

Sequence Organization and Conservation in *sh2/a1*-Homologous Regions of Sorghum and Rice

Mingsheng Chen, Phillip SanMiguel and Jeffrey L. Bennetzen

Genetics Program and Department of Biological Sciences, Purdue University, West Lafayette, Indiana 47907

Manuscript received July 3, 1997

Accepted for publication September 22, 1997

ABSTRACT

Previously, we have demonstrated microcolinearity of gene composition and orientation in *sh2/a1*-homologous regions of the rice, sorghum, and maize genomes. However, the *sh2* and *a1* homologues are only about 20 kb apart in both rice and sorghum, while they are separated by about 140 kb in maize. In order to further define sequence organization and conservation in *sh2/a1*-homologous regions, we have completely sequenced a 42,446-bp segment of sorghum DNA. Four genes were identified: a homologue of *sh2*, two homologues of *a1*, and a putative transcriptional regulatory gene. A solo long terminal repeat of the retroelement *Leviathan* was detected between the two *a1* homologues, and eight miniature inverted repeat transposable elements were found in this region. Comparison of the sorghum sequence with the sequence of the homologous segment from rice indicated that only the identified genes were evolutionarily conserved between these two species, which have evolved independently for over 50 million years. The introns of the *a1* homologues have evolved faster than the introns of the *sh2* homologue. The *a1* tandem duplication appears to be an ancient event that may have preceded the ancestral divergence of maize, sorghum, and rice.

MOST higher plants have complex genomes that contain large amounts of repetitive DNA. In maize, for instance, genes probably make up less than 20% of the nuclear genome. In most plants, the nature and organization of the nongenic regions remain largely unknown.

Rice, sorghum, and maize are three grass species of great agronomic value. They have very different haploid genome sizes; 430 Mbp pairs for rice, 750 Mbp pairs for sorghum, and 2500 Mbp pairs for maize (Arumuganathan and Earle 1991). Sorghum and maize are close relatives, having diverged from a common ancestor about 15–20 million years ago. Rice has undergone independent descent from sorghum and maize for over 50 million years. Rice and sorghum are true diploids, while maize is a segmental allotetraploid (Gaut and Doebley 1997). Despite their many years of independent evolution, maize, sorghum, and rice have retained a very similar gene repertoire (Bennetzen and Freeling 1993). They also have maintained extensive genetic map colinearity (Ahn *et al.* 1993; Hulbert *et al.* 1990; Moore *et al.* 1995; Whitkus *et al.* 1992).

Previously, we demonstrated microcolinearity in the *sh2/a1*-homologous regions of the rice, sorghum, and maize genomes (Chen *et al.* 1997). While the physical distance between *sh2* and *a1* is about 140 kb in maize, it is about 20 kb in both sorghum and rice. The transcriptional orientation of *sh2* and *a1* homologues is conserved

among rice, sorghum, and maize. The *sh2/a1*-homologous region in rice has been completely sequenced (Chen and Bennetzen 1996). Three candidate genes were identified: Sh2, A1, and gene X.¹ Gene X appears to encode a transcription factor, as indicated by putative protein-protein interaction and zinc finger domains. Eight miniature inverted repeat transposable elements (MITEs) were identified in this region.

To characterize the sequence organization of the *sh2/a1*-homologous region in sorghum further and to investigate the nature and evolution of sequences in *sh2/a1*-homologous regions, we have completely sequenced the *sh2/a1*-homologous region in sorghum. Comparison of the sequence in these homologous regions of rice and sorghum reveals that the genes are conserved, but that the intergenic spaces are highly diverged.

MATERIALS AND METHODS

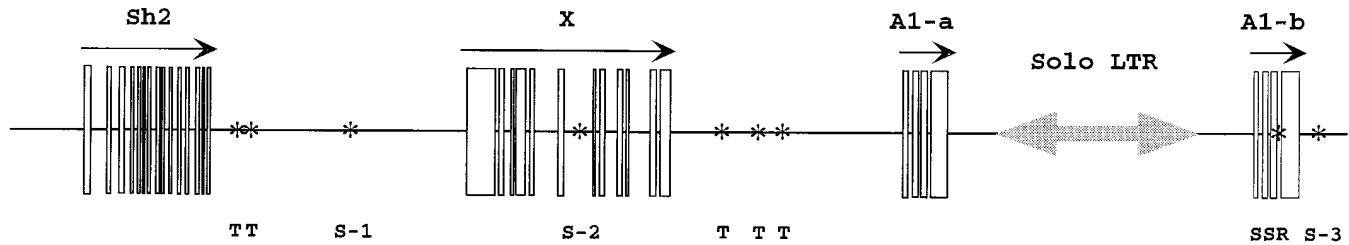
Materials: Genomic DNA was from *Sorghum bicolor* line BTx623. Sorghum bacterial artificial chromosome clones and subclones were as previously described (Chen *et al.* 1997). Restriction enzymes were from Promega (Madison, WI), New England Biolabs (Beverly, MA), and GIBCO (Grand Island, NY). DNA sequencing reagents were from Pharmacia (Piscataway, NJ).

DNA isolation and gel blot hybridization: Genomic DNA isolation, restriction enzyme digestions, agarose gel electro-

¹ Gene symbols referring to known maize genes appear as lower case italics. However, there are no Sh2, A1, or X genes certified in sorghum or rice. Hence, we use only operating names (Sh2, A1-a, A1-b, and X) for these putative loci, without italics. We cannot name sorghum or rice genes with only indirect information.

Corresponding author: Jeffrey L. Bennetzen, Department of Biological Sciences, Purdue University, West Lafayette, IN 47907
E-mail: maize@bilbo.bio.purdue.edu

Sorghum



Rice

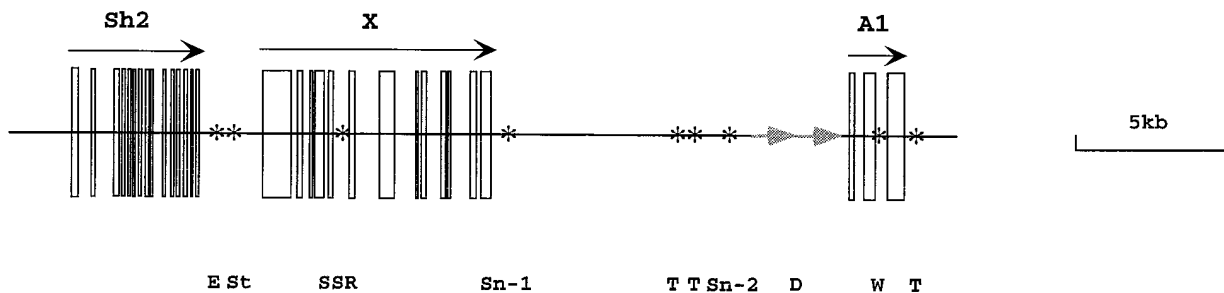


Figure 1.—Maps of *sh2/a1*-homologous regions of sorghum and rice. Lines represent the entire sequenced 42,446 bp of sorghum and 30,034 bp of rice (Chen and Bennetzen 1996). Empty boxes represent putative exons, and arrows above the boxes indicate their transcriptional orientation. The two shaded opposite arrows represent a solo LTR of the retroelement *Leviathan* (Bennetzen 1996) in sorghum. The two direct shaded arrows indicate a tandem duplication of 1432 bp in rice (the orientation of the arrows is arbitrary). Asterisks with letters below them indicate other sequence features. D, duplication; E, *Explorer* element; S-1, S-2, S-3, putative MITEs of sorghum; Sn-1, *Snabo-1*; Sn-2, *Snabo-2*; SSR, simple sequence repeat; St, *Stowaway* element; T, *Tourist* element; W, *Wanderer* element.

phoresis, filter replication, and hybridization were all as previously described (Chen *et al.* 1997).

DNA sequencing: A $\gamma\delta$ transposon-facilitated DNA sequencing strategy was used (Strathmann *et al.* 1991). Cycle sequencing products were run on an express automated sequencer (ALF; Pharmacia [Piscataway, NJ]).

Sequence analysis: All sequence analysis was done using the software package from the University of Wisconsin Genetics Computer Group, Madison (Devereaux *et al.* 1984), through the Purdue University AIDS Center Laboratory for Computational Biology.

Sequence homologies were defined using the Pileup program with a penalty of 3.00 for gaps and of 0.10 for each bp of gap length. The presented nucleotide identities counted each sequence difference as a single change and each gap or insertion as a single change. We also estimated pairwise sequence divergence rates using only nongapped positions and generated highly similar results (data not shown).

RESULTS

The sequence of the *sh2/a1*-homologous region of sorghum: The completed sequence of the *sh2/a1*-homo-

logous region in sorghum (Figure 1) is 42,446 bp (GenBank accession number AF010283). Sequence data redundancy (sum of total raw sequence data divided by length of the region sequenced) was 3.2, and all re-

TABLE 1

Comparison of the *sh2* homologues in sorghum, rice, and maize genomes

			Percent Rice Sh2	Percent Maize <i>sh2</i>
Sorghum Sh2	Nucleotide	Exon	83	95
		Intron	54	88
	Protein	Identity	77	91
			Similarity	87
Maize <i>sh2</i>	Nucleotide	Exon	82	
		Intron	57	
	Protein	Identity	77	
				Similarity

Zmash2	M Q F A L A L D T N	S G P H Q I . R S C	E G D G I D R L . E	K L S I G G R K Q E	38
Sbish2	M Q F S L A S D A N	S G P H P I R R S C	E G G G I D R L . E	R L S I G G S K Q E	39
Osash2	M Q F M M P L D T N	A C A Q P M R R A G	E G A G T E R L M E	R L N I G G M T Q E	40
Zmash2	K A L R N R C F G G	R V A A T T Q C I L	T S D A C P E T L H	S Q T Q S S R K N Y	78
Sbish2	K A L R N R C F G G	R V A A T T Q C I L	T S D A C P E T L H	F Q T Q S S R K S Y	79
Osash2	K A L R K R C F G D	G V T G T A R C V F	T S D A D R D T P H	L R T Q S S R K N Y	80
Zmash2	A D A N R V S A I I	L G G G T G S Q L F	P L T S T R A T P A	V P V G G C Y R L I	118
Sbish2	A D A N H V S A I I	L G G G T G S Q L F	P L T S T R A T P A	V P V G G C Y R L I	119
Osash2	A D A S H V S A V I	L G G G T G V Q L F	P L T S T R A T P A	V P V G G C Y R L I	120
Zmash2	D I P M S N C F N S	G I N K I F V M S Q	F N S T S L N R H I	H R T Y L E G G I N	158
Sbish2	D I P M S N C F N S	G I N K I F V M T Q	F N S T S L N R H I	H R T Y L G G E I N	159
Osash2	D I P M S N C F N S	G I N K I F V M T Q	F N S A S L N R H I	H H T Y L G G G I N	160
Zmash2	F A D G S V Q V L A	A T Q M P E E P A G	W F Q G T A D S I R	K F I W V L E D Y Y	198
Sbish2	F A D G S V Q V L A	D T Q M P E E P D G	W F Q G T A D S V R	K F I W V L E D Y Y	199
Osash2	F T D G S V Q V L A	A T Q M P D E P A G	W F Q G T A D A I R	K F M W I L E D H Y	200
Zmash2	S H K S I D N I V I	L S G D Q L Y R M N	Y M E L V Q K H V E	D D A D I T I S C A	238
Sbish2	N H K S I E H I V I	L S G D Q L Y Q M N	Y M E L V Q K H V E	D N A D I T V S C A	239
Osash2	N Q N N I E H V V I	L C G D Q L Y R M N	Y M E L V Q K H V D	D N A D I T I S C A	240
Zmash2	P V D E S R A S K N	G L V K I D H T G R	V L Q F F E K P K G	A D L N S M R V E T	278
Sbish2	P V D E S R A S N N	G L V K C D H T G R	V L Q F F E K P K G	A D L N S M R V D T	279
Osash2	P I D G S R A S D Y	G L V K F D D S G R	V I Q F L E K P E G	A D L E S M K V D T	280
Zmash2	N F L S Y A I D D A	Q K Y P Y L A S M G	I Y V F K K D A L L	D L L K S K Y T Q L	318
Sbish2	N F L S Y A I G D A	Q K Y Q Y I A S M G	I Y V F K K D A L L	D L L K S K Y T Q L	319
Osash2	S F L S Y A I D D K	Q K Y P Y I A S M G	I Y V L K K D V L L	D I L K S K Y A H L	320
Zmash2	H D F G S E I L P R	A V L D H S V Q A C	I F T G Y W E D V G	T I K S F F D A N L	358
Sbish2	H D F G S E I L P R	A V L E H N V Q T C	I F M G Y W E D V G	T I K S F F D A N L	359
Osash2	Q D F G S E I L P R	A V L E H N V K A C	V F T E Y W E D I G	T I K S F F D A N L	360
Zmash2	A L T E Q P S K F D	F Y D P K T P F F T	A P R C L P P T Q L	D K C K M K Y A F I	398
Sbish2	A L T E Q P S K F D	F Y D P K T P F F T	A P R Y L P P T Q L	D K C K I K D A S I	399
Osash2	A L T E Q P P K F E	F Y D P K T P F F T	S P R Y L P P A R L	E K C K I K D A I I	400
Zmash2	S D G C L L R E C N	I E H S V I G V C S	R V S S G C E L K D	S V M M G A D T Y E	438
Sbish2	S D G C L L R E C S	I E H S V I G V C S	R V S Y G C E L K D	C V M M G A D I Y E	439
Osash2	S D G C S F S E C T	I E H S V I G I S S	R V S I G C E L K D	T M M M G A D Q Y E	440
Zmash2	T E E E A S K L L L	A G K V P V G I G R	N T K I R N C I I D	M N A R I G K N V V	478
Sbish2	T E E E A S K L L L	A G E V P V G I G G	N T K I R N C I I D	I N A R I G K N V V	479
Osash2	T E E E T S K L L F	E G K V P I G I G E	N T K I R N C I I D	M N A R I G R N V I	480
Zmash2	I T N S K G I Q E A	D H P E E G Y Y I R	S G I V V I L K N A	T I N D G S V I * 516	
Sbish2	I T N S K G I Q E A	D H P E E G Y Y I K	S G I V V I L K N A	T I K D G S V I * 517	
Osash2	I A N T Q G V Q E S	D H P E E G Y Y I R	S G I V V I L K N A	T I K D G T V I * 518	

Figure 2.—Comparison of predicted amino acid sequence encoded by *sh2* homologues of sorghum, rice, and maize. Asterisks represent stop codons. Osash2, *sh2* homologue of rice (Chen and Bennetzen 1996); Sbish2, *sh2* homologue of sorghum; Zmash2, *sh2* homologue of maize (Hannah and Shaw 1992).

Zmaal1 M E	R G A G A S E K G T	V L V T G A S G F V	G S W L V M K L L Q	32
Sbial1-b	M N G G A S V K G P	V V V T G A S G F V	G S W L V M K L L Q	30
Sbial1-a	M G E V V A T W E A	T E G G A G V K G P	V V V T G A S G F L	G S W L V M K L L Q	40
Osaal1 M G E A V K G P	V V V T G A S G F V	G S W L V M K L L Q	28
Zmaal1	A G Y T V R A T V R	D P A N V G K T K P	L M D L P G A T E R	L S I W K A D L A E	72
Sbial1-b	A G Y T V R A T V R	D P A N V V K T K P	L L D L P G A T E R	L S L W K A D L A V	70
Sbial1-a	A G Y T V R A T V R	D P A N V V K T K P	L L D L P G A T E R	L S L W K A D L A D	80
Osaal1	A G Y T V R A T V R	D P S N V G K T K P	L L E L A G S K E R	L T L W K A D L G E	68
Zmaal1	E G S F H D A I R G	C T G V F H V A T P	M D F L S K D P E N	E V I K P T V E G M	112
Sbial1-b	E G S F D D A I R G	C T G V F H V A T P	M D F E S K D P E N	E V I K P T V E G M	110
Sbial1-a	E G S F D D A I R G	C T G V F H V A T P	M D F E S K D P E N	E V I K P T V E G M	120
Osaal1	E G S F D A A I R G	C T G V F H V A T P	M D F E S E D P E N	E V I K P T V E G M	108
Zmaal1	I S I M R A C K E A	G T V R R I V F T S	S A G T V N L E E R	Q R P V Y D E E S W	152
Sbial1-b	I S I M R A C K E A	G T V R R I V F T S	S A G T C N I E E W	R K P V Y D E D N W	150
Sbial1-a	M S I M R A C K E A	G T V R R I V F T S	S A G T V N I E E R	Q R P V Y D Q D N W	160
Osaal1	L S I M R A C R D A	G T V K R I V F T S	S A G T V N I E E R	Q R P S Y D H D D W	148
Zmaal1	T D V D F C R R V K	M T G W M Y F V S K	T L A E K A A L A Y	A A E H G L D L V T	192
Sbial1-b	T D V D F C R R V K	M T G W M Y F V S K	T L A E K A A L A Y	A A E H G M E L V T	190
Sbial1-a	S D V D F C Q R V K	M T G W M Y F V S K	S L A E K A A M A Y	A A E H G L D F I S	200
Osaal1	S D I D F C R R V K	M T G W M Y F V S K	S L A E K A A M E Y	A R E H G L D L I S	188
Zmaal1	I I P T L V V G P F	I S A S M P P S L I	T A L A L I T G N A	P H Y S I L K Q V Q	232
Sbial1-b	I I T T L V V G P F	L S T G M P P S M I	T R L A L V T G N E	A H Y S I L K Q V Q	230
Sbial1-a	I I P T L V V G P F	L S A G M P P S L I	T A L A L V T G N E	A H Y S I L K Q V Q	240
Osaal1	V I P T L V V G P F	I S N G M P P S H V	T A L A L L T G N E	A H Y S I L K Q V Q	228
Zmaal1	L I H L D D L C D A	E I F L F E N P A A	A G R Y V C S S H D	V T I H G L A A M L	272
Sbial1-b	F V H L D D L C D A	H I F L F E H P A A	A G R Y V C S S C D	T T I H D L A A M L	270
Sbial1-a	F V H L D D L C D A	H L F L F E H P A A	A G R Y V C S S H D	A T I H G L A A M L	280
Osaal1	F V H L D D L C D A	E I F L F E S P E A	R G R Y V C S S H D	A T I H G L A T M L	268
Zmaal1	R D R Y P E Y D V P	Q R F P . . . G I Q	D D L Q P V R F S S	K K L Q D L G F T F	309
Sbial1-b	R D R Y P E Y D I P	E R F P A G T G I E	D D L Q M V H M S A	K K L Q D L G F T F	310
Sbial1-a	R D R Y P E Y D I P	E R F P . . . G I E	D D L Q P V H F S S	K K L L D H G F T F	317
Osaal1	A D M F P E Y D V P	R S F P G I D A . .	D H L Q P V H F S S	W K L L A H G F R F	306
Zmaal1	R Y K T L E D M F D	A A I R T C Q E K G	L I P L	A T A A G G D G F A	343
Sbial1-b	R Y T R M E D M Y D	D A I R T C R E K G	L I P L	A A A G R D D G S A	344
Sbial1-a	K Y T . V E D M F D	A A I R M C R E K G	L I P L	A T A G G G R A L P	350
Osaal1	R Y T . L E D M F E	A A V R T C R E K G	L L P P L P P P P T	T A V A G G D G S A	345
Zmaal1	S V R A P G E T E A	T I G A * 357	
Sbial1-b	S V R A P G E R D V	T A T A G G D V S A	P V R A P G G E T D	V T I G A * 379	
Sbial1-a	* 350	
Osaal1	G V A G E K E P I L	G R G T G T A V G A	E T E A L V K * 372	

Figure 3.—Comparison of predicted amino acid sequence encoded by *a1* homologues of sorghum, rice, and maize. Asterisks represent stop codons. Dotted lines represent gaps. Sbial1-a, A1-a homologue of sorghum; Sbial1-b, A1-b homologue of sorghum; Osaal1, *a1* homologue of rice (Chen and Bennetzen 1996); Zmaal1, *a1* homologue of maize (Schwarz-Sommer *et al.* 1987).

gions were sequenced on both strands.

On the basis of comparisons to the maize and rice sequences, this region contains four genes; an *sh2* homologue, a gene homologue to gene X of rice (Chen

and Bennetzen 1996), and two *a1* homologues. The *sh2* homologue in sorghum appears to have 15 protein-encoding exons. The two *a1* homologues have four exons each. The distance between the 3' end of the *sh2*

homologue and the 5' end of the closer *a1* homologue is 22,025 bp. Gene X was identified between the *sh2* homologue and the nearer *a1* homologue. It appears to have 12 protein-encoding exons. A solo long terminal repeat (LTR) of the retroelement *Leviathan* (Bennetzen 1996) was identified between the two *a1* homologues. Five *Tourists* were observed and three putative MITEs were also detected. A simple sequence repeat (SSR) with 21 ATs is present in the third intron of the A1-b homologue.

Comparison of *sh2* homologues in sorghum, rice, and maize: The *sh2* homologue in sorghum apparently encodes 517 amino acids, compared to 516 and 518 amino acids in maize and rice, respectively. Sequence identity between the *sh2* homologues in sorghum, rice, and maize is listed in Table 1. An amino acid sequence comparison is shown in Figure 2.

As expected, protein-encoding portions of exons are significantly more conserved than introns within these exons. The *sh2* homologue in sorghum is more homologous to the *sh2* homologue in maize than to the *sh2* homologue in rice for both exons and introns. The degree of homology between the *sh2* homologues of sorghum and rice is very similar to the degree of homology between the *sh2* homologues of maize and rice.

Comparison of *a1* homologues in sorghum, rice, and maize: The A1-a and A1-b genes in sorghum apparently encode peptides of 350 and 379 amino acids, respectively, compared to orthologous 357 and 372 amino acid peptides in maize and rice. These substantial differences in predicted protein size are primarily due to variations in stop codon location at the C terminus of each putative peptide (Figure 3).

Sequence comparisons of the *a1* homologues in the sorghum, rice, and maize genomes are displayed in Table 2. An amino acid sequence alignment is shown in

Figure 3. As found previously, exons are more conserved than introns. The sequence homology is higher between sorghum and maize than between sorghum and rice or maize and rice.

Each of the A1-a and A1-b sorghum genes has four exons, the same number as in maize, but one more than in rice. The *Wanderer* element located in the second intron of the rice *a1* homologue is not present in either sorghum *a1* homologue.

Gene X between the *sh2* and *a1* homologues: Gene X was identified between the *sh2* and *a1* homologues by comparison to the sequences of the *sh2/a1*-homologous region in rice. Gene X has 12 exons in sorghum. It is predicted to encode an 895 amino acid protein, compared to 1070 amino acids in rice gene X. The putative zinc finger motif identified in the 3' end of the rice version is not present in the sorghum sequence. At the nucleotide level, gene X of sorghum and gene X of rice share, 82% and 57%, respectively, identity of exons and introns. At the protein level, they are 79% identical and 85% similar (Figure 4). The SSRs in the fifth intron of gene X in rice are not present in the gene X sequence in sorghum. A putative MITE is positioned in the sixth intron of sorghum gene X, but it is absent in the rice gene X sequence.

Sequence comparison of *sh2/a1*-homologous regions in the rice and sorghum genomes: We compared the entire sequenced regions in rice and sorghum to define other conserved components in the *sh2/a1*-homologous regions. A window size of 50 and a stringency of 40 were used to reveal homologies $\geq 80\%$ (Figure 5). By this criterion, only the protein-encoding portions of the *sh2* homologues, gene X, and the *a1* homologues were detected.

Tandem duplication of the *a1* homologue in sorghum: The *a1* homologue was directly duplicated in the *sh2/a1*-homologous region of sorghum. The distance between the putative stop codon of A1-a and the presumptive start codon of A1-b is 9756 bp. To characterize the features of the duplication, we made a comparison between these two *a1* homologues. A window size of 50 and a stringency of 33.3 (66.7% homology) were used to analyze the sequences (Figure 6). The comparison indicated that only the *a1* coding regions are highly homologous. The 5' end and the 3' end are not homologous, except that there are an 89-bp stretch with 72% identity on the 5' end of each *a1* homologue (201–289 bp upstream of the putative A1-a start codon and 281–357 bp upstream of the apparent A1-b start codon) and a 94-bp stretch with 80% identity on the 3' end of each *a1* homologue (82–175 bp downstream of the putative A1-a stop codon and 184–270 bp downstream of the presumptive A1-b stop codon). Further analysis revealed that the 5' and 3' end homologies are conserved in the respective 5' and 3' termini of the *a1* gene in maize. Putative TATA and CAAT boxes are located within the 5' end homologies, and potential poly-

TABLE 2
Comparison of the *a1* homologues in rice, maize,
and sorghum genomes

			Percent Rice A1	Percent Maize a1	Percent Sorghum A1-b
Sorghum A1-a	Nucleotide	Exon	84	91	92
		Intron	52	55	62
	Protein	Identity	80	85	87
		Similarity	89	92	92
Sorghum A1-b	Nucleotide	Exon	78	90	
		Intron	54	69	
	Protein	Identity	72	83	
		Similarity	81	90	
Maize <i>a1</i>	Nucleotide	Exon	82		
		Intron	52		
	Protein	Identity	75		
		Similarity	83		

Sbigenex M S H S D E D S E I S D S E I D E Y A M S G G L K V R N N G E N Y S C L F C S S K K K N N Y S K S S L V Q H A S G V S A A P N R K A K	48
Osagenex D K F Y A R L V A G E F K V K D G Q S G L K V R N N G E N Y S C F F C S G K K K K D F N I N N L I Q H A S G V C A A S N R Q A K	74
Sbigenex	E K A H R A L F K Y L K N D L A K S S E P Q P L V I P F V E P Q P L Q N R D E K F V W P W M G I L V N V P T E W K D G R Q I G F S G N R L K E Q L S H	123
Osagenex	D K A T H R A L A K H L K N G L T K S S G Q Q S Q T A V E G Q Q S Q T A V E P Q P L F N R D E K F V W P W M G V L V N V P T E W K D G R Q I G R S G N H L K E Q L S R	149
Sbigenex	F C P L K V I P L W T F R G H T G N A I V F F G K D W M G F R N A R T F E S H F A A G G F G K K D W T G K K N Q G S E L Y G W L A R A E D Y N S P G I	198
Osagenex	F C P L K I I P L W N F R G H S G N A I V E F G K D W H G F R N A L A F E D Y F C K E G Y G K R D W K E K Q N Q G S H L F G W V A R A E D H T S P G L	224
Sbigenex	I A D Y L R K N G D L K S V N D L A K E G A R K T D R L V A N L A N Q I E V K N R Y L Q E L E S K Y S E T T A S L E K M M G O R E O L L Q S Y N E E I	273
Osagenex	I G D H L R K N G D L K T I N D L E N E G A R K T D K L V A N L A N Q I E V K N R H L Q E L E V T Y N E R T T S L E K M M G O R E O L L Q K Y N E E I	299
Sbigenex	S K M Q Q L A R R H R K M Q Q L A Q R H S Q K V I D E N Q K S Q K I I D E N Q K L R S E L E A K M N D L D V R S K Q L D E L A A K S D Y D R R N L E Q E K Q K N T L E Q Q	348
Osagenex	R K M Q Q L A Q R H R K M Q Q L A Q R H S Q K I I D E N Q K S Q K I I D E N Q K L R S E L E S K M S E L N T R S K E L D E I A A K S D Y D R R I I D Q E K Q K N T L E Q E	374
Sbigenex	K A D E N V L K L R R A D E N V L K L R E K H A A L K K I L E K E A A V K K I L M L E Q Q L D A K Q K L E L E I Q Q L K K L E L D I Q Q L K G K L K V M E H M P G D E D S A S K N K I N E L S E A L Q E K I D E L	423
Osagenex	R A D E N V L K L R R A D E N V L K L R E K E A A V K K I L E K E A A V K K I L M L E Q Q V D A K Q K L E L D I Q Q L K K L E L D I Q Q L K G K L E V M K H M P G D E D S A L L K N K I D E L S E E L Q E K M D E L	449
Sbigenex	D G M E S L N Q T L D A M E S L N Q T L V I K E S K S N I E V I K E R K S N T E L Q E A R K E L E N M Q D A R K E L E N V C G O A H I G I K R M G E L D L K A F S K A C Q K E R T E D A F V T A A F L C S K W E	497
Osagenex	D A M E S L N Q T L D A M E S L N Q T L V I K E R K S N T E V I K E R K S N T E M Q D A R K E L E N M Q D A R K E L E N V Y G Q S H I G I K R M G E L D L E A F S K A C R K M S S E E D A E I T A A I L C S K W Q	529
Sbigenex	A E I K N P D W H P A E I K N P D W H P F R V E I I E D D A F R F E I I E D D A K L R A L K E E H G K L Q E L L K E E H G E E I Y A L V T K A L L E I N E Y K S K G S Y P V G E L W N F K E N R K V T L K E A V Q F	572
Osagenex	A E I K N P D W H P A E I K N P D W H P F R F E I I E D D A F R F E I I E D D A K L Q E L L K E E H G K L Q E L L K E E H G E D I Y F L V R D A L L E I N E Y N P S G R F F V G E L W N F K D K R K A T L K E T V Q F	599
Sbigenex	H R S G R Q Q W Q R T T E E T R P A S T T A P G S D E A C G I C R E K F G M G G I C R E K F G M G G W A G A S S D F V N L P C E H A F H A N C V L A W F Y K G N T C P V C S H D V C G Q L L Y W	749
Osagenex	H R S G R Q Q W Q R T T E E T R P A S T T A P G S D E A C G I C R E K F G M G G I C R E K F G M G G W A G A S S D F V N L P C E H A F H A N C V L A W F Y K G N T C P V C S H D V C G Q L L Y W	749
Sbigenex	S E A L K S F L D H S E A L K S F L D H I P V S S V K L D G I P V S S I R L D G S V L G A L D A M Y S V P D A I D S M Y S N A A G A V I V R S G V A G A V I V D V V Q S S F G K Y V D R D I G F V E F S S L V L W A L E I A W L A K	649
Osagenex	S E A L K S F L D H S E A L K S F L D H I P V S S I R L D G I P V S S I R L D G S V P D A I D S M Y S V P D A I D S M Y R S G V A G A V I V R S G V A G A V I V D V V R T S F G K F V D R D I G F V E F P S L V L W A I E I A W L A K	824
Sbigenex	L F L W E P F F P V R F A N V R K P V L V Y S D Q T L A D G V Y S D Q T L A D G L H I L S K E K M G L H I L S K E K I G V A V I D R K T S C V A V I D R K T S C L I G S I Q C S D L Y L F L D D S S L F S K R T T A E D S S P P Q Q N	799
Osagenex	S F L W E P F F P V R F A N T T K P V S V Y S D Q T L A D G V Y S D Q T L A D G L H I L S K E K I G L H I L S K E K I G V A V I D R K T S C V A V I D R K T S C L I G S I Q C S D L Y Q L L D D S S L F R N R N T E N S S A S G G Q N	974
Sbigenex	I L A L R N R Q P S V L S L R T G Q R I M V G L P A T N L K T A G L P V T N R K S D T L K Q A M E K S D T L K Q A M E K L T T S R S S C S F I V D E Q G H V E G V V T T R D I I S V F S P P C M D S R I D G G T F	874
Osagenex	V L S L R T G Q R I M V G L P A T N L K T A G L P V T N R K S D T L K Q A M E K S D T L K Q A M E K L T T S R S S C S F I V D E H G R V E G V V T A R D I I S V F S P P C M D S R I D G G T F	1049
Sbigenex	F S A L E Q A G C F S A A L A Q T G C R V E N G Q M I Q N R V E H G Q M I Q N S * 895 S * 1070	
Osagenex	F S A A L A Q T G C R V E H G Q M I Q N S * 1070	

Figure 4.—Comparison of predicted amino acid sequence encoded by gene X orthologues of sorghum and rice. Asterisks indicate stop codons. Dotted lines represent gaps. The gap in position 572–573 in Sbigenex represents the putative zinc finger motif missing in the sorghum gene X homologue. Osagenex, gene X homologue of rice (Chen and Bennetzen 1996); Sbigenex, gene X homologue of sorghum.

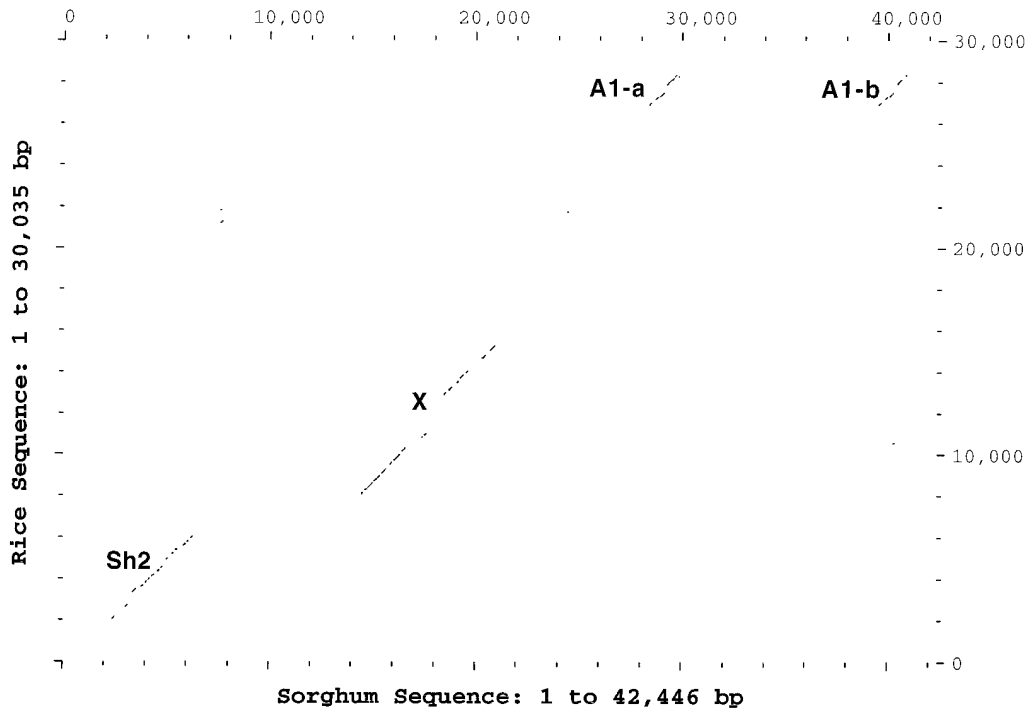


Figure 5.—DNA sequence comparison of *sh2/a1*-homologous regions of rice and sorghum.

adenylation signals are located within the 3' end homologies.

A solo LTR between the two *a1* homologues: A solo LTR of the retroelement *Leviathan* (Bennetzen 1996) was detected between the two *a1* homologues (Figure 1). It is 4479 bp long and flanked by a 5-bp CTACA target site duplication. This solo LTR shares 53% identity with the 5' LTR of the *Leviathan* element described previously (Bennetzen 1996). The solo LTR in our clone might have resulted from an unequal recombination event between the two LTRs of a complete *Leviathan* during the cloning process. To investigate this possibility, we used a 2.5-kb *KpnI/Pml* single-copy DNA fragment between the solo LTR and A1-b as a probe in gel blot hybridization experiments. For every restriction enzyme used (*Bam*HI, *Eco*RI, *Eco*RV, *Hind*III, *Sac*I, and *Xho*I), the hybridization pattern was the same for the sorghum genomic DNA and the sorghum BAC DNA (data not shown). It appears that this solo LTR is present, in this form, between the *a1* homologues in the sorghum genome.

MITEs in the sorghum *sh2/a1* region: Five *Tourist* transposable elements were identified in this region. Two are 3' to the *sh2* homologue, while three are between gene X and A1-a (Figure 1). Three putative MITEs were observed (S-1, S-2, and S-3 in Figure 1). S-1 is located between the *sh2* homologue and gene X. It is 678 bp long, with perfect, 43-bp terminal inverted repeats.

S-2 is positioned in the sixth intron of gene X. It has two homologues in the EMBL/GenBank database. One is in the second intron of the nucleic-acid-binding protein gene of maize (Cook and Walder 1992), while the other is in the 5' end of the 27-kD zein gene of maize

(Das *et al.* 1991). S-2 and the element in the nucleic-acid-binding protein gene have 70% homology over 124 bp. S-2 and the element in the 27-kD zein gene have 61% homology over 136 bp. S-2 can form a fold-back structure with 70% complementarity over 103 bp.

S-3 is positioned in the 3' end of A1-b. It has one homologue in the GenBank database, which is located in the 3' end of the *ltp2* (lipid transfer protein) gene of sorghum (Pelése-Siebenbourg *et al.* 1994). They share 78% identity over 147 bp. S-3 can form a snap-back structure with 78% complementarity over 98 bp.

Target site duplications were not identified for S-1, S-2, or S-3. Terminal repeats (either inverted or direct) were also not identified for S-2 or S-3.

DISCUSSION

We have completely sequenced the *sh2/a1*-homologous region in sorghum. Mapping studies (A. Melake-Berhan and J. L. Bennetzen, unpublished results) position these two loci in an orthologous location to their homologues in maize. Analysis of sequence data identified four putative genes: a homologue of maize *sh2*, a homologue of rice gene X, and two homologues of maize *a1*. The order and orientation of the genes are conserved among sorghum and rice, but whether gene X exists in the orthologous position of maize is not known. Sequence comparison of this region in rice and sorghum revealed that only these genes are conserved. All the other sequence features (such as MITEs, SSRs, and other noncoding regions) are not recognizably conserved.

Maize and sorghum are close relatives despite their

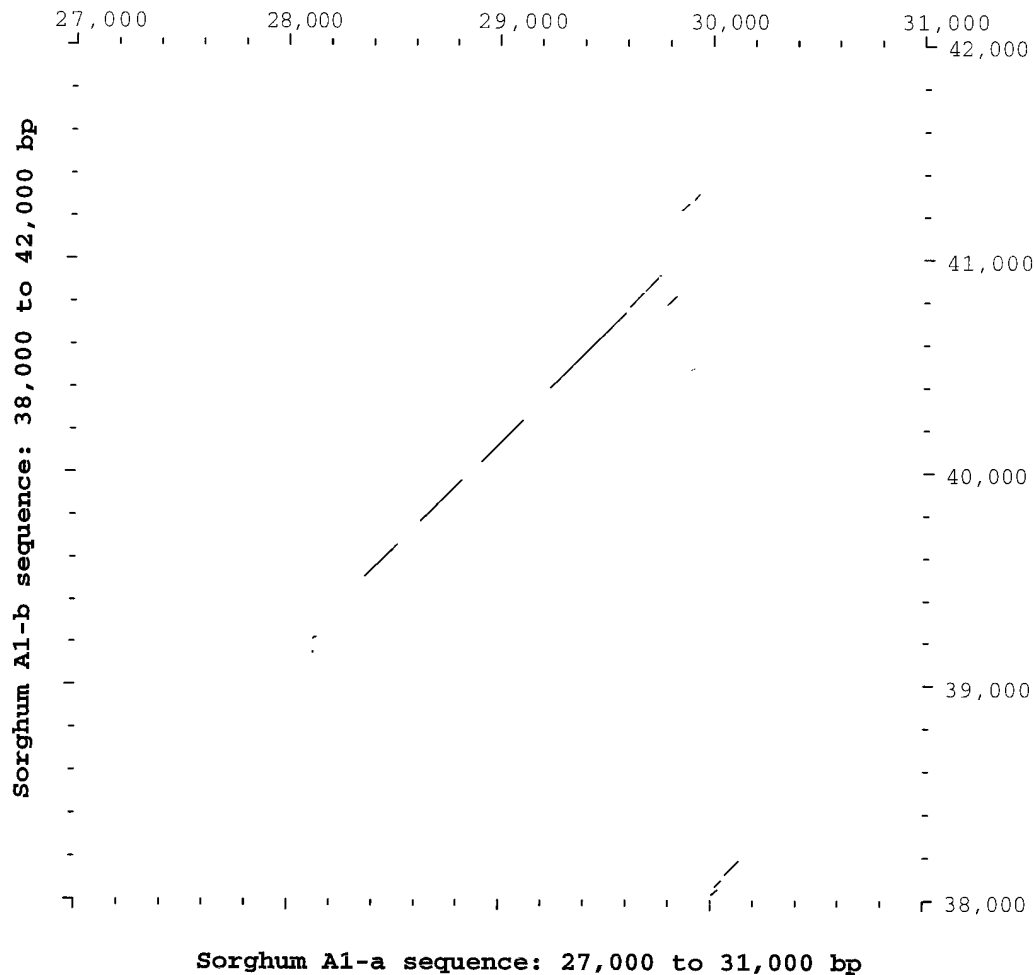


Figure 6.—DNA sequence comparison of the tandem duplication of *a1* homologue in sorghum.

different genome sizes. They have retained the same gene content and gene order (Whitkus *et al.* 1992), although maize has a 3.5-fold larger genome (Arumuganathan and Earle 1991). In the *Adh1* region of maize, intergenic regions are composed of a small set of nested retroelements that make up over 50% of the maize genome (SanMiguel *et al.* 1996). In the 42-kb region that we have sequenced in sorghum, only a 4.5-kb solo LTR of the retroelement *Leviathan* was identified, accounting for less than 15% of this region. In maize, the physical distance between *sh2* and *a1* is 140 kb (Civardi *et al.* 1994), whereas in sorghum it is only 22 kb. Consequently, nested retrotransposons between the genes in maize probably account for the majority of the over sixfold size difference in this region.

Comparison of the sorghum and maize sequences reveals that the introns in the *sh2* homologue are more conserved than the introns of the *a1* homologues. The noncoding regions of these genes apparently have different rates of divergence during evolution, despite their tight linkage. We know that the *sh2* homologue in rice is transcriptionally active because exons of the *sh2* genomic sequence were 100% identical to a cDNA sequence from rice (Chen and Bennetzen 1996). It is

not clear whether the *a1* homologues in sorghum and rice are expressed. However, the high homology shared with the maize *a1* homologue, especially within the exons, suggests that they are active. The reason for the different rates of divergence for introns of closely linked loci and their possible implications in evolution remain unknown.

Gene X was first identified in the rice *sh2/a1*-homologous region (Chen and Bennetzen 1996). It is transcriptionally active in rice and encodes a putative transcription factor. By comparison to the rice sequence, we have identified a gene X homologue in this region in sorghum. Surprisingly, the putative zinc finger motif is missing in the sorghum sequences. This putative zinc finger motif belongs to the ring finger family and might participate in protein-protein interactions (Berg and Shi 1996). The implications of the absence of the zinc finger motif on the function of sorghum gene X remain to be investigated, but the absence of this motif suggests that the sorghum gene may have fewer protein interactions. We have identified an expressed gene X relative in maize by searching a cDNA database (B. Bowen, personal communication), which revealed a maize endosperm-specific homologue. Therefore, gene

X homologues exist in the sorghum, rice, and maize genomes. They are located between the *sh2* and *a1* homologues in rice and sorghum and may also be in this orthologous position in maize.

We have characterized a direct tandem duplication of *a1* homologues in sorghum. The exons of these genes are highly homologous, but the 5' and 3' non-coding regions are not conserved. This tandem duplication presumably arose by unequal recombination, but the duplicated region has apparently undergone extensive subsequent divergence. The similarity in sequence between A1-a and A1-b is not much greater than their homology to the *a1* locus of maize, suggesting that these duplicated genes have diverged for about 15–20 million years. The timing of the duplication event could have been more ancient, however, because subsequent unequal conversions might have permitted a reasonable duration of concerted evolution (Elder and Turner 1995). Laughnan (1952, 1955) demonstrated that some *a1* alleles in maize also exist as directly duplicated tandems. We have recently observed that the *a1* tandem duplication is present in the *sh2/a1*-homologous region in rice, where the two *a1* homologues are separated by about 5 kb (M. Chen and J. L. Bennetzen, unpublished data). These results suggest that the *a1* tandem duplication may have occurred before the divergence of maize, rice, and sorghum from a common ancestor. Alternatively, independent tandem duplication of *a1* orthologues must be a relatively frequent event.

We thank J. Wendel for helpful comments on this manuscript and S. Frank for technical assistance. This research was supported by United States Department of Agriculture grant #94-37300-0299.

LITERATURE CITED

- Ahn, S., J. A. Anderson, M. E. Sorrells and S. D. Tanksley, 1993 Homoeologous relationships of rice, wheat and maize chromosomes. *Mol. Gen. Genet.* **241**: 483–490.
- Arumuganathan, K., and E. D. Earle, 1991 Nuclear DNA content of some important plant species. *Plant Mol. Biol. Rep.* **9**: 208–218.
- Bennetzen, J. L., 1996 The contributions of retroelements to plant genome organization, function and evolution. *Trends Micro.* **4**: 347–353.
- Bennetzen, J. L., and M. Freeling, 1993 Grasses as a single genetic system: genome composition, colinearity and compatibility. *Trends Genet.* **9**: 259–261.
- Berg, J. M., and Y. Shi, 1996 The galvanization of biology: a growing appreciation for the roles of zinc. *Science* **271**: 1081–1085.
- Chen, M., and J. L. Bennetzen, 1996 Sequence composition and organization in the *sh2/a1*-homologous region of rice. *Plant Mol. Biol.* **32**: 999–1001.
- Chen, M., P. SanMiguel, A. C. De Oliveira, S.-S. Woo, H. Zhang *et al.*, 1997 Microcolinearity in the *sh2*-homologous regions of the maize, rice and sorghum genomes. *Proc. Natl. Acad. Sci. USA* **94**: 3431–3435.
- Civardi, L., Y. Xia, K. J. Edwards, P. S. Schnable and B. J. Nikolau, 1994 The relationship between genetic and physical distances in the cloned *a1-sh2* interval of the *Zea mays* L. genome. *Proc. Natl. Acad. Sci. USA* **91**: 8268–8272.
- Cook, W. B., and J. C. Walker, 1992 Identification of a maize nucleic acid-binding protein (NBP) belonging to a family of nuclear-encoded chloroplast proteins. *Nucleic Acids Res.* **20**: 359–364.
- Das, O. P., E. Poliak, K. Ward and J. Messing, 1991 A new allele of the duplicated 27KD zein locus of maize generated by homologous recombination. *Nucleic Acids Res.* **19**: 3325–3330.
- Devereaux, J., P. Haeblerli and O. Smithies, 1984 A comprehensive set of sequence analysis programs for the VAX. *Nucleic Acids Res.* **12**: 387–395.
- Elder, J. F., and B. J. Turner, 1995 Concerted evolution of repetitive DNA sequences in eukaryotes. *Quart. Rev. Biol.* **70**: 297–320.
- Gaut, B. S., and J. F. Doebley, 1997 DNA sequence evidence for the segmental allotetraploid origin of maize. *Proc. Natl. Acad. Sci. USA* **94**: 6809–6814.
- Hannah, L. C., and J. R. Shaw, 1992 Genomic nucleotide sequence of a wild-type *shrunken-2* allele of *Zea mays*. *Plant Physiol.* **98**: 1214–1216.
- Hulbert, S. H., T. E. Richter, J. D. Axtell and J. L. Bennetzen, 1990 Genetic mapping and characterization of sorghum and related crops by means of maize DNA probes. *Proc. Natl. Acad. Sci. USA* **87**: 4251–4255.
- Laughnan, J. R., 1952 The action of allelic forms of the gene *A* in maize: IV. On the compound nature of *A^b* and the occurrence and action of its *A^d* derivatives. *Genetics* **37**: 375–395.
- Laughnan, J. R., 1955 Structural and functional aspects of the *A^b* complexes in maize: I. Evidence for structural and functional variability among complexes of different geographic origin. *Proc. Natl. Acad. Sci. USA* **41**: 78–84.
- Moore, G., K. M. Devos, Z. Wang and M. D. Gale, 1995 Grasses, line up and form a circle. *Curr. Biol.* **5**: 737–739.
- Pelase-Siebenbourg, F., C. Caelles, J. C. Kader, M. Delseny and P. Puigdomenech, 1994 A pair of genes coding for lipid-transfer proteins in *Sorghum vulgare*. *Gene* **148**: 305–308.
- Sanmiguel, P., A. Tikhonov, Y.-K. Jin, N. Motchoulskaia, D. Zakharov *et al.*, 1996 Nested retrotransposons in the intergenic regions of the maize genome. *Science* **274**: 765–768.
- Schwarz-Sommer, Z., N. Shepherd, E. Tacke, A. Gierl, W. Rohde *et al.*, 1987 Influence of transposable elements on the structure and function of the *A1* gene of *Zea mays*. *EMBO J.* **6**: 287–294.
- Strathmann, M., B. A. Hamilton, C. A. Mayeda, M. I. Simon, E. M. Meyerowitz *et al.*, 1991 Transposon-facilitated DNA sequencing. *Proc. Natl. Acad. Sci. USA* **88**: 1247–1250.
- Whitkus, R., J. Doebley and M. Lee, 1992 Comparative genome mapping of sorghum and maize. *Genetics* **132**: 1119–1130.

Communicating editor: J. A. Birchler