

# Supporting Information

## SI Material and Methods

### *Nucleic acid sequence alignments.*

Three new parameters were recently defined in Salse *et al.* (1) to increase the stringency and significance of BLAST sequence alignment by parsing BLASTN results and rebuilding HSPs (High Scoring Pairs) or pairwise sequence alignments. The first parameter, AL (Aligned length), corresponds to the sum of all HSP lengths. The second, CIP (cumulative identity percentage) corresponds to the cumulative percent of sequence identity obtained for all the HSPs ( $CIP = \sum nb \text{ ID by HSP} / AL \times 100$ ). The third parameter, CALP is the cumulative alignment length percentage. It represents the sum of the HSP lengths (AL) for all the HSPs divided by the length of the query sequence ( $CALP = AL / \text{Query length}$ ). The CIP and CALP parameters allow the identification of the best alignment, *i.e.* the highest cumulative percentage of identity in the longest cumulative length, taking into account all HSPs obtained for any pairwise alignment. These parameters were applied to all the BLAST alignments that were performed in the present study.

### *Wheat, barley, rice, maize sorghum sequence databases.*

The 6,426 **wheat** ESTs representing 15,569 non redundant loci that were assigned to deletion bins by Qi *et al.* (2) were downloaded from the GrainGene web site (<http://wheat.pw.usda.gov/>). Since some ESTs may correspond to non overlapping 3' and 5' ends of the same cDNA sequence, a unigene set of 5,003 mapped wheat ESTs was established (1) and used in the current analysis. The **barley** data correspond to 1,015 mapped ESTs (3) the 2,755 mapped ESTs downloaded from the 'HarvEST' web site [<http://www.harvest-web.org>]. A unigene set of 3,423 non-redundant barley mapped ESTs was constructed from these data set following the same method used for constructing the Wheat unigene set (with 95% CIP, 85% CALP as thresholds criteria). The sequences of the 12 **rice** pseudomolecules (build 5) were downloaded from the 'tigr' web site ([www.tigr.org](http://www.tigr.org)) and the 41,046 genes identified by the annotation were used for the analysis. The sequences of the 10 **sorghum** chromosomes were downloaded from the 'phytozome' web site ([www.phytozome.net](http://www.phytozome.net)) and the 34,008 gene models identified by the annotation were used for the analysis. The sequences of the 10 **maize** chromosomes were downloaded from the 'genome arizona' web site ([www.genome.arizona.edu/](http://www.genome.arizona.edu/)) and the 60,284 gene models (STS, Sequence Tag Sites) identified by the annotation were used for the analysis as described in Wei *et al.* (4).

### *Identification of duplicated and syntenic regions.*

Duplications were identified through intra-genomic sequence comparisons. Paralogous pairs were identified with 70% CIP and 70% CALP cut off values after the alignment of the annotated genes available for each genomes, (1). Syntenic regions were identified through inter-genomic sequence comparisons. In this case, orthologous pairs were identified using 60% CIP and 70% CALP cut off values after the alignment of the annotated genes available for each genome. (1). Paralogous pairs identified during the duplication analyses as well as orthologous gene pairs identified from the colinearity analysis were statistically validated with the CloseUp software (5) as described in Salse *et al.* (1). Two criteria were established to define paralogous or orthologous regions taking into account the physical size (Size), number of annotated genes (Gnumber) and number of orthologous or paralogous couples (Cnumber) in the regions. The two criteria are the Density Ratio (DR) and the Cluster Ratio (CR) (3):  $DR = [(Size_1 + Size_2) / (2 \times Cnumber)] \times 100$  and  $CR = [(2 \times Cnumber) / (Gnumber_1 + Gnumber_2)] \times 100$ . The Density Ratio (DR) considers the number of links between two regions (duplicated or syntenic) as a function of the size of the considered blocks. The Cluster Ratio (CR) considers the number of links between two regions (duplicated or syntenic) as a function of the number of annotated genes available in the considered blocks. Statistically significant collinear or duplicated regions are associated with the highest Density Ratio and Cluster Length values. The remaining collinear or duplicated regions were considered as artificial, *i.e.* obtained at random considering the number of links between two regions (duplicated or syntenic) characterized by a physical size and number of annotated genes available (1).

### ***Graphical display***

Colinearity was visualized graphically using the Genome Pixelizer software ([http://niblrrs.ucdavis.edu/GenomePixelizer/GenomePixelizer\\_Welcome.html](http://niblrrs.ucdavis.edu/GenomePixelizer/GenomePixelizer_Welcome.html)). Duplications were visualized using the Circos software (<http://mkweb.bcgsc.ca/circos/>). Three data sheets comprising (1) marker information (name, position on chromosome, chromosome number), (2) link information (data obtained for the statistical identification of paralogous or orthologous sequence pairs) and (3) graphical information (number, size of chromosomes) were provided to the two softwares.

### ***CIP/CALP distribution for paralogous and orthologous gene pairs.***

We classically performed the sequence divergence as well as speciation event analysis based on the rate of nonsynonymous ( $Ka$ ) vs. synonymous ( $Ks$ ) substitutions calculated with MEGA-3 (6). The average substitution rate ( $r$ ) of  $6.5 \times 10^{-9}$  substitutions per synonymous site per year for grasses is classically used to calibrate the ages of the considered gene (7, 8). The time ( $T$ ) since gene insertion is then classically estimated using the formula  $T = Ks / r$  (8). We

have recently performed this strategy to study the sequence divergence between paralogs and orthologs for cereals species and subspecies (1, 9, 10). However, as it has been already reported in cereal (1) and vertebrate genome (11) analyses, this strategy is not the most appropriate when looking and comparing sequence divergence for old speciation events such as the 7 shared cereal duplications dating back to more than 50-70 My in the current study. As a consequence, here, sequence similarity has been accurately investigated based on sequence alignment identity [see Jaillon *et al.* (11) for teleost fish *Tetraodon nigroviridis* duplication analysis]. In the current analysis we then investigated the sequence divergence among the paralogous and orthologous gene pairs based on the average CIP/CALP value ranking from 70 to 100%. For example, an average CIP/CALP value of 100% corresponds to two orthologous or paralogous sequences that have 100% identity over 100% of the entire sequence length.

#### ***Ancestor genome reconstruction.***

The dot plot (12) cluster approach was used to support ancestral genome reconstruction by taking into account the syntenic and duplication relationships that were characterized independently between the 5 cereal genomes. First, a series of dot plots was created for each of the orthologous cereal chromosome segments using rice as a template (12 dot plots). Then, the dot plots were assembled into 5 non-redundant groups of duplicates based on information from the 7 shared duplications. Dot plots were represented using the SpotFire software, TIBCO (<http://spotfire.tibco.com>). The analysis resulted in the definition of 22 blocks within the 5 ancestral proto-chromosomes delineated by the paleoduplication boundaries. Within each chromosome blocks, the gene number was estimated as the non-redundant number of genes (conserved between at least two genomes) present on the dot plot diagonal. The size of the ancestral genome was then estimated on the basis of the longest shared sequence conserved between at least two genomes.

#### **SI Material and Methods references**

1- Salse J, *et al.* (2008) Identification and characterization of shared duplications between rice and wheat provide new insight into grass genome evolution. *Plant Cell* 20(1):11-24.

2- Qi LL, *et al.* (2004). A chromosome bin map of 16,000 expressed sequence tag loci and distribution of genes among the three genomes of polyploid wheat. *Genetics* 168:701-712.

3- Stein N, *et al.* (2007) A 1,000-loci transcript map of the barley genome: new anchoring points for integrative grass genomics. *Theor Appl Genet* 114:823-839.

- 4- Wei F, *et al.* (2007) Physical and genetic structure of the maize genome reflects its complex evolutionary history. *PLoS Genet* 3(7):e123.
- 5- Hampson SE, Gaut BS, Baldi P (2005). Statistical detection of chromosomal homology using shared-gene density alone. *Bioinformatics* 21:1339-48.
- 6- Kumar S, Tamura K, Nei M (2004) MEGA3: Integrated software for Molecular Evolutionary Genetics Analysis and sequence alignment. *Brief Bioinform* 5:150-163.
- 7- Gaut BS, Morton BR, McCaig BC, Clegg MT (1996) Substitution rate comparisons between grasses and palms: synonymous rate differences at the nuclear gene *Adh* parallel rate differences at the plastid gene *rbcL*. *Proc Natl Acad Sci USA* 93(19):10274-9.
- 8- SanMiguel P, Gaut BS, Tikhonov A, Nakajima Y, Bennetzen JL (1998) The paleontology of intergene retrotransposons of maize. *Nat Genet* 20(1):43-5.
- 9- Chantret N, *et al.* (2005) Molecular basis of evolutionary events that shaped the hardness locus in diploid and polyploid wheat species (*Triticum* and *Aegilops*). *Plant Cell* 17(4):1033-45.
- 10- Salse J, *et al.* (2008) New insights into the origin of the B genome of hexaploid wheat: Evolutionary relationships at the SPA genomic region with the S genome of the diploid relative *Aegilops speltoides*. *BMC Genomics* 9(1):555.
- 11- Jaillon O, *et al.* (2004) Genome duplication in the teleost fish *Tetraodon nigroviridis* reveals the early vertebrate proto-karyotype. *Nature* 431(7011):946-957.
- 12- Sonnhammer EL, Durbin R (1995) A dot-matrix program with dynamic threshold control suited for genomic DNA and protein sequence analysis. *Gene* 167:1–10.

**Table S1: Orthologous relationships and shared duplications between the rice, wheat, barley, sorghum and maize genomes.**

Colinear rice, wheat, barley, sorghum and maize chromosomes are displayed on the same lines in the table. The number of orthologs identified between rice and the 4 other genomes is indicated in brackets beside each chromosome. Intra-genomic Duplications (ID) are indicated with black brackets while duplications identified with the Dual Synteny (DS) approach are indicated with red brackets. Dotted black brackets illustrated the chromosome fusions identified in maize and sorghum. Dotted and black squares represent lineage specific events in the Triticeae (translocation between chromosomes 4-5) and maize (tetraploidisation event) genomes, respectively. The 5 ancestral chromosomes (A5, A7, A11, A8, A4) defined by the shared duplications identified in the four genomes are displayed on the left side with the same color code as used in the figures in the text. The 2 duplications shared between the Triticeae and rice chromosomes that do not share a common ancestry but are found on orthologous chromosomes in both species are indicated with black double arrows. The numbers of orthologs, paralogs shared and lineage specific duplications are indicated at the end of the table for the 5 genomes.

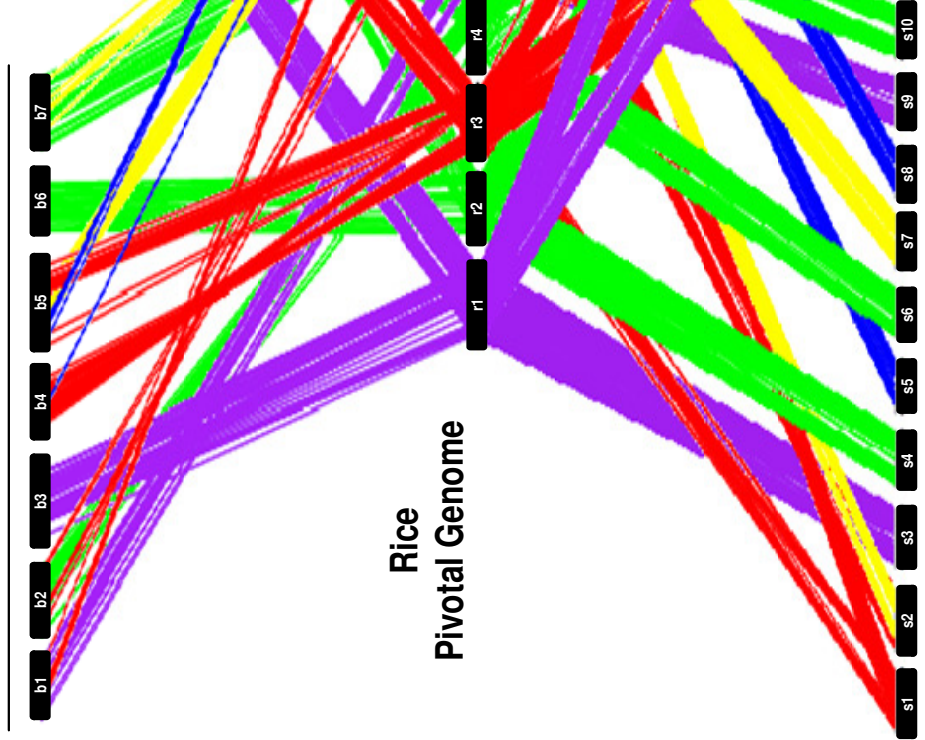
Unigene Species	41046 Rice	5003 Wheat	3423 Barley	34008 Sorghum	60284 Maize
<b>A5</b>	r5(3298)	w1 (67)	b1(23) b3(4)	s9(446) s1(16)	m6(187) m8(133) m10(20)
	r1(5313)	w3 (176)	b3(47) b1(3)	s3(1024) s9(11)	m3(417) m8(296) m1(19)
<b>A7</b>	r3(4559)	w4 (112) w5 (26)	b4(43) b5(20) b7(5) b3(3)	s1(1049) s6(9) s2(9)	m1(515) m9(158) m5(119) m7(13) m2(10)
	r7(3270)	w2 (39)	b2(14) b7(3)	s2(448) s6(11)	m7(209) m2(113) m5(7)
	r10(2404)	w1 (17)	b1(10) b5(4) b4(3)	s1(263) s3(8)	m1(125) m5(48) m9(22) m2(5)
<b>A11</b>	r11(2936)	w4 (15)	b4(2) b1(2)	s5(145) s8(18)	m4(68) m2(32) m1(7)
	r12(2610)	w5 (30)	b5(11) b4(1)	s8(180) s3(12)	m3(49) m10(45) m1(43) m5(10)
<b>A8</b>	r8(2905)	w7 (45)	b7(11) b6(6)	s7(379) s1(16)	m1(99) m4(98) m10(41) m6(24) m2(9)
	r9(2399)	w5 (11)	b5(22) b7(1)	s2(340) s1(8)	m7(160) m2(87) m5(7)
<b>A4</b>	r2(4319)	w6 (108)	b6(40) b2(3)	s4(821) s1(16)	m5(346) m4(267) m9(9)
	r4(3613)	w2 (89)	b2(40) b4(1)	s6(595) s10(9)	m2(274) m10(170) m5(10)
	r6(3420)	w7 (92)	b7(26) b5(2)	s10(457) s4(15)	m9(141) m6(121) m5(47) m4(18)
<b>Orthologs</b>		<b>827</b>	<b>309</b>	<b>6147</b>	<b>4454</b>
<b>Paralogs</b>	<b>383</b>	<b>102</b>	<b>38</b>	<b>390</b>	<b>3469</b>
<b>Shared Duplications</b>	<b>7</b>	<b>7</b>	<b>5</b>	<b>5</b>	<b>5</b>
<b>Lineage-specific duplications</b>	<b>3-12; 4-8; 4-10</b>	<b>3-5; 1-7; 3-7</b>	<b>1-5; 2-5; 3-7; 5-6</b>	<b>1-3; 1-10; 3-4</b>	<b>15 Tetraploidy + 1-2; 3-4</b>
<b>Total duplications</b>	<b>10</b>	<b>10</b>	<b>9</b>	<b>8</b>	<b>17</b>

**Figure S1 : Enlarged resolution of Figure 1A.**

Schematic representation of the 11,737 orthologs identified between the rice chromosomes (r1 to r12) used as a reference, and the barley (b1 to b7), wheat (w1 to w7), sorghum (s1 to s10), and maize (m1 to m10) chromosomes. Each line represents an orthologous gene. The 5 different colors used to represent the blocks reflect the origin from the 5 ancestral proto-chromosomes (1).

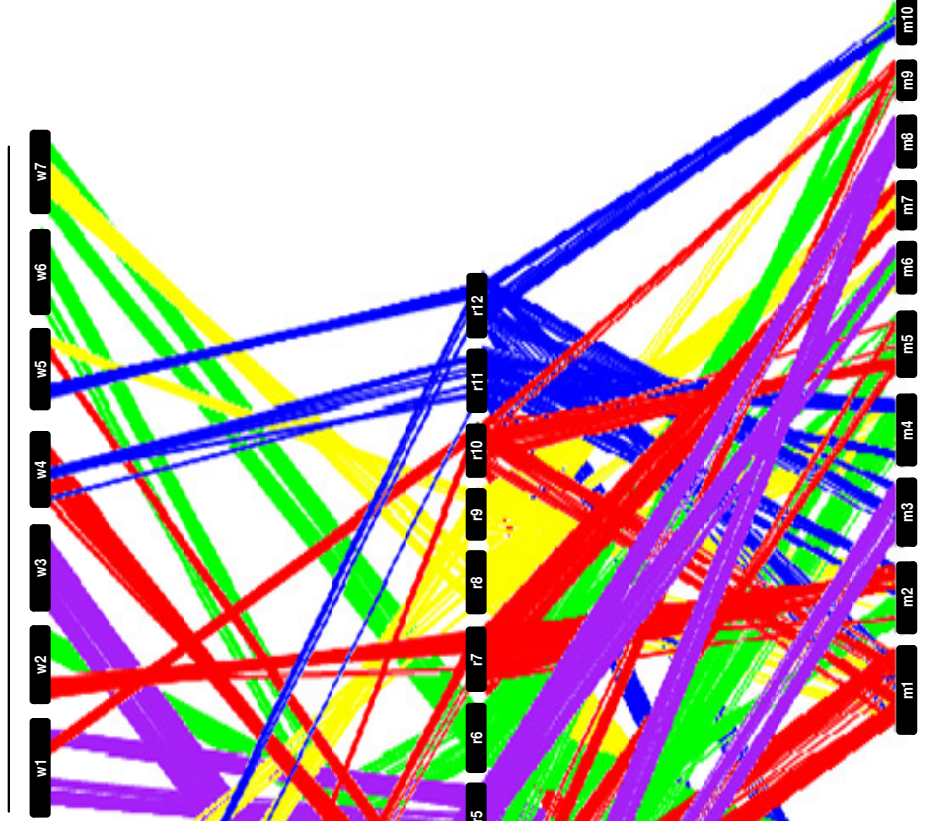
### Barley

[13 blocks – 309 orthologs]



### Wheat

[13 blocks – 827 orthologs]



### Rice Pivotal Genome

### Sorghum

[12 blocks – 6 147 orthologs]



### Maize

[30 blocks – 4 454 orthologs]

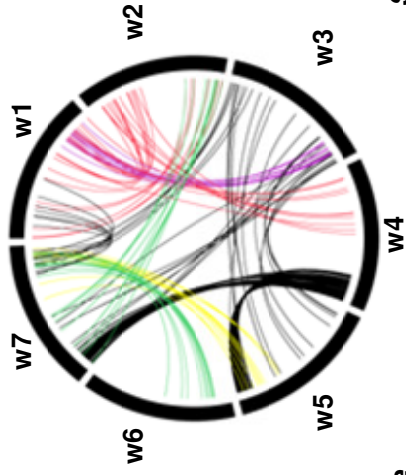




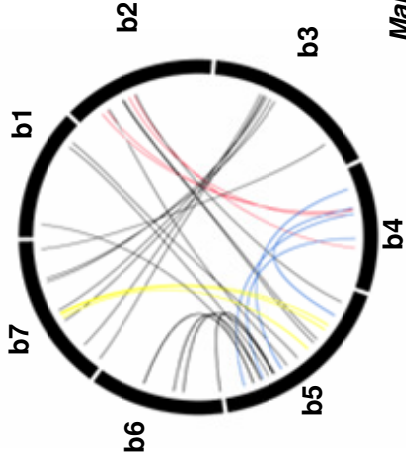
**Figure S2 : Enlarged resolution of the article Figure 1B.**

Schematic representation of the 4,382 paralogous pairs identified within the rice (r1 to r12), barley (b1 to b7), wheat (w1 to w7), sorghum (s1 to s10), and maize (m1 to m10) genomes. Each line represents a duplicated gene. The 5 different colors used to represent the blocks reflect the origin from the 5 ancestral proto-chromosomes.

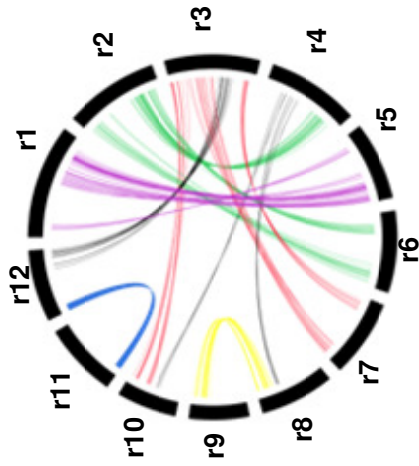
**Wheat**  
[10 blocks – 102 paralogs]



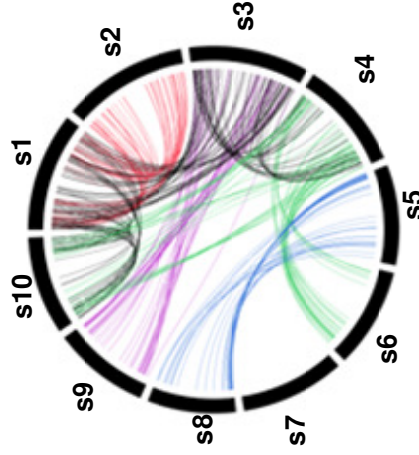
**Barley**  
[9 blocks – 38 paralogs]



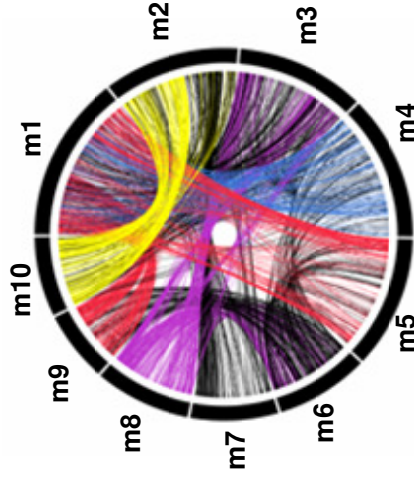
**Rice**  
[10 blocks – 383 paralogs]



**Sorghum**  
[8 blocks – 390 paralogs]

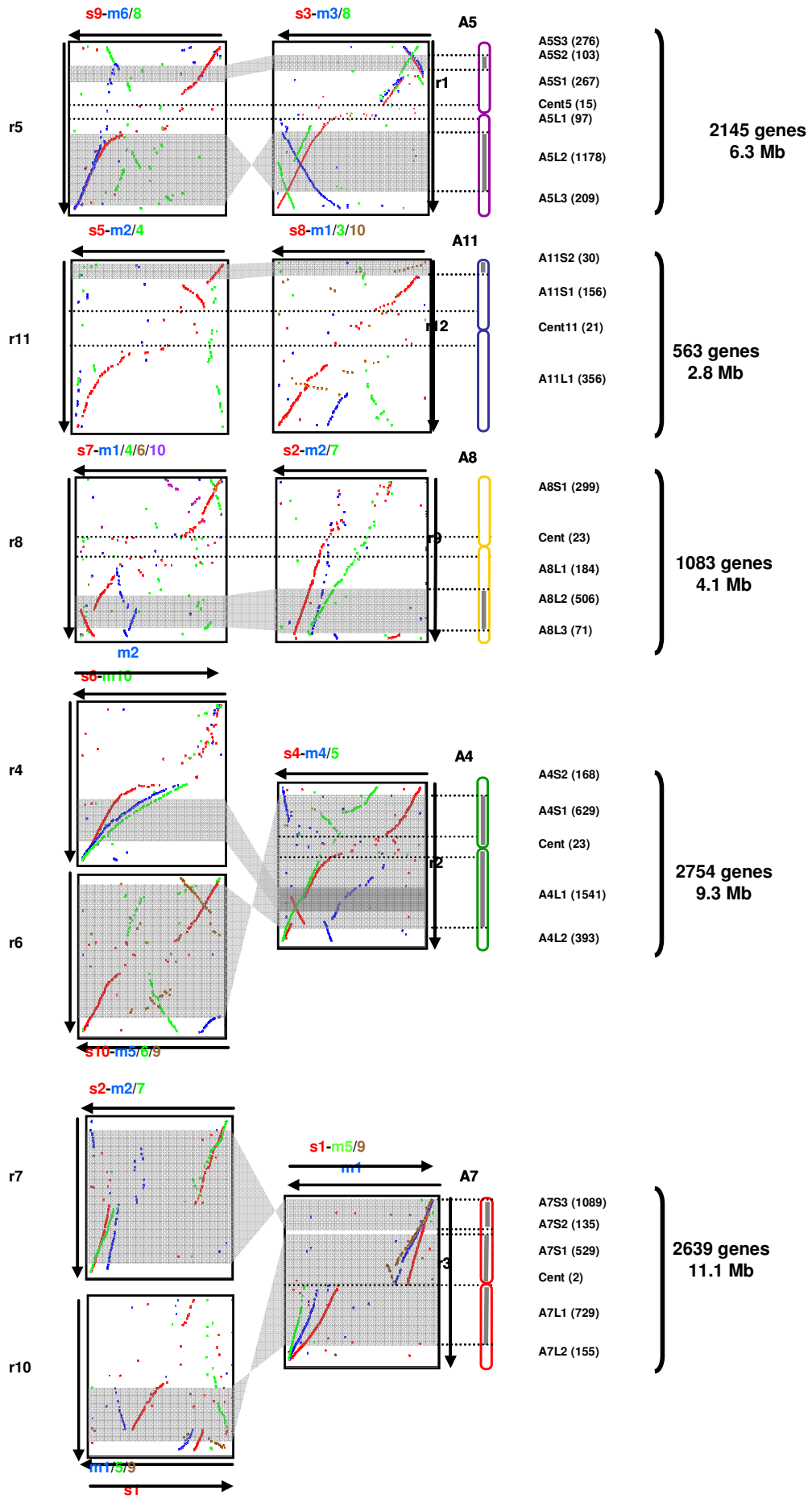


**Maize**  
[17 blocks – 3 469 paralogs]



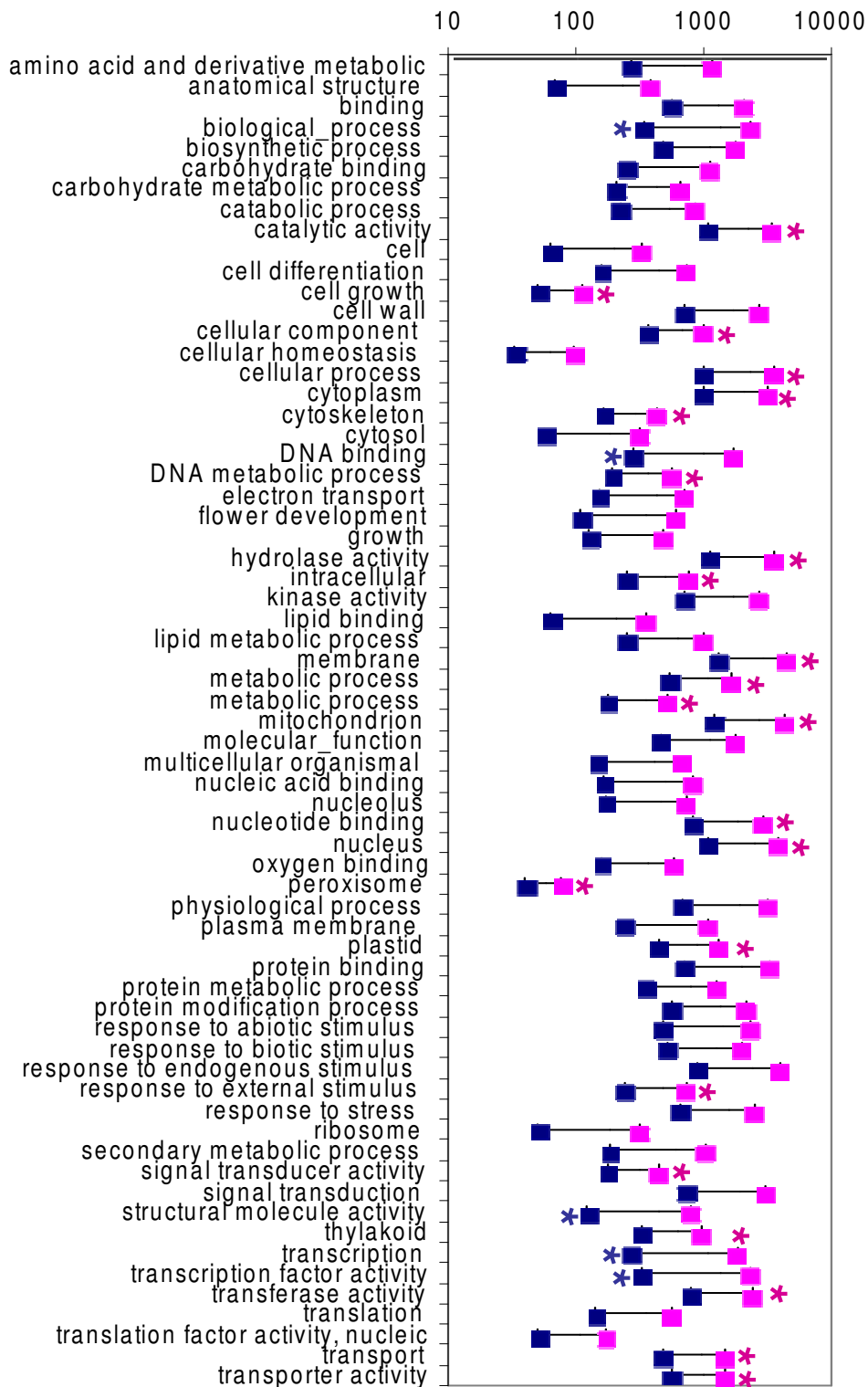
**Figure S3 : Enlarged resolution of the article Figure 2A.**

The synteny between rice (r), considered as the reference sequence (vertical), and maize (m) or sorghum (s) (horizontal) is shown as 12 dot-plots. The rice/sorghum synteny is depicted with 12 red dot-plots. The synteny between rice and maize is displayed as 12 blue and green dot-plots (reflecting the tetraploid nature of the maize genome). The seven paleo-duplications are indicated by grey blocks within the dot plots. 22 ancestral proto-chromosome blocks (from A5S3 to A7L2) were identified with respect to paleo-duplication boundaries shown with grey blocks on A5, A11, A8, A4, A7 harboring different colored blocks reflecting the origin from the 5 ancestral proto-chromosomes. The number of conserved genes (cumulative diagonal dot-plots) as well as the physical size (cumulative coding sequence length) of each proto-chromosome block is shown in parenthesis on the right end side of the figure.



**Figure S4 : Evolution of the cereal proto-chromosome gene functions.**

Graphical representation of the number of genes (in log scale) identified for the 65 Gene Ontology (GO provided on the left end side) classes established for the 41,046 rice genes (<http://gnn.tigr.org/tdb/e2k1/osa1/GO.retrieval.shtml>). The pink distribution represents the GO distribution for the rice genome and the blue distribution the number of genes for the corresponding GO classes observed for the ancestral genome gene content, *i.e.* genes that are conserved between at least two of the cereal genomes analysed in this work. Blue and pink stars indicate the GO classes that are statistically underrepresented in the ancestor and in the actual cereal genomes, respectively.



■ Ancestor

■ Rice