DOI [10.1038/nature12132](https://doi.org/10.1038/nature12132).

- Jeffreys AJ, Neumann R, Wilson V. 1990.** Repeat unit sequence variation in minisatellites: a novel source of DNA polymorphism for studying variation and mutation by single molecule analysis. *Cell* **60**(3):473–485 DOI [10.1016/0092-8674\(90\)90598-9](https://doi.org/10.1016/0092-8674(90)90598-9).
- Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. 2005.** Repbase update, a database of eukaryotic repetitive elements. *Cytogenetic and Genome Research* **110**:462–467 DOI [10.1159/000084979](https://doi.org/10.1159/000084979).
- Kimura M. 1980.** A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *Journal of Molecular Evolution* **16**:111–120 DOI [10.1007/BF01731581](https://doi.org/10.1007/BF01731581).
- Koike A, Nakai K, Takagi T. 2002.** The origin and evolution of eukaryotic protein kinases. *Genome Letters* **1**:83–104 DOI [10.1166/gl.2002.010](https://doi.org/10.1166/gl.2002.010).
- Koonin EV, Csuros M, Rogozin IB. 2013.** Whence genes in pieces: reconstruction of the exon-intron gene structures of the last eukaryotic common ancestor and other ancestral eukaryotes. *Wiley Interdisciplinary Reviews: RNA* **4**:93–105 DOI [10.1002/wrna.1143](https://doi.org/10.1002/wrna.1143).
- Kurtz S. 2003.** The Vmatch large scale sequence analysis software. Ref type: computer program.
- Li F, Fan G, Lu C, Xiao G, Zou C, Kohel RJ, Ma Z, Shang H, Ma X, Wu J, Liang X, Huang G, Percy RG, Liu K, Yang W, Chen W, Du X, Shi C, Yuan Y, Ye W, Liu X, Zhang X, Liu W, Wei H, Wei S, Huang G, Zhang X, Zhu S, Zhang H, Sun F, Wang X, Liang J, Wang J, He Q, Huang L, Wang J, Cui J, Song G, Wang K, Xu X, Yu JZ, Zhu Y, Yu S. 2015.** Genome sequence of cultivated Upland cotton (*Gossypium hirsutum* TM-1) provides insights into genome evolution. *Nature Biotechnology* **33**:524–530 DOI [10.1038/nbt.3208](https://doi.org/10.1038/nbt.3208).
- Li F, Fan G, Wang K, Sun F, Yuan Y, Song G, Li Q, Ma Z, Lu C, Zou C, Chen W, Liang X, Shang H, Liu W, Shi C, Xiao G, Gou C, Ye W, Xu X, Zhang X, Wei H, Li Z, Zhang G, Wang J, Liu K, Kohel RJ, Percy RG, Yu JZ, Zhu YX, Wang J, Yu S. 2014.** Genome sequence of the cultivated cotton *Gossypium arboreum*. *Nature Genetics* **46**:567–572 DOI [10.1038/ng.2987](https://doi.org/10.1038/ng.2987).
- Liu X, Zhao B, Zheng HJ, Hu Y, Lu G, Yang CQ, Chen JD, Chen JJ, Chen DY, Zhang L, Zhou Y, Wang LJ, Guo WZ, Bai YL, Ruan JX, Shangguan XX, Mao YB, Shan CM, Jiang JP, Zhu YQ, Jin L, Kang H, Chen ST, He XL, Wang R, Wang YZ, Chen J, Wang LJ, Yu ST, Wang BY, Wei J, Song SC, Lu XY, Gao ZC, Gu WY, Deng X, Ma D, Wang S, Liang WH, Fang L, Cai CP, Zhu XF, Zhou BL, Chen ZJ, Xu SH, Zhang YG, Wang SY, Zhang TZ, Zhao GP, Chen XY. 2015.** *Gossypium barbadense* genome sequence provides insight into the evolution of extra-long staple fiber and specialized metabolites. *Scientific Reports* **5**:14139 DOI [10.1038/srep14139](https://doi.org/10.1038/srep14139).
- Liu Y, Peng R, Liu F, Wang X, Cui X, Zhou Z, Wang C, Cai X, Wang Y, Lin Z, Wang K. 2016.** A *Gossypium* BAC clone contains key repeat components distinguishing sub-genome of allotetraploidy cottons. *Molecular Cytogenetics* **9**:27 DOI [10.1186/s13039-016-0235-y](https://doi.org/10.1186/s13039-016-0235-y).
- Liu Z, Liu Y, Liu F, Zhang S, Wang X, Lu Q, Wang K, Zhang B, Peng R. 2018.** Genome-wide survey and comparative analysis of long terminal repeat (LTR)

- retrotransposon families in four gossypium species. *Scientific Reports* **8**:9399 DOI [10.1038/s41598-018-27589-6](https://doi.org/10.1038/s41598-018-27589-6).
- Lu H, Cui X, Liu Z, Liu Y, Wang X, Zhou Z, Cai X, Zhang Z, Guo X, Hua J, Ma Z, Wang X, Zhang J, Zhang H, Liu F, Wang K. 2018.** Discovery and annotation of a novel transposable element family in *Gossypium*. *BMC Plant Biology* **18**:307 DOI [10.1186/s12870-018-1519-7](https://doi.org/10.1186/s12870-018-1519-7).
- Ma J, Bennetzen JL. 2004.** Rapid recent growth and divergence of rice nuclear genomes. *Proceedings of the National Academy of Sciences of the United States of America* **101**:12404–12410 DOI [10.1073/pnas.0403715101](https://doi.org/10.1073/pnas.0403715101).
- Marchler-Bauer A, Bo Y, Han L, He J, Lanczycki CJ, Lu S, Chitsaz F, Derbyshire MK, Geer RC, Gonzales NR, Gwadz M, Hurwitz DI, Lu F, Marchler GH, Song JS, Thanki N, Wang Z, Yamashita RA, Zhang D, Zheng C, Geer LY, Bryant SH. 2017.** CDD/SPARCLE: functional classification of proteins via subfamily domain architectures. *Nucleic Acids Research* **45**:D200–D203 DOI [10.1093/nar/gkw1129](https://doi.org/10.1093/nar/gkw1129).
- McClintock B. 1950.** The origin and behavior of mutable loci in maize. *Proceedings of the National Academy of Sciences of the United States of America* **36**:344–355 DOI [10.1073/pnas.36.6.344](https://doi.org/10.1073/pnas.36.6.344).
- Morgante M. 2006.** Plant genome organisation and diversity: the year of the junk!. *Current Opinion in Biotechnology* **17**:168–173 DOI [10.1016/j.copbio.2006.03.001](https://doi.org/10.1016/j.copbio.2006.03.001).
- Munoz-Lopez M, Garcia-Perez J. 2010.** DNA transposons: nature and applications in genomics. *Current Genomics* **11**:115–128 DOI [10.2174/1389202107908886871](https://doi.org/10.2174/1389202107908886871).
- Okonechnikov K, Golosova O, Fursov M, Team U. 2012.** Unipro UGENE: a unified bioinformatics toolkit. *Bioinformatics* **28**(8):1166–1167 DOI [10.1093/bioinformatics/bts091](https://doi.org/10.1093/bioinformatics/bts091).
- Paterson AH, Wendel JF, Gundlach H, Guo H, Jenkins J, Jin D, Llewellyn D, Showmaker KC, Shu S, Udall J, Yoo MJ, Byers R, Chen W, Doron-Faigenboim A, Duke MV, Gong L, Grimwood J, Grover C, Grupp K, Hu G, Lee TH, Li J, Lin L, Liu T, Marler BS, Page JT, Roberts AW, Romanel E, Sanders WS, Szadkowski E, Tan X, Tang H, Xu C, Wang J, Wang Z, Zhang D, Zhang L, Ashrafi H, Bedon F, Bowers JE, Brubaker CL, Chee PW, Das S, Gingle AR, Haigler CH, Harker D, Hoffmann LV, Hovav R, Jones DC, Lemke C, Mansoor S, Rahman MU, Rainville LN, Rambani A, Reddy UK, Rong JK, Saranga Y, Scheffler BE, Scheffler JA, Stelly DM, Triplett BA, Van Deynze A, Vaslin MFS, Waghmare VN, Walford SA, Wright RJ, Zaki EA, Zhang T, Dennis ES, Mayer KFX, Peterson DG, Rokhsar DS, Wang X, Schmutz J. 2012.** Repeated polyploidization of *Gossypium* genomes and the evolution of spinnable cotton fibres. *Nature* **492**:423–427 DOI [10.1038/nature11798](https://doi.org/10.1038/nature11798).
- Piegu B, Guyot R, Picault N, Roulin A, Saniyal A, Kim H, Collura K, Brar DS, Jackson S, Wing RA, Panaud O. 2006.** Doubling genome size without polyploidization: dynamics of retrotransposition-driven genomic expansions in *Oryza australiensis*, a wild relative of rice. *Genome Research* **16**:1262–1269 DOI [10.1101/gr.5290206](https://doi.org/10.1101/gr.5290206).
- Sanmiguel P, Bennetzen JL. 1998.** Evidence that a recent increase in maize genome size was caused by the massive amplification of intergene retrotransposons. *Annals of Botany* **82**:37–44 DOI [10.1006/anbo.1998.0746](https://doi.org/10.1006/anbo.1998.0746).

- Saranga Y.** 2007. Special issue: a century of wheat research—from wild emmer discovery to genome analysis. *Israel Journal of Plant Sciences* 55:207–313 DOI 10.1560/IJPS.55.3-4.207.
- Senchina DS, Alvarez I, Cronn RC, Liu B, Rong J, Noyes RD, Paterson AH, Wing RA, Wilkins TA, Wendel JF.** 2003. Rate variation among nuclear genes and the age of polyploidy in *Gossypium*. *Molecular Biology and Evolution* 20:633–643 DOI 10.1093/molbev/msg065.
- Treangen TJ, Salzberg SL.** 2012. Repetitive DNA and next-generation sequencing: computational challenges and solutions. *Nature Reviews Genetics* 13:36–46 DOI 10.1038/nrg3117.
- Wang K, Guo W, Zhang T.** 2007. Detection and mapping of homologous and homoeologous segments in homoeologous groups of allotetraploid cotton by BAC-FISH. *BMC Genomics* 8:178 DOI 10.1186/1471-2164-8-178.
- Wang K, Huang G, Zhu Y.** 2016. Transposable elements play an important role during cotton genome evolution and fiber cell development. *Science China Life Sciences* 59:112–121 DOI 10.1007/s11427-015-4928-y.
- Wang KB, Wang WK, Wang CY, Song GL, Cui RX, Li SH, Zhang XD.** 2001. Studies of FISH and karyotype of *Gossypium barbadense*. *Yi chuan xue bao = Acta genetica Sinica* 28:69–75.
- Wang K, Wang Z, Li F, Ye W, Wang J, Song G, Yue Z, Cong L, Shang H, Zhu S, Zou C, Li Q, Yuan Y, Lu C, Wei H, Gou C, Zheng Z, Yin Y, Zhang X, Liu K, Wang B, Song C, Shi N, Kohel RJ, Percy RG, Yu JZ, Zhu YX, Cang J, Yu S.** 2012. The draft genome of a diploid cotton *Gossypium raimondii*. *Nature Genetics* 44:1098–1103 DOI 10.1038/ng.2371.
- Wang K, Wendel JF, Hua J.** 2018. Designations for individual genomes and chromosomes in *Gossypium*. *Journal of Cotton Research* 1:3 DOI 10.1186/s42397-018-0002-1.
- Wang M, Tu L, Yuan D, Zhu D, Shen C, Li J, Liu F, Pei L, Wang P, Zhao G, Ye Z, Huang H, Yan F, Ma Y, Zhang L, Liu M, You J, Yang Y, Liu Z, Huang F, Li B, Qiu P, Zhang Q, Zhu L, Jin S, Yang X, Min L, Li G, Chen LL, Zheng H, Lindsey K, Lin Z, Udall JA, Zhang X.** 2019. Reference genome sequences of two cultivated allotetraploid cottons, *Gossypium hirsutum* and *Gossypium barbadense*. *Nature Genetics* 51:224–229 DOI 10.1038/s41588-018-0282-x.
- Wendel JF, Cronn RC.** 2001. Polyploidy and the evolutionary history of cotton. *Advances in Agronomy* 78:139–186 DOI 10.1016/S0065-2113(02)78004-8.
- Wendel JF, Flagel LE, Adams KL.** 2012. Jeans, genes, and genomes: cotton as a model for studying polyploidy. In: *Polyploidy and genome evolution*. Berlin, Heidelberg: Springer 181–207 DOI 10.1007/978-3-642-31442-1\_10.
- Wessler SR.** 2006. Transposable elements and the evolution of eukaryotic genomes. *Proceedings of the National Academy of Sciences of the United States of America* 103:17600–17601 DOI 10.1073/pnas.0607612103.
- Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Cappy P, Chalhoub B, Flavell A, Leroy P, Morgante M, Panaud O, Paux E, SanMiguel P, Schulman AH.** 2007. A unified

- classification system for eukaryotic transposable elements. *Nature Reviews Genetics* **8**:973–982 DOI [10.1038/nrg2165](https://doi.org/10.1038/nrg2165).
- Wickham H.** 2016. ggplot2: elegant graphics for data analysis. Berlin, Heidelberg: Springer.
- Xie Y, Dong F, Hong D, Wan L, Liu P, Yang G.** 2012. Exploiting comparative mapping among Brassica species to accelerate the physical delimitation of a genic male-sterile locus (BnRf) in Brassica napus. *Theoretical and Applied Genetics* **125**:211–222 DOI [10.1007/s00122-012-1826-6](https://doi.org/10.1007/s00122-012-1826-6).
- Yuan D, Tang Z, Wang M, Gao W, Tu L, Jin X, Chen L, He Y, Zhang L, Zhu L, Li Y, Liang Q, Lin Z, Yang X, Liu N, Jin S, Lei Y, Ding Y, Li G, Ruan X, Ruan Y, Zhang X.** 2015. The genome sequence of Sea-Island cotton (*Gossypium barbadense*) provides insights into the allopolyploidization and development of superior spinnable fibres. *Scientific Reports* **5**:17662 DOI [10.1038/srep17662](https://doi.org/10.1038/srep17662).
- Zhang T, Hu Y, Jiang W, Fang L, Guan X, Chen J, Zhang J, Saski CA, Scheffler BE, Stelly DM, Hulse-Kemp AM, Wan Q, Liu B, Liu C, Wang S, Pan M, Wang Y, Wang D, Ye W, Chang L, Zhang W, Song Q, Kirkbride RC, Chen X, Dennis E, Llewellyn DJ, Peterson DG, Thaxton P, Jones DC, Wang Q, Xu X, Zhang H, Wu H, Zhou L, Mei G, Chen S, Tian Y, Xiang D, Li X, Ding J, Zuo Q, Tao L, Liu Y, Li J, Lin Y, Hui Y, Cao Z, Cai C, Zhu X, Jiang Z, Zhou B, Guo W, Li R, Chen ZJ.** 2015. Sequencing of allotetraploid cotton (*Gossypium hirsutum* L. acc. TM-1) provides a resource for fiber improvement. *Nature Biotechnology* **33**:531–537 DOI [10.1038/nbt.3207](https://doi.org/10.1038/nbt.3207).
- Zhang X, Tolzmann CA, Melcher M, Haas BJ, Gardner MJ, Smith JD, Feagin JE.** 2011. Branch point identification and sequence requirements for intron splicing in plasmodium falciparum. *Eukaryotic Cell* **10**:1422–1428 DOI [10.1128/EC.05193-11](https://doi.org/10.1128/EC.05193-11).
- Zhao X, Ji Y, Ding X, Stelly DM, Paterson AHJPMB.** 1998a. Macromolecular organization and genetic mapping of a rapidly evolving chromosome-specific tandem repeat family (B77) in cotton (*Gossypium*). *Plant Molecular Biology* **38**(6):1031–1041 DOI [10.1023/A:1006073116627](https://doi.org/10.1023/A:1006073116627).
- Zhao XP, Si Y, Hanson RE, Crane CF, Price HJ, Stelly DM, Wendel JF, Paterson AH.** 1998b. Dispersed repetitive DNA has spread to new genomes since polyploid formation in cotton. *Genome Research* **8**(5):479–492 DOI [10.1101/gr.8.5.479](https://doi.org/10.1101/gr.8.5.479).