

# Identification and characterization of *CACTA* transposable elements capturing gene fragments in maize

LI Qing<sup>1</sup>, LI Lin<sup>1</sup>, DAI JingRui<sup>1</sup>, LI JianSheng<sup>1</sup> & YAN JianBing<sup>1,2†</sup>

<sup>1</sup>National Maize Improvement Center of China, China Agricultural University, Beijing 100193, China;

<sup>2</sup>International Maize and Wheat Improvement Center (CIMMYT), Apdo. Postal 6-641, Mexico 06600, Mexico

**Transposable elements (TEs)-mediated gene sequence movement is thought to play an important role in genome expansion and origin of genes with novel functions. In this study, a gene, *HGGT*, involved in vitamin E synthesis was used in a case study to discover and characterize transposons carrying gene fragments in maize. A total of 69 transposons that are distributed across the 10 chromosomes and have an average length of 3689 bp were identified from the maize sequence database by using the BLAST search algorithm. Three of these carry gene fragments from the progenitor *HGGT* gene, while the rest (66) contain gene fragments from other cellular genes. Nine of the 69 transposons contain fragments derived from two locations in the genome. By querying the maize Expressed Sequence Tag (EST) database, we found that at least thirteen out of the 69 TEs had corresponding transcripts. More interestingly, two transposons that carry gene fragments from two different chromosomal loci could be expressed as chimeric transcripts.**

*CACTA* elements, gene fragment, genome distribution, chimeric transcripts, maize

Transposable elements (TEs), which move from place to place in the genome and often make duplicates of themselves, have tremendous impact on genome evolution, origin of novel genes and phenotypic variation<sup>[1-3]</sup>. Based on the form of their transposition intermediate, TEs can be divided into two classes<sup>[4]</sup>: Class I elements such as LINE and SINE that are transposed via an RNA intermediate and class II elements such as *Mu* and *Helitron* elements that are transposed directly from DNA to DNA. *En/Spm* elements belong to class II and are widely distributed in genomes of plants like maize<sup>[5,6]</sup>, sorghum<sup>[7]</sup>, rice<sup>[8]</sup>, soybean<sup>[9]</sup> and others. *En/Spm* elements are also referred to as *CACTA* elements because they carry the conserved 5 bp sequence 'CACTA' in the outermost terminal inverted repeats (TIRs)<sup>[10]</sup>. *CACTA* elements are characterized by the generation of a 3 bp target site duplication (TSD) upon insertion and the widespread occurrence of sub-terminal repeats (STRs). The STRs are usually regarded as the binding sites of transposase and are thus related to the transposition of

TEs<sup>[11]</sup>.

The enhancer (*En*) and suppressor-mutator (*Spm*) elements in maize that were independently identified by Peterson<sup>[5]</sup> and McClintock<sup>[6]</sup> are genetically and molecularly homologous<sup>[12-16]</sup>. Both are autonomous elements since they encode the transposase that is required for transposition. The corresponding non-autonomous inhibitor (*I*)<sup>[5]</sup> and defective *Spm* (*dSpm*)<sup>[16]</sup> are deletion derivatives of the two autonomous elements<sup>[13-15]</sup>, and are capable of transposition only in the presence of a functional transposase. In addition to those non-autonomous elements derived from autonomous elements, other types of non-autonomous elements that carry gene segments have also been found to be widely distributed across species. The first elements discovered

Received October 22, 2008; accepted December 14, 2008

doi: 10.1007/s11434-009-0061-2

†Corresponding author (email: yjianbing@cau.edu.cn)

Supported by National Natural Science Foundation of China (Grant No. 30500322) and National Hi-Tech Research and Development Program of China (Grant Nos. 2006AA10Z183, 2006AA10A107)

were the *Mu* elements in maize<sup>[17]</sup>, and, subsequently, *Mutator*-like elements (MULE) carrying gene segments were identified in rice<sup>[18–20]</sup>, *Arabidopsis*<sup>[21]</sup>, *Lotus*<sup>[22]</sup> and melon<sup>[23]</sup>. *Helitron*, which was recently identified in maize, is also capable of acquiring gene fragments<sup>[24–26]</sup>. *CACTA* elements that capture host DNA segments have also been identified in Japanese morning glory<sup>[27,28]</sup>, soybean<sup>[9]</sup> and *Antirrhinum*<sup>[29]</sup>, but there were no reports of such elements in maize. The objectives of this study were to identify and characterize *CACTA* elements carrying gene fragments in maize and elucidate their role in helping to understand the mechanism of genome expansion and origin of genes with novel functions.

## 1 Materials and methods

### 1.1 Determination of TE feature sequences

The repeat sequences of TE were identified using the CENSOR software<sup>[30]</sup>. Analysis of TSD and TIRs were done manually based on their respective characteristics, that is, TSD are usually 3 bp in length and the left and right TIRs are usually reverse complementary to each other. STRs were analyzed with the DNAMAN software (Version 5.2.2, Lynnon Biosoft, Canada).

### 1.2 The genome-wide search of TEs

Two methods were used to search the maize B73 genome for TEs carrying gene fragments, namely a BLAST method and a search method based on previously characterized TIRs and TSD. In the BLAST method: the repeat sequences identified by CENSOR in T2, T3 and T9 (see results) were aligned using MUSCLE<sup>[31]</sup>. A consensus sequence (Figure 1 (a)) obtained with BioEdit<sup>[32]</sup> software was used to query the HTGS database of maize (GenBank Release164,

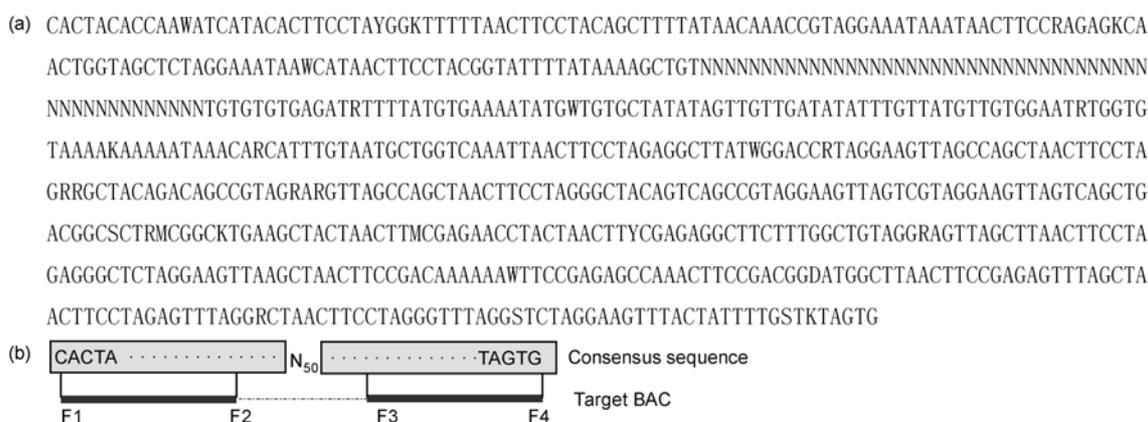
downloaded in June 2008) to identify sequences starting with *CACTA* and ending with *TAGTG*. These sequences generated four position numbers in the BACs where they are incorporated (Figure 1 (b)): F1, F2 (the end position of 5' alignment), F3 (the start position of 3' alignment) and F4. Sequences that matched the criteria outlined below were recovered: F1<F2<F3<F4 or F1>F2>F3>F4, no gaps and  $|F3 - F2| - 1 < 10$  kb.

In the search method based on previously characterized TIRs and TSD, two short sequences  $N_1N_2N_3$ -*CACTA* and *TAGTGN* $_1N_2N_3$  (N represents any of the four nucleotides, where the two  $N_1N_2N_3$  are the same) were screened for across the whole genome in order to retrieve sequences with a length of less than 10 kb between the two short sequences. If the first and last 14 bp (excluding  $N_1N_2N_3$ ) of the retrieved sequences were reverse complementary to each other, and if either the first 17 bp were  $N_1N_2N_3$ *CACTACACCAA*AAT or the last 17 bp were *ATTTTGGTGTAGTGN* $_1N_2N_3$ , the sequences would be used in the non-redundant analysis.

A non-redundant analysis was performed under the condition that TEs were identical across the entire length, or if the identity was less than 100%, the two BACs containing the TEs had to overlap on the maize physical map (<http://www.genome.arizona.edu/fpc/WebAGCoL/maize/WebFPC/>), the TSD must be the same, and the length of the TEs must be similar ( $\pm 1$ ).

### 1.3 Annotation and copy number estimation

Online Blast2GO (<http://www.blast2go.de/>) was utilized to annotate all the TEs with cutoff *E* value  $\leq -10$  (database updated in June 2008). We identified candidate homologous TEs if they were annotated to a similar protein. The identified TEs were used as query sequences



**Figure 1** Consensus sequence used in the BLAST method (a) and schematic diagram of the typical blast result (b).

to blast against the maize HTGS database (<http://www.ncbi.nlm.nih.gov/>). TEs that match each other were defined as copies of the same element.

#### 1.4 PCR amplification of TEs

Primers were designed using Primer3 (<http://frodo.wi.mit.edu/>). PCR was performed in a 15  $\mu$ L reaction volume containing 40 ng DNA, a final concentration of  $1 \times$  Taq buffer (including  $Mg^{2+}$ ), 200  $\mu$ mol/L each dNTP, 0.4  $\mu$ mol/L each primer and 0.5 U Taq polymerase. The PCR program is as follows: 1 cycle at 94°C for 3 min; 94°C, 30 s, 58°C, 50 s, 72°C, 1 min, 35 cycles; and a final extension cycle at 72°C for 10 min. Sequencing was carried out at the SunBio Company. Primers used in this study were: 1F: 5'-GATTGTGGCTATGTATTAGGG-3', 1R: 5'-AGGAAAGAGGCATGGATTA-3', 2F: 5'-CGCTGCTTCTCGTCTCCAC-3', 2R: 5'-TCGTATGCAATCGCCTGA-3', 3F: 5'-CCCGCTGCTTCTCATC-TCC-3' and 3R: 5'-GAGGAATGCCTGCGTTGTC-3'.

#### 1.5 Phylogenetic analysis

The 100 bp sequences from both the 5' and 3' ends of the TEs annotated to transposase or cellular genes were extracted and fused to form a 200 bp sequence, which was then aligned using the multiple sequence alignment tool, MUSCLE<sup>[31]</sup>. The maximum likelihood method in the PHYLIP package<sup>[33]</sup> was used to construct a phylogenetic tree with a bootstrap value of 100. A consensus tree was constructed using MEGA3.0<sup>[34]</sup> with a cutoff of value greater than 50.

#### 1.6 Expression analysis

The method used for expression analysis was a modified version of the method described by Jiang et al.<sup>[18]</sup>. The unmasked sequences with CENSOR were used as queries to search against the maize EST database (updated in June 2008). The TE was considered to have an EST

match if, (i) the sequence similarity between the TE and the EST was greater than 99% with a coverage of >95% of the entire EST length, and (ii) the best match of the EST in the maize HTGS database was the original TE and no other better matches existed in the whole genome.

## 2 Results

### 2.1 Identification of *CACTA* elements carrying gene fragments in maize

(i) Parts of the *HGGT* gene were captured by *CACTA* elements. Homogentisic acid geranylgeranyl transferase (*HGGT*) catalyzes the committed step of vitamin E biosynthesis<sup>[35]</sup>. The *HGGT* gene is located on chromosome 9 in maize genome. The BAC clone, AC210150, contains the full-length sequence of *HGGT* gene and another three BACs contain part of *HGGT* gene (Figure 2). For example, AC215796, a BAC clone located on chromosome 9, contains *HGGT* gene sequences from intron 2 to intron 6 and intron 11 to the 3' UTR (the two parts were directly joined with no other interleaved sequences). The similarities between each of two segments and the progenitor gene are 92% and 96%, respectively (Figure 2).

The upstream and downstream 4 kb sequences flanking the homologous regions in the 3 BACs show high levels of similarity to ENSPM2\_ZM, a class of *CACTA* elements, implying that they might also be TEs. The typical 3 bp TSD and 14 bp TIRs were identified in both AC215501 and AC214002. AC215796 also contains an identifiable TSD and TIRs, although some variations were observed (Table 1). The first 5 bp of one of the two TIRs are CACTA, demonstrating they are likely members of *CACTA* elements. The three TEs were named T2, T3 and T9 according to their location on the chromosomes, and are 4202, 3623 and 3349 bp in length, re-



**Figure 2** Schematic diagram of matches from the NCBI maize HTGS database that match the maize *HGGT* gene. Filled green boxes represent exons, introns are shown by “V” lines. The red lines indicate the sequences that show similarity to the progenitor *HGGT* gene, and the grey lines represent the sequences deleted from AC215796. The position of matched sequences in the BAC and their percent identity with the progenitor *HGGT* gene are shown on the right. The alternatively spliced introns of AC215501 and AC215796 are highlighted by yellow and pink lines, respectively.

**Table 1** Features of T2, T3, T9 and ENSPM2\_ZM

Name	BAC	Position	Length (bp)	TSD	TIR <sup>a)</sup>
T2	AC215501	34906–	4202	GAC	<u>CACTACACCAA</u> AAT
		39107		GAC	ATTTTGGT <u>G</u> TAGTG
T3	AC214002	186704–	3623	CAA	<u>CACTACACCA</u> AATAT
		190326		CAA	ATTTTGGT <u>G</u> TAGTG
T9	AC215796	17640–	3349	TGG	<u>CCATA</u> CACCAAAT
		20988		TGA	ATTTTGGT <u>T</u> AGTA
ENSPM2_ZM <sup>b)</sup>	–	–	7628	–	<u>CACTACACCAA</u> AAT ATTTTGGT <u>G</u> TAGTG

a) The underlined 5 nucleotides represent the conserved sequence for which *CACTA* elements were named. The italic letters indicate mutated nucleotides compared with the T2 sequence. All the sequences were in 5' to 3' orientation. b) ENSPM2\_ZM is a consensus sequence deposited in the CENSOR database.

spectively. The TIRs of T2, T3, T9 and ENSPM2\_ZM are highly similar (Table 1). The 8 bp motif TAACTTCC, occurring at both the 5' and 3' ends as direct or inverted repeats (Figure 5), are identical, indicating that T2, T3 and T9 belong to ENSPM2\_ZM family.

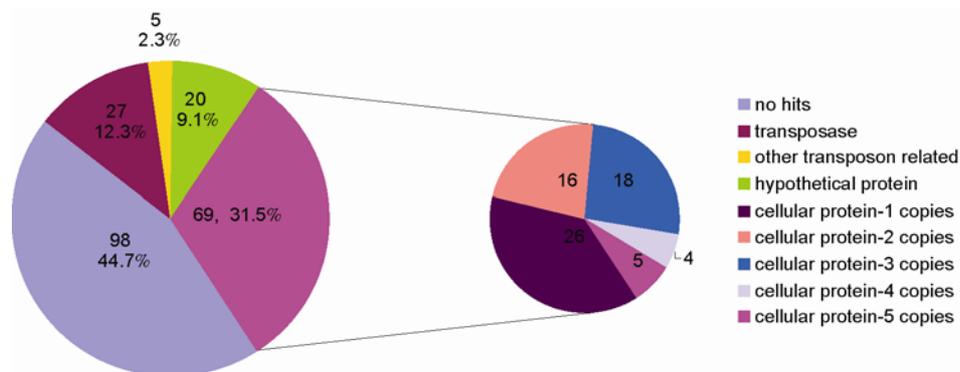
(ii) *CACTA* elements that carry gene fragments in B73 genome. A total of 238 TEs were identified in B73 genome using the BLAST method while 233 TEs were tagged in the whole-genome scan using the sequences of previously characterized TIRs and TSD. As most of the TEs identified using the two different methods were overlapped, a total of 219 non-redundant TEs (including T2, T3 and T9) were obtained and these were designated ES2-1 to ES2-219.

As shown in Figure 3, at  $E \leq -10$ , 69 TEs (31.5%) contained sequences corresponding to previously identified proteins involved in 43 different cellular functions, 32 TEs (14.6%) contained sequences similar to transposase or other transposon-related proteins, 20 TEs (9.1%) were homologous to predicted or hypothetical proteins and 98 TEs (44.7%) showed no similarity to either transposon-related, cellular or hypothetical proteins. Therefore, at least 69 ES2 TEs carrying gene fragments were identified in the B73 genome.

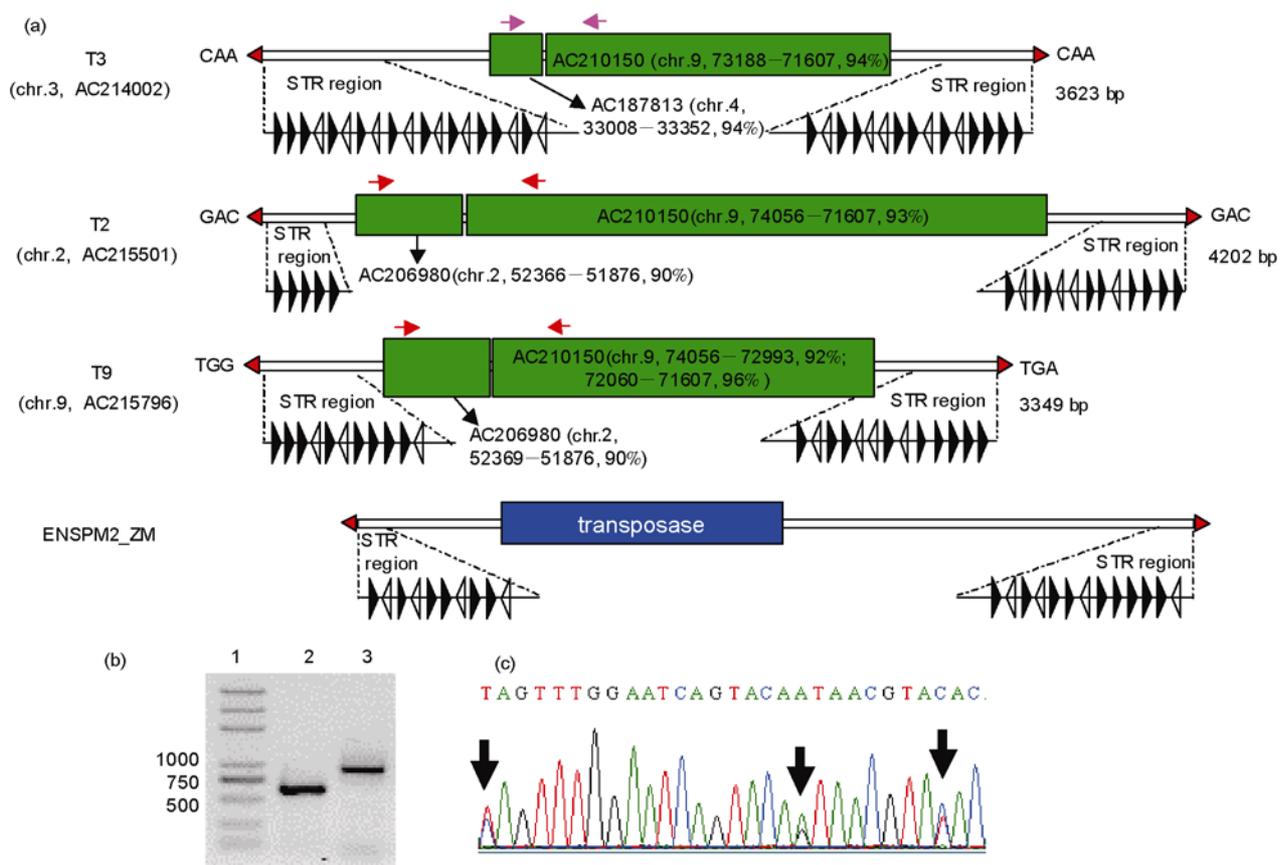
## 2.2 Characteristics of *CACTA* elements carrying gene fragments

(i) Distribution of *CACTA* elements carrying gene fragments in the maize genome. The 69 TEs identified in this study were distributed across the 10 maize chromosomes. Chromosome 3 and chromosome 4 had the highest number of TEs, with each having 11 TEs, while chromosome 8 had the lowest number, having only 2 TEs (Figure 4(a)). The 3 bp TSD was rich in A/T nucleotides. The 20 bp sequences flanking the TSD were also rich in A/T nucleotides with an exception of the innermost nucleotide tending to be G or C (Figure 4(c)). The accumulated length of the 69 TEs is 254.5 kb (sequences between but excluding the TSD), and accounted for 0.01% of the maize genome (254.5 kb/2500 Mb). Of these TEs, 65 had a length between 2 and 5 kb, while four were longer than 5 kb (Figure 4(b)), and the average length was 3689 bp.

The copy number of each TE ranged from one to five (Figure 3). Twenty-six of the 69 TEs were single-copy elements, eight had two copies, six had three copies, one TE had 4 and 5 copies, respectively. TEs with more than five copies were not identified. The average copy num-

**Figure 3** Annotation of ES2 transposable elements.



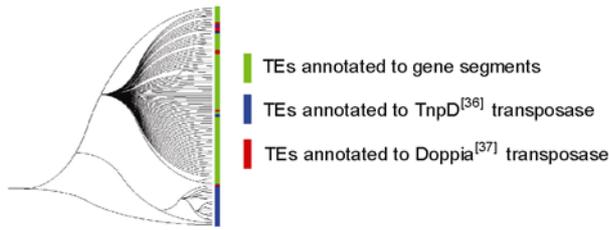


**Figure 5** Transposable element structure and PCR amplification. (a) TIRs are indicated by red filled triangles. The filled and open black triangles represent the 8 bp direct and inverted STRs, respectively. Green boxes represent the gene fragments carried by each TE, each fragment's corresponding progenitor sequences and the percent identity between them are also shown. Pink and red arrows indicate primer pair 1 and primer pair 2, respectively. (b) PCR products of primer pair 1 (lane2) and primer pair 2 (lane3) amplified with B73 genome. Lane1 is a DL2000 Plus DNA Marker. (c) Sequencing chromatogram of primer pair 2. Black arrows represent double peaks due to nucleotide acid differences between T2 and T9.

**Table 2** Transposable elements with two protein hits

ID	BAC <sup>a)</sup>	Chromosome	Position <sup>b)</sup>	Length (bp)	Copy	E value <sup>c)</sup>	Annotation
ES2-91	AC197362	8	101962–105741	3780	3	8.00×10 <sup>-36</sup> 1.00×10 <sup>-16</sup>	sigma factor protein CIPK-like protein
ES2-115	AC149034	NA	110312–114443	4132	1	2.00×10 <sup>-19</sup> 1.00×10 <sup>-17</sup>	DEAD-box ATP-dependent RNA helicase Os03g0276600
ES2-134	AC214230	NA	81652–86041	4390	2	1.00×10 <sup>-34</sup> 2.00×10 <sup>-31</sup>	hypothetical protein OsI_007935 U3 snoRNA-associated protein 11
ES2-142	AC186590	6	23277–28334	5058	5	1.00×10 <sup>-50</sup> 3.00×10 <sup>-24</sup>	putative protein kinase AtSIK putative 3-ketoacyl-CoA thiolase
ES2-149	AC208430	5	156174–161550	5377	4	7.00×10 <sup>-29</sup> 2.00×10 <sup>-24</sup>	putative kinase endoplasmatic reticulum retrieval protein
ES2-154	AC212442	4	171638–178544	6907	2	2.00×10 <sup>-29</sup> 3.00×10 <sup>-16</sup>	one repeat MYB transcriptional factor U3 snoRNA-associated protein 11

a) Accession number of the BAC in GenBank; b) start and end position of the TEs in the BAC; c) E value of annotation.



**Figure 6** Phylogenetic analysis.

in T9 (Figure 7).

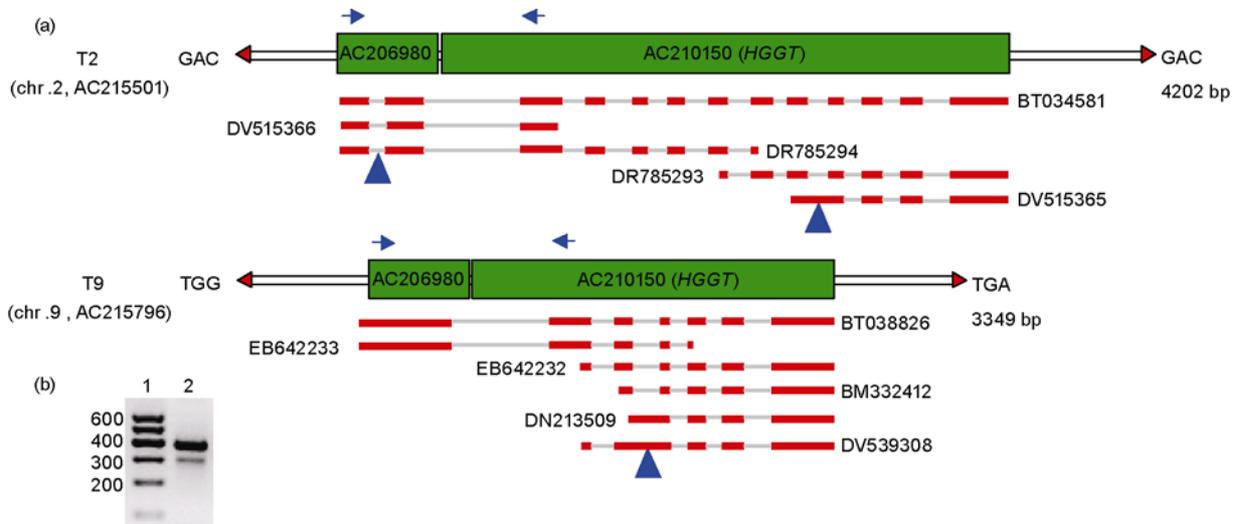
Eleven of the other 66 TEs carrying gene fragments were represented by ESTs (Table 3), of which three had only one EST, and eight had two or more ESTs. ES2-67 was represented by 13 ESTs. The number of corresponding ESTs may reflect the level of transcription. Similar to T2 and T9, alternative splicing phenomenon was observed in 7 of the 8 TEs which were represented by at least 2 ESTs (Table 3). In addition, three TEs, (ES2-115, ES2-142, ES2-154), each associated with the sequences of two genes, are capable of transcription, however, no chimeric transcripts were identified. This lack of chimeric transcripts may be due to the limited number of ESTs in the database.

### 3 Discussion

#### 3.1 *CACTA* elements carrying gene segments exist in the maize genome

In this study, we found that parts of the *HGGT* gene had

been acquired by 3 *CACTA* elements. Compared with other non-autonomous TEs derived from autonomous elements, the 3 *CACTA* elements lacked any remnants of transposase, but the characteristic TIRs and TSD were retained. The internal sequences were composed of identifiable fragments from at least 2 different genomic regions, the *HGGT* gene and sequences from either chromosome 2 or chromosome 4. The chimeric ESTs in the database, the PCR product with the genomic DNA and the RT-PCR product with the cDNA provide strong evidence that TEs carrying gene segments from different genomic regions exist in the maize genome and are not artefacts of sequence assembly errors. The identification of an additional 66 TEs distributed across the 10 chromosomes of maize by querying the sequences of 16,531 maize BACs suggests that TEs carrying gene fragments widely occur in the maize genome. With an average length of 154 kb per BAC<sup>[38]</sup>, a total of 2546 Mb of sequence (the maize genome is estimated to be 2300–2700 Mb) was used in the whole-genome scan. Considering possible BAC overlaps, the genome coverage may be overestimated, thus, the 69 TEs obtained in this study represent a conserved estimate of the expected number of TEs in the maize genome. As the TEs contain both exons and introns, gene fragments capture is likely to involve the acquisition of genomic DNA rather than cDNA copies. The ability of TEs to capture gene seg-



**Figure 7** Transcriptional activities of transposable elements. (a) Red triangles represent TIRs, green boxes indicate the gene segments captured by TEs. Exons are shown as red lines, while introns are indicated by grey lines. The GenBank accession number of each EST is also shown. Alternatively spliced introns in different ESTs are indicated by blue triangles. The blue arrows indicate primer pair 3. (b) RT-PCR products of primer pair 3 amplified in the immature ears of inbred zong3 (lane2). Lane1, DNA marker. The 400 and 300 bp bands correspond to T9 and T2 transcripts, respectively.

**Table 3** Transposable elements with transcriptional activities

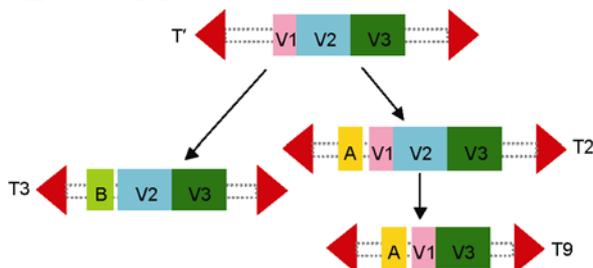
ID	BAC <sup>a)</sup>	Chromosome	Position <sup>b)</sup>	Length (bp)	Copy	EST <sup>c)</sup>	AS <sup>d)</sup>	E value <sup>e)</sup>	Annotation
ES2-59	AC190847	NA	96586–99784	3199	1	<b>EE332027</b>	n	1.00×10 <sup>-34</sup>	peroxisomal 3-ketoacyl-CoA thiolase 2 precursor
ES2-116	AC209072	5	22161–26344	4184	3	CD941388	n	4.00×10 <sup>-22</sup>	HAD-superfamily hydrolase
ES2-52	AC187920	NA	72752–76143	3392	1	EE047636/DN216834	y	2.00×10 <sup>-21</sup>	cellulose synthase catalytic subunit
ES2-36	AC210796	2	156116–158925	2810	1	DV174513/EC875310/EE028235 /EB159347/DY623477/CF636757 /DN230273/DN229640	y	2.00×10 <sup>-25</sup>	putative auxin efflux carrier
ES2-90	AC194109	8	87232–91228	3997	1	DN560136/DV174286/DT948532 /DY686177/DT948531/CO533420 / <b>CO533419</b> /CD996297/DV174287	y	5.00×10 <sup>-10</sup>	26S proteasome regulatory particle non-ATPase subunit12
ES2-67	AC215995	9	81662–85106	3445	3	EE287615/ <b>BM080212</b> /BE511517 /EC894882/EC894881/DR811361 /CD670441/ <b>DY400559</b> /BQ667757 /BE639892/CA830085/BE453784 /DW864676	y	2.00×10 <sup>-18</sup>	KH domain-containing protein-like
ES2-32	AC213592	3	97729–100343	2615	2	EC894675/DN218861/DR958808 /CD650968/DV173879/DN219854 /DW933796/CD650773/DV173878	y	8.00×10 <sup>-22</sup>	importin alpha 1b
ES2-49	AC193518	2	71305–74396	3092	2	EE177428/DN216421/ <b>BM333663</b> / <b>DN226876</b> /DR827220/DR827219	y	8.00×10 <sup>-22</sup>	importin alpha 1b
ES2-115 <sup>b)</sup>	AC149034	NA	110312–114443	4132	1	BM429047/DY398718	n	2.00×10 <sup>-19</sup> 1.00×10 <sup>-17</sup>	DEAD-box ATP-dependent RNA helicase Os03g0276600
ES2-142 <sup>b)</sup>	AC186590	6	23277–28334	5058	3	DR801744/DV511238 /DR801743/DR831406	y	1.00×10 <sup>-50</sup> 3.00×10 <sup>-24</sup>	putative protein kinase AtSIK putative 3-ketoacyl-CoA thiolase
ES2-154 <sup>b)</sup>	AC212442	4	171638–178544	6907	2	CD960507	n	3.00×10 <sup>-16</sup> 2.00×10 <sup>-29</sup>	U3 snoRNA-associated protein 11 one repeat MYB transcriptional factor

a) Accession number of the BAC in GenBank. b) The start and end position of the TEs in the BAC. c) Accession number of ESTs in GenBank, ESTs with a polyA tail were bold-faced, ESTs which could be translated into functional protein in one of the six ORFs were italic. d) AS, alternative splicing; y, yes; n, no. e) E value of annotation. f) TEs with two protein hits.

ments is remarkable, for example, the complete sequences from intron 2 to the 3' UTR of the progenitor *HGGT* gene had been captured by T2.

### 3.2 *CACTA* elements and gene amplification

The *CACTA* elements carrying fragments of the *HGGT* gene had three copies in the maize genome. Three possible explanations could be proposed for this phenomenon: firstly, genome duplication; secondly, the insertion of *CACTA* elements into other TEs<sup>[39]</sup>; thirdly, transposition driven by transposase. The fact that the TSD of the 3 TEs is different excludes the first two possibilities. Under the third conditions, two different processes could be postulated for the amplification of the original progenitor transposon T' (Figure 8): First, T' acquired sequence from intron 2 to the 3' UTR of the progenitor *HGGT* gene (Figure 8, V1, V2 and V3; Figure 2). It was then transposed twice. At one time, T' acquired sequences from BAC-AC187813 (B) located on chromosome 4 with the loss of V1 (intron 2 to intron 5) to form T3. At another time, T' acquired sequences from BAC-AC206980 (A) located on chromosome 2 to form T2, which then transposed to a new location losing V2 (intron 6 to intron 11) to form T9. Under the second situation, the formation of T2 and T9 are similar to the process above, the difference is that the formation of T3 is through direct acquisition of sequences from the progenitor *HGGT* gene (V2 and V3) and is not dependent on T'. The fact that the 3' region was almost split from the same position in the 3 TEs (Figure 2) gives weight to the first situation.



**Figure 8** Schematic diagram of amplification of T2, T3 and T9. TIRs are indicated by red filled triangles. V1, V2 and V3 represent stretches of sequences from progenitor *HGGT* gene. Sequences from other regions are indicated by B and A.

In addition to T2, T3 and T9, other TEs with multiple copies were also identified as having a distinct TSD. Therefore, it is quite possible that these TEs had transposition activity in the past. All the TEs identified in this study belong to the family of ENSPM2\_ZM elements,

which cannot encode transposase, and thus can only be transposed in the presence of corresponding autonomous elements. No such autonomous elements have been reported previously. Some ES2 TEs contained sequences that can be annotated to transposase *Doppia*<sup>[37]</sup>, implying that *Doppia* could be the autonomous elements. However, the difference in TIRs and the supposed transposase binding site-STR<sup>[11]</sup> between ES2 TEs and *Doppia* elements clearly lead us to reject this proposal. Like other cellular gene fragments, the *Doppia* fragments might have been captured by ES2 TEs. The fact that TEs with *Doppia* fragments and TEs with cellular gene fragments could be grouped together into a branch strongly supports this inference (Figure 6). The movement of gene fragments mediated by TEs has a great impact on the host genome. This movement could contribute to the lack of gene collinearity between or within species as found in *Helitron* elements<sup>[25,26]</sup>. The high similarity between the gene fragments and the progenitor host genes caused great difficulties in gene cloning and functional analysis.

### 3.3 Gene fragments captured by *CACTA* elements and pseudogenes

Of the 69 TEs that contained sequences annotated to cellular proteins, 13 (18.8%), are represented by ESTs in the queried databases. In addition, the fact that T2 and T9 carry sequences from two different chromosomal loci and that they could be expressed as chimeric transcripts is a circumstance that aids the origin of novel genes. Jiang et al.<sup>[18]</sup> analyzed over 3000 Pack-MULEs in rice and concluded that more than 10% of them might have been functionally constrained. By contrast, the study by Juretic et al.<sup>[20]</sup> suggested that virtually all TEs contain pseudogenic features such as frameshifts and premature stop codons. In this study, all the ESTs have features preventing them from encoding functional proteins. However, this does not imply that there was no function. The transcripts could regulate expression of the progenitor gene either at the RNA level by RNAi, or at the protein level by competing with the original functional protein for substrates<sup>[20]</sup>. Further research on TEs carrying gene fragments should provide more information and help to better understand genome structure, evolution and function in maize.

*The authors thank Drs. George Mahuku and Trushar Shah for helping to improve the English writing.*

- 1 Bennetzen J L. Transposable element contributions to plant gene and genome evolution. *Plant Mol Biol*, 2000, 42: 251–269
- 2 Bennetzen J L. Transposable elements, gene creation and genome rearrangement in flowering plants. *Curr Opin Genet Dev*, 2005, 15: 621–627
- 3 Wessler S R. Phenotypic diversity mediated by the maize transposable elements *Ac* and *Spm*. *Science*, 1988, 242(4877): 399–405
- 4 Feschotte C, Jiang N, Wessler S R. Plant transposable elements: where genetics meets genomics. *Nat Rev Genet*, 2002, 3: 329–341
- 5 Peterson P A. A mutable pale green locus in maize. *Genetics*, 1953, 38: 682–683
- 6 McClintock B. Mutations in maize and chromosomal aberrations in *Neurospora*. *Carnegie Inst Washington Year Book*, 1954, 53: 254–260
- 7 Chopra S, Brendel V, Zhang J, et al. Molecular characterization of a mutable pigmentation phenotype and isolation of the first active transposable element from *Sorghum bicolor*. *Proc Natl Acad Sci USA*, 1999, 96(26): 15330–15335
- 8 Motohashi R, Ohtsubo E, Ohtsubo H. Identification of Tnr3, a suppressor-mutator/enhancer-like transposable element from rice. *Mol Gen Genet*, 1996, 250(2): 148–152
- 9 Zabala G, Vodkin L O. The *wp* mutation of *Glycine max* carries a gene-fragment-rich transposon of the CACTA superfamily. *Plant Cell*, 2005, 17: 2619–2632
- 10 Tian P F. Progress in plant CACTA elements. *Acta Genetica Sinica*, 2006, 33(9): 765–774
- 11 Gierl A, Lütticke S, Saedler H. TnpA product encoded by the transposable element *En-1* of *Zea mays* is a DNA binding protein. *EMBO J*, 1988, 7(13): 4045–4053
- 12 Peterson P A. A relationship between the *Spm* and *En* control systems in maize. *Am Nat*, 1965, 99(908): 391–398
- 13 Pereira A, Schwarz-Sommer Z, Gierl A, et al. Genetic and molecular analysis of the Enhancer (*En*) transposable element system of *Zea mays*. *EMBO J*, 1985, 4(1): 17–23
- 14 Gierl A, Schwarz-Sommer Z, Saedler H. Molecular interactions between the components of the *En-I* transposable element system of *Zea mays*. *EMBO J*, 1985, 4(3): 579–583
- 15 Pereira A, Cuyppers H, Gierl A, et al. Molecular analysis of the *En/Spm* transposable element system of *Zea mays*. *EMBO J*, 1986, 5(5): 835–841
- 16 Masson P, Surosky R, Kingsbury J A, et al. Genetic and molecular analysis of the *Spm*-dependent *a-m2* alleles of the maize *a* locus. *Genetics*, 1987, 177: 117–137
- 17 Talbert L E, Chandler V L. Characterization of a highly conserved sequence related to mutator transposable elements in maize. *Mol Biol Evol*, 1988, 5(5): 519–529
- 18 Jiang N, Bao Z, Zhang X, et al. Pack-MULE transposable elements mediate gene evolution in plants. *Nature*, 2004, 431: 569–573
- 19 Ohtsu K, Hirano H Y, Tsutsumi N, et al. *Anaconda*, a new class of transposon belonging to the *Mu* superfamily, has diversified by acquiring host genes during rice evolution. *Mol Genet Genomics*, 2005, 274: 606–615
- 20 Juretic N, Hoen D R, Huynh M L, et al. The evolutionary fate of MULE-mediated duplications of host gene fragments in rice. *Genome Res*, 2005, 15: 1292–1297
- 21 Hoen D R, Park K C, Elrouby N, et al. Transposon-mediated expansion and diversification of a family of *ULP*-like genes. *Mol Biol Evol*, 2006, 23(6): 1254–1268
- 22 Holligan D, Zhang X, Jiang N, et al. The transposable element landscape of the model legume *Lotus japonicus*. *Genetics*, 2006, 174: 2215–2228
- 23 Leeuwen H V, Monfort A, Puigdomenech P. *Mutator*-like elements identified in melon, Arabidopsis and rice contain ULPI protease domains. *Mol Genet Genomics*, 2007, 277: 357–364
- 24 Gupta S, Gallavotti A, Stryker G A, et al. A novel class of *Helitron*-related transposable elements in maize contain portions of multiple pseudogenes. *Plant Mol Biol*, 2005, 57: 115–127
- 25 Lai J, Li Y, Messing J, et al. Gene movement by *Helitron* transposons contributes to the haplotype variability of maize. *Proc Natl Acad Sci USA*, 2005, 102: 9068–9073
- 26 Morgante M, Brunner S, Pea G, et al. Gene duplication and exon shuffling by *helitron*-like transposons generate intraspecies diversity in maize. *Nat Genet*, 2005, 37(9): 997–1002
- 27 Takahashi S, Inagaki Y, Satoh H, et al. Capture of a genomic *HMG* domain sequence by the *En/Spm*-related transposable element *Tpn1* in the Japanese morning glory. *Mol Gen Genet*, 1999, 261: 447–451
- 28 Kawasaki S, Nitasaka E. Characterization of *Tpn1* family in the Japanese morning glory: *En/Spm*-related transposable elements capturing host genes. *Plant Cell Physiol*, 2004, 45(7): 933–944
- 29 Roccaro M, Li Y, Sommer H, et al. ROSINA (RSI) is part of a CACTA transposable element, *TamRSI*, and links flower development to transposon activity. *Mol Genet Genomics*, 2007, 278: 243–254
- 30 Kohany O, Gentles A J, Hankus L, et al. Annotation, submission and screening of repetitive elements in Repbase: RepbaseSubmitter and Censor. *BMC Bioinformatics*, 2006, 25(7): 474
- 31 Edgar R C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucl Acid Res*, 2004, 32(5): 1792–1797
- 32 Hall T A. BioEdit: a user-friendly biological sequence alignment editor and analysis program for Windows 95/98/NT. *Nucl Acids Symp*, 1999, 41: 95–98
- 33 Felsenstein J. PHYLIP (Phylogeny Inference Package), version 3.6. (Department of Genetics, University of Washington, Seattle, 2000)
- 34 Kumar S, Tamura K, Nei M. MEGA3: integrated software for molecular evolutionary genetics analysis and sequence alignment. *Brief Bioinform*, 2004, 5(2): 150–163
- 35 Cahoon E B, Hall S E, Ripp K G, et al. Metabolic redesign of vitamin E biosynthesis in plants for tocotrienol production and increased antioxidant content. *Nat Biotechnol*, 2003, 21: 1082–1087
- 36 Frey M, Reinecke J, Grant S, et al. Excision of the *En/Spm* transposable element of *Zea mays* requires two element-encoded proteins. *EMBO J*, 1990, 9(12): 4037–4044
- 37 Bercury S D, Panavas T, Irenze K, et al. Molecular analysis of the *Doppia* transposable element of maize. *Plant Mol Biol*, 2001, 47: 341–351
- 38 Yim Y S, Davis G L, Duru N A, et al. Characterization of three maize bacterial artificial chromosome libraries toward anchoring of the physical map to the genetic map using high-density bacterial artificial chromosome filter hybridization. *Plant Physiol*, 2002, 130: 1686–1696
- 39 Jameson N, Georgelis N, Fouladbash E, et al. *Helitron* mediated amplification of cytochrome P450 monooxygenase gene in maize. *Plant Mol Biol*, 2008, 67: 295–304