

Structural and functional partitioning of bread wheat chromosome 3B

Frédéric Choulet,^{1,2*} Adriana Alberti,³ Sébastien Theil,^{1,2} Natasha Glover,^{1,2} Valérie Barbe,³ Josquin Daron,^{1,2} Lise Pingault,^{1,2} Pierre Sourdille,^{1,2} Arnaud Couloux,³ Etienne Paux,^{1,2} Philippe Leroy,^{1,2} Sophie Mangenot,³ Nicolas Guilhot,^{1,2} Jacques Le Gouis,^{1,2} Francois Balfourier,^{1,2} Michael Alaux,⁴ Véronique Jamilloux,⁴ Julie Poulain,³ Céline Durand,³ Arnaud Bellec,⁵ Christine Gaspin,⁶ Jan Safar,⁷ Jaroslav Dolezel,⁷ Jane Rogers,⁸ Klaas Vandepoele,⁹ Jean-Marc Aury,³ Klaus Mayer,¹⁰ Hélène Berges,⁵ Hadi Quesneville,⁴ Patrick Wincker,^{3,11,12} Catherine Feuillet^{1,2}

We produced a reference sequence of the 1-gigabase chromosome 3B of hexaploid bread wheat. By sequencing 8452 bacterial artificial chromosomes in pools, we assembled a sequence of 774 megabases carrying 5326 protein-coding genes, 1938 pseudogenes, and 85% of transposable elements. The distribution of structural and functional features along the chromosome revealed partitioning correlated with meiotic recombination. Comparative analyses indicated high wheat-specific inter- and intrachromosomal gene duplication activities that are potential sources of variability for adaption. In addition to providing a better understanding of the organization, function, and evolution of a large and polyploid genome, the availability of a high-quality sequence anchored to genetic maps will accelerate the identification of genes underlying important agronomic traits.

Bread wheat (*Triticum aestivum* L.) is a staple food for 30% of the world population. It is a hexaploid species ($6x = 2n = 42$, AABBDD) that originates from two interspecific hybridizations estimated to have taken place ~0.5 million and 10,000 years ago (1). The predicted closest extant representatives of the ancestral parental diploid species ($2n = 14$) are *Triticum urartu* (A genome), *Aegilops speltoides* (S genome related to the B genome), and *Aegilops tauschii* (D genome). Each of the three ancestral genomes is about 5.5 Gb in size and, therefore,

results in a highly redundant 17-Gb hexaploid genome with three homeologous sets of seven chromosomes (1A to 7A, 1B to 7B, and 1D to 7D), each carrying highly similar gene copies. Moreover, most of the genome was shaped by the amplification of transposable elements (TEs) that include highly repeated families and sequences (2). This high redundancy has complicated the assembly of a complete and properly ordered reference sequence of the bread wheat genome. A fully sequenced genome enables scientists and breeders to have access to a complete gene set, with the gene order along each chromosome, and to identify candidate genes between markers associated with important traits. It also enables the identification of recent duplicates, which may be involved in species-specific evolution (3), and tracing of their evolutionary history. Before obtaining a full genome sequence, the wheat gene space has been investigated through various genome and transcriptome survey sequencing approaches and through microarray hybridizations (4–7). Recently, whole-genome shotgun sequencing of cultivar Chinese Spring using Roche/454 technology and syntenic-driven assembly yielded ~95,000 gene models [$N50 = 0.9$ kb (8)]. Furthermore, the gene space of the diploid wild relatives *Ae. tauschii* (DD) and *T. urartu* (AA) has also been assembled and led to describe 43,150 and 34,879 genes, respectively (9, 10). Although these sequences are useful templates for marker design and comparative analyses, as a result of assembly limitations of short-read-based sequencing (11, 12), they are still very fragmented, and a large fraction of the genes are unanchored to chromosomes. The maize (*Zea mays*) and potato (*Solanum tuberosum*) sequencing projects, both representing species with highly

repetitive genomes, were able to avoid overfragmentation by combining multiple sequencing technologies and through the use of DNA libraries with a diversity of insert sizes (13, 14).

The International Wheat Genome Sequencing Consortium (IWGSC) road map focuses on physically mapping and obtaining a high-quality reference sequence of each of the 21 individual wheat chromosomes rather than approaching the hexaploid genome as a whole. This strategy relies on flow-sorting individual chromosomes and/or chromosome arms from ditelosomic lines of the cultivar Chinese Spring to construct bacterial artificial chromosome (BAC) libraries (15). The largest chromosome is 3B (~1 Gb). It was the first chromosome for which a BAC library was constructed (16) and a physical map achieved (17). A pilot sequencing study on 13 contigs (2) suggested that genes tend to be mainly clustered into small islands, the presence of a twofold gene density increase from the centromere toward the telomeres, and a high proportion of nonsynthetic genes interspersed within a conserved ancestral grass gene backbone. It provided a proof of principle for this strategy and opened the way for producing a reference sequence of the large and polyploid wheat genome.

Sequencing and construction of a pseudomolecule

We used a hybrid sequencing and BAC pooling strategy to sequence 8452 BAC clones from the minimal tiling path (MTP) that was established during the construction of the chromosome 3B physical map (4, 18). After the integration of BAC-end sequences, manual curation of the scaffolding, gap filling, and correction of potential sequencing errors (18), we obtained a final assembly of 2808 scaffolds representing 833 Mb with a $N50$ of 892 kb (i.e., half of the chromosome sequence is assembled in scaffolds larger than 892 kb). We estimated that about 6% of the chromosome sequence was not present in the MTP BAC-based assembly through comparison with the 546,922 contigs assembled from whole-chromosome shotgun sequencing of flow-sorted 3B DNA (19). This suggests that the size of chromosome 3B is nearly 886 Mb—that is, about 11% smaller than originally predicted (16, 20). We built a pseudomolecule of chromosome 3B by ordering 1358 scaffolds along the chromosome using an ordered set of 2594 anchor single-nucleotide polymorphism (SNP) markers. The pseudomolecule represents 774.4 Mb (93% of the complete sequence), with a scaffold $N50$ of 949 kb (table S1). The order of markers was determined by linkage analysis of a recombinant inbred line (RIL) population derived from a cross between *T. aestivum* cultivars Chinese Spring (reference sequence) and Renan (a French elite cultivar) and refined by integrating linkage disequilibrium data from two panels and physical BAC contig information (18). This sequence corresponds to an annotation-directed improved high-quality draft (21) situated between the high-quality finished rice

¹Institut National de la Recherche Agronomique (INRA) UMR1095, Genetics, Diversity and Ecophysiology of Cereals, 5 Chemin de Beaulieu, 63039 Clermont-Ferrand, France.

²University Blaise Pascal, UMR1095, Genetics, Diversity and Ecophysiology of Cereals, 5 Chemin de Beaulieu, 63039 Clermont-Ferrand, France. ³Commissariat à l'Énergie Atomique et aux Énergies Alternatives, Direction des Sciences du Vivant, Institut de Génétique, Genoscope, 2 Rue Gaston Crémieux, 91000 Evry, France. ⁴INRA, UR1164 Unité de Recherche Génomique Info Research Unit in Genomics-Info, INRA de Versailles, Route de Saint-Cyr, 78026 Versailles, France. ⁵Centre National des Ressources Génétiques Végétales, INRA UPR 1258, 24 Chemin de Borde Rouge, 31326 Castanet-Tolosan, France. ⁶Biométrie et Intelligence Artificielle, INRA, Chemin de Borde Rouge, BP 27, 31326 Castanet-Tolosan, France. ⁷Centre of the Region Haná for Biotechnological and Agricultural Research, Institute of Experimental Botany, Slechtitelu 31, CZ-78371 Olomouc, Czech Republic. ⁸The Genome Analysis Centre, Norwich Research Park, Norwich NR4 7UH, UK. ⁹Department of Plant Systems Biology (VIB) and Department of Plant Biotechnology and Bioinformatics (Ghent University), Technologiepark 927, 9052 Gent, Belgium. ¹⁰Munich Information Center for Protein Sequences, Institute of Bioinformatics and Systems Biology, Helmholtz Zentrum Muenchen, D-85764 Neuherberg, Germany. ¹¹CNRS UMR 8030, 2 Rue Gaston Crémieux, 91000 Evry, France. ¹²Université d'Evry, CP5706 Evry, France.

*Corresponding author. E-mail: frederic.choulet@clermont.inra.fr

genome sequence and the improved draft maize genome (13).

Annotation of genes, transcribed loci, and transposable elements

Gene modeling led to the prediction of 7264 coding loci on the 3B pseudomolecule (Table 1), including 5326 with a functional structure and 1938 (27%) likely corresponding to pseudogenes. An additional 251 gene models and 188 pseudogenes were annotated in unanchored scaffolds. RNA-Seq data revealed that 71.4% of the predicted genes/pseudogenes are transcribed and led to the identification of 3692 unannotated transcribed loci that may encode functional non-coding RNAs or unknown proteins, hereafter referred to as novel transcribed regions (NTRs) (Table 1). In addition, 791 highly conserved non-coding RNA genes involved in RNA maturation and protein synthesis were also predicted (Table 1). Chromosome 3B appears to contain a high number of small nuclear RNA genes (U1 to U6) including nine U1-snRNAs (small nuclear RNAs), seven of which are tandemly duplicated. As a comparison, there are 14 U1-snRNAs in the entire *Arabidopsis thaliana* genome (www.plantgdb.org). The higher number of U1-snRNAs may reflect a higher level of duplication in the wheat genome. We found 53,288 complete and 181,058 truncated copies of TEs, belonging to 485 TE families and representing 85% (640 Mb) of the 3B pseudomolecule, through a similarity-search approach. Further de novo repeat detection (18) identified 3.6% putatively new TEs.

We estimated the putative location of the centromeric region by plotting the density of the long terminal repeat retrotransposons (LTR-RTs) CRW (centromeric retrotransposons of wheat) and Quinta along the pseudomolecule. These LTR-RTs are recognized by the centromere-specific histone CenH3 and thus are centromere-functional sequences (22). Two major peaks covering a region of 122 Mb (265 to 387 Mb) (fig. S1)—which includes 1 Mb previously shown as interacting with histone CenH3 (22) and encompassing the centromere of the orthologous rice chromosome 1 (23)—were identified. This region was defined as the centromeric-pericentromeric region of chromosome 3B. A strong correlation has been observed between the size of the centromeres and the chromosomes in grasses (24), and it is likely that large chromosomes have centromeres larger than 10 Mb. This may be critical to ensure the structural rigidity of the pericentromeric regions needed for kinetochore co-orientation (25). Marker assignment to either the short or the long arm indicated the presence of a break point between 349.4 and 350.0 Mb that might be the position of the core centromere.

Variability in recombination rate and gene density along the chromosome

We found 787 crossover (CO) events on chromosome 3B in the Chinese Spring x Renan population, with on average 2.6 COs per chromosome per individual, which is similar to maize [2.7 to

3.4 (26)]. Distribution of meiotic recombination rates revealed extreme variations along the chromosome. Whereas the average recombination rate is 0.16 cM/Mb, actual values range from 0 to 2.30 cM/Mb (per 10-Mb window) (Fig. 1A). Segmentation analysis (18) revealed partitioning with the two distal regions of 68 Mb (region R1) and 59 Mb (region R3) on the short and long arms, respectively, showing recombination rates of 0.60 cM/Mb and 0.96 cM/Mb on average, and a large proximal region of 648 Mb (region R2) spanning the centromere with an average recombination rate of 0.05 cM/Mb (Table 2 and Fig.

1A). This provides insight into the actual physical size of the highly recombinogenic regions previously detected at the end of the wheat chromosomes (27, 28). When a narrower window of 1 Mb was used, variations ranged from 0 to 12 cM/Mb (Fig. 1A), a range similar to that observed in maize [0.8 to 11.5 cM/Mb (26)] and sorghum [0 to 10 cM/Mb (29)]. All crossover events occurred in only 13% of the chromosomes in our population of 305 individuals. The largest region totally deprived of recombination corresponds to 150 Mb and includes the putative 122-Mb centromeric-pericentromeric region. This was

Table 1. General features of the 3B pseudomolecule. LINEs, long interspersed nuclear elements; SINEs, short interspersed nuclear elements; rRNA, ribosomal RNA; snoRNA, small nucleolar RNA; TIRs, terminal inverted repeats.

Pseudomolecule sequence		774,434,471		
Protein-coding genes	Length (bps)	46.16%		
	G+C content	All	Full genes	Pseudogenes
	No. of genes	7264	5326	1938
Noncoding RNA genes	Average size (bps) of coding sequences (± standard deviation)	1095 ± 807	1187 ± 821	840 ± 710
	Average number of exons (± standard deviation)	4.2 ± 4.4	4.4 ± 4.6	3.6 ± 3.8
	Gene density (kb ⁻¹)	107	145	400
	No. of expressed genes	5185	4125	1060
	No. of genes with alternative splicing	3185	2596	589
	% genes with alternative splicing	61	63	56
	Average no. isoforms/expressed gene	5.8	5.8	5.8
	NTRs	3692		
	tRNA	589		
	5S rRNA	85		
	Others (snRNA, snoRNA)	117		
	Total	791		
Transposable elements (TEs)				
Class I	Copia	15.6%		
	Gypsy	46.9%		
	Unclassified			
	LTR-retrotransposons	3.5%		
	LINEs	1.2%		
	SINEs	0.01%		
	Total class I	67.1%		
	CACTA	16.4%		
	Harbinger	0.19%		
	Mariner	0.19%		
Class II	Mutator	0.43%		
	hAT	0.02%		
	Unclassified class II with TIRs	0.22%		
	Unclassified class II	0.10%		
	Helitron	0.01%		
	Total class II	17.6%		
	Unclassified repeats	0.81%		
	Total TEs	85.5%		

Downloaded from https://www.science.org at University of East Anglia on August 12, 2022

confirmed by the linkage disequilibrium (LD) pattern (fig. S2). Twenty-two regions showed a recombination ratio higher than 1.6 cM/Mb—i.e., >10 times the average for this chromosome—and thus may contain recombination hot spots (Fig. 1A). However, no significant correlation was observed between the recombination rate and gene content, coding DNA, or TE content of these regions.

The 7264 genes are not evenly distributed, and gene density is increasing on both arms along the centromere-telomere axis, correlating with the distance to the centromere ($r_s = 0.79$, $P < 2.2 \times 10^{-16}$, $R^2 = 0.61$) (Fig. 1B). Using a 10-Mb window, the average gene density estimate is 9 ± 5 genes/Mb, ranging from 1.3 in the centromeric-pericentromeric region up to 27.9 at the most telomeric end of the short arm, a pattern commonly observed in grass genomes. Variation of the gene density in wheat chromosome 3B is higher than for chromosomes in the more compact rice genome (30); lower than in sorghum, where genes are mostly found in the telomeric regions (31); and in the same

range as in maize, which also contains a high percentage of TEs (13). Segmentation analysis revealed five major regions with contrasted gene densities (Fig. 1B) and a fourfold gradient of the gene density—i.e., twice as many as suggested by the pilot study on chromosome 3B (2). The distal segments exhibiting the highest gene density (19 genes per Mb) correspond nearly to the highly recombinogenic R1 and R3 regions (Fig. 1A). The R2 region was subdivided into three segments, with the lowest gene density (5 genes per Mb) in a 234-Mb segment encompassing the centromeric-pericentromeric region. As previously suggested (2, 4), there is no large region completely devoid of coding sequence (maximum of 3.7 Mb). We confirmed that the intergenic distances (IGDs) are extremely variable (average 104 ± 190 kb) and that a majority (73%) of the genes are organized in small islands, or insulae (32). This suggests that most of the intergenic regions are under selective constraint prevented from TE insertion. Indeed, only 29% of the IGDs are larger than 104 kb, but they

account for 81% of the chromosome size, demonstrating that TE-mediated genome expansion likely occurred within a limited number of intergenic regions.

Relationships between gene expression, function, and chromosome location

Of all annotated genes on chromosome 3B, 71.4% are expressed in at least one of the 15 conditions analyzed [five organs at three developmental stages each (table S2)], 33% in all conditions, and 5% in one only (fig. S3). On average, genes are expressed in 10.8 of 15 conditions (considering all predictions), and expressed genes are transcribed into 5.8 alternative transcripts, or isoforms. Both the expression breadth and the average number of isoforms are distributed unevenly along the chromosome, with a clear decrease of the two parameters toward the telomeres (Fig. 1, C and D). Segmentation revealed distal segments with boundaries similar to that of regions R1 and R3 and with genes expressed in fewer conditions than in the proximal region: 8.7 versus 11.7, respectively ($P < 2.2 \times 10^{-16}$, Welch *t* test) (Table 2). Similarly, the average number of alternative transcripts is higher in the proximal (6.5) than in the distal (4.3) regions ($P < 2.2 \times 10^{-16}$, Welch *t* test) (Table 2).

Gene ontology (GO) term enrichment was estimated for the R1, R2, and R3 regions (18) (tables S3 to S5). The distal regions are enriched in many GO categories, some being related to adaptation (response to abiotic stimulus or response to stress). Well-known examples of genes related to adaptation are those involved in resistance to pathogens. Chromosome 3B carries 171 genes putatively associated with disease resistance (18), and their distribution is highly biased, with 135 (79%) of them located in the distal regions (whereas these regions contain just 33% of the gene set). Such uneven distribution and the correlation with the distribution of crossovers suggest that meiotic recombination acts as a main driver for creating variability in distal regions of chromosome 3B.

To investigate whether such partitioning is a common pattern of large plant genomes, we analyzed the distribution of the gene expression breadth in maize and barley, which both exhibit large genome size (>1 Gb) and increased recombination rates at chromosomal extremities (33, 34). In barley, segmentation analysis of the seven chromosomes based on recombination data identified the same pattern as on chromosome 3B, with two highly recombinogenic distal regions and a large nonrecombinogenic region. Using expression data of eight conditions (34), we also observed that the two high-recombination distal regions carry genes expressed in fewer conditions than those carried by the low-recombination proximal regions (5.9 versus 6.7; $P = 2.2 \times 10^{-16}$, Welch *t* test) (fig. S4A). Using GO terms, we found a significant enrichment of these regions in the categories “cell death” and “defense response,” which support previous findings that barley disease resistance genes are clustered in the distal regions (34). In contrast, in maize, although we

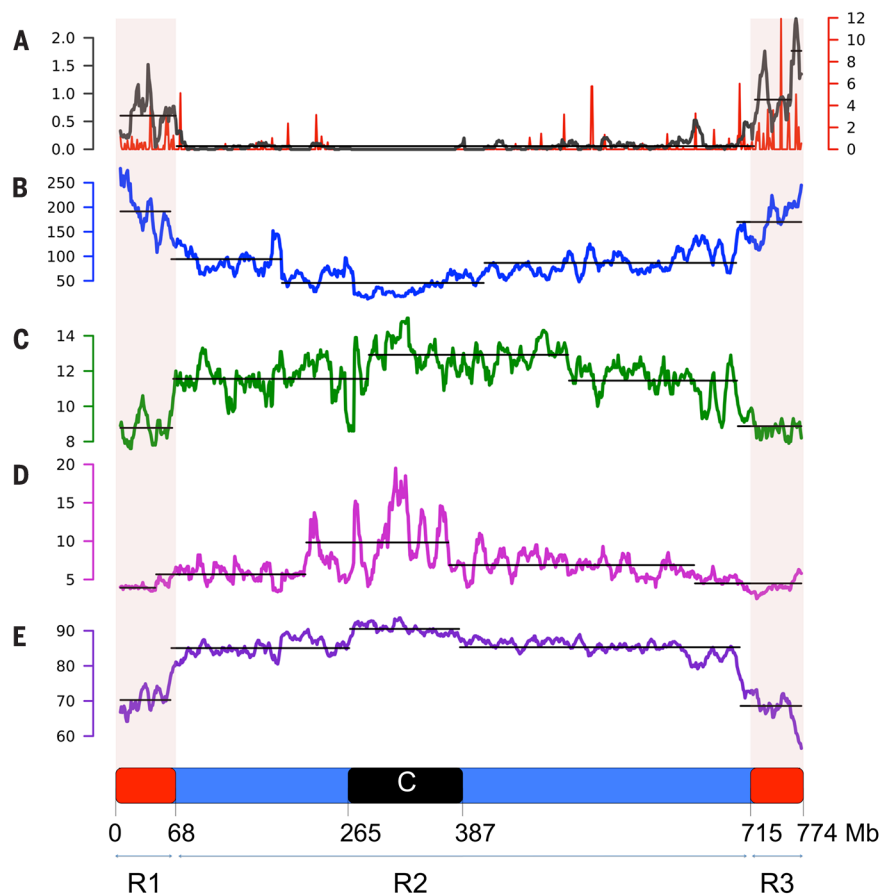


Fig. 1. Structural and functional partitioning of wheat chromosome 3B. Distribution and segmentation analysis of (A) meiotic recombination rate (cM/Mb in sliding window of 10 Mb in black and 1 Mb in red); (B) gene density (CDS/10 Mb); (C) expression breadth; (D) average number of alternative spliced transcripts per expressed gene; and (E) TE content along the 3B pseudomolecule. Distal regions of the chromosome R1 and R3 are represented in red. In (C), the centromeric/pericentromeric region is in black. The borders of these regions are indicated in Mb. Sliding window size: 10 Mb; step: 1 Mb.

observed partitioning of the recombination rate, no gene-expression partitioning was detected using RNA-Seq data of 18 conditions (35). Overall, the expression breadth is 13.2 and 12.7 in high- and low-recombination regions, respectively, with chromosome-specific patterns (fig. S4B). Nevertheless, high-recombination regions are also enriched in GO categories “cell death” and “defense response to fungus and bacteria,” suggesting that such genes are consistently found in distal recombinogenic regions in large plant genomes. These results suggest that the partitioning observed on wheat chromosome 3B is conserved in the *Triticeae* and may not reflect a general pattern of large genomes. Alternatively, it is possible that the active rearrangements observed in the maize genome have modified this pattern. Additional evidence will come once other large plant genomes (>1 Gb) are sequenced and analyzed.

Uneven distribution of transposable elements

LTR-RTs represent 66% of the chromosome 3B sequence (gypsy, 47%; copia, 16%; unclassified, 3%) (Table 1), which is slightly lower than the ~75% of LTR-RT identified in the whole maize genome (36). Only 4% (3 out of 85) of the LTR-RT families were found in single copies, compared with 41% in maize (36) and 48% in the rice genome (37). Sixteen percent of the sequence is composed of class II DNA transposons that mostly correspond to CACTA elements (Table 1), compared with 3.2% in the maize genome (13). Only six families account for 50% of the wheat chromosome 3B TE fraction, as previously suggested from partial sequence analyses (2, 38) and from observations in other large genomes (36). However, in contrast to the maize genome, in which most of the intact elements are found in 1 to 10 copies (36), the majority of the TE families annotated on chromosome 3B have a higher number of copies (10 to 1000 copies). Estimated insertion dates for the most abundant LTR-RT families showed a major peak at 1.5 million years (My) but also quite specific patterns of TE activity for each family (fig. S5). Our data support the hypothesis (2, 38) that most of the transposable elements that shaped the B genome were inserted before polyploidization [0.5 million years ago (Ma)] and have been less active since then. Distribution of recently inserted elements revealed that TE insertion occurred at a similar rate in the distal and proximal regions. In contrast, older insertions (>1.5 Ma) were 1.7 times as abundant in the R2 region compared with the R1 and R3 regions, suggesting a higher rate of TE elimination in the distal ends of chromosome 3B.

The TE density distribution was not random (Fig. 1E), with a lower density in the R1 (73%) and R3 (68%) regions compared with the R2 region (88%) (Table 2). The 122-Mb centromeric-pericentromeric region displayed the highest density (93%) of TEs. Beneath the global TE distribution pattern, each superfamily presents its own specificities (fig. S6). For example, CACTA transposons are more abundant in the distal gene-rich regions

(Table 2), supporting in situ hybridization findings at the whole-genome level (39). In addition, the distribution of TE families varied on the basis of their relative distance to genes (18) (fig. S7). DNA transposons Mutator, Harbinger, and MITEs are found close to genes, whereas LTR-RTs and CACTAs tend to be located at much larger distances from the genes. For instance, the 17,479 annotated MITEs were found to be significantly associated with genes ($r = 0.89$; $P < 1 \times 10^{-10}$), as previously observed in plant genomes (40).

Synteny between chromosome 3B and related grass genomes

Comparative genomics in grasses has been used to define syntenic relationships between different species (41, 42) and to provide insight into their evolution since the divergence from a common ancestor 50 to 70 Ma (43). We compared the wheat chromosome 3B genes (Ta3B) with the closest sequenced relative, *Brachypodium distachyon* [common ancestor, 32 to 39 Ma (44)], and with one representative of each of the *Ehrhartoideae* and *Panicoideae* grass subfamilies: *Oryza sativa indica* [rice (30)] and *Sorghum bicolor* (31), respectively. Wheat chromosomes of group 3 are syntenic with chromosome 1 of rice (Os1), chromosome 3 of sorghum (Sb3), and the distal parts of *B. distachyon* chromosome 2 (Bd2). We first

investigated potential gene loss after polyploidization by using the conserved and syntenic genes found on chromosomes Os1, Bd2, and Sb3. These represent the grass core genes that are expected to be present on wheat homeologous group 3, unless they have been lost by fractionation after polyploidization. The finding that 94% of the conserved genes are also present on the 3B sequence (Fig. 2A), which represents 94% of the chromosome (see above), suggests that no major gene loss has occurred in the B subgenome yet. This is confirmed at the whole-genome level by the results of the chromosome survey sequences (19). In contrast, 2065 genes on chromosome 3B (34.6%, including pseudogenes) shared similarity with genes on nonorthologous chromosomes in the other grass genomes. This proportion of nonsyntenic genes is much higher than the 5% (between 149 and 207) of nonsyntenic genes found in the other grass species analyzed (Fig. 2A and table S6). It confirms previous results showing substantial modifications and rearrangements of the wheat gene space (2). When looking at the conservation of the gene order, collinear genes represent 42 to 68% of the genes present on Os1, Bd2, and Sb3, whereas they represent less than 30% of the Ta3B genes (including pseudogenes) (table S7 and fig. S8). The spatial distribution of syntenic and nonsyntenic genes

Table 2. Distribution of features in the three regions of chromosome 3B as defined from the recombination segmentation along the chromosome.

	R1	R2	R3
Size (Mb)	68	648	59
Recombination rate (cM/Mb)	0.60	0.05	0.96
Genes			
Predicted gene density (Mb ⁻¹)	19	7	19
Number of predicted genes/pseudogenes	1318	4845	1,101
Full genes (no.)	910 (69%)	3682 (76%)	734 (67%)
Pseudogenes/gene fragments (no.)	408 (31%)	1163 (24%)	367 (33%)
Mean intergenic distance (kb)	49	130	52
Expressed predicted genes (no.)	823 (62%)	3629 (75%)	733 (67%)
Expressed full genes (no.)	621	2963	541
Expressed pseudogenes/fragments (no.)	202	666	192
Average expression breadth (per expressed gene; /15)	8.8	11.7	8.6
Average FPKM (per expressed gene)	141	255	156
Average number of isoforms (per expressed gene)	4.2	6.5	4.4
Proportion of nonsyntenic genes* (%)	44	28	53
Proportion of intrachromosomally duplicated genes* (%)	49	33	42
Proportion of tandemly duplicated genes* (%)	24	14	22
Proportion of dispersed duplicated genes* (%)	26	18	20
Proportion of interchromosomally duplicated genes* (%)	36	33	37
Transposable elements (%)	73.0	88.3	68.4
Copia (%)	14.7	15.8	14.1
Gypsy (%)	31.7	50.3	27.1
CACTA (%)	18.7	15.9	19.5

*Number of duplicated genes (filtered set, including pseudogenes) divided by the total number of genes in each region.

along the 3B pseudomolecule (Fig. 2B) shows an increased proportion of nonsynthetic genes in the R1 (44%) and R3 (53%) regions compared with the R2 region (28%) (Table 2). This supports the hypothesis that accelerated evolution occurred in the wheat lineage compared with other grasses (2, 45, 46), with insertions of nonsynthetic genes intercalated in the ancestral grass genome backbone through gene duplications or translocations that preferentially occurred in the distal recombining regions.

Origin and evolution of nonsynthetic genes

With such a high proportion of nonsynthetic genes, one key question is whether these genes are under selection pressure or in the process of becoming pseudogenes. On the basis of the coding sequence structure, 32% of the nonsynthetic genes (versus 17% of syntenic genes) were annotated as likely pseudogenes or gene fragments. This ratio is not surprising, given that TE activity can duplicate gene fragments that are dead upon arrival. Expression patterns revealed that a majority of the nonsynthetic genes (69% versus

82% of syntenic genes) are expressed in at least one condition tested (table S8), thereby suggesting that a large fraction of these relocated genes are unlikely to be pseudogenes and may contribute to recent wheat genome evolution and, therefore, to adaptation. Interestingly, a majority (51%) of the genes expressed in a single condition corresponds to nonsynthetic genes, whereas 80% of the genes that are expressed in all 15 conditions are syntenic genes (fig. S9). This suggests that nonsynthetic genes are involved in specific processes that may be related to adaptation, whereas syntenic genes tend to be associated with essential biological processes. This hypothesis is supported by the fact that putative resistance genes identified on chromosome 3B are mainly nonsynthetic genes (18). In addition, GO term enrichment of nonsynthetic genes revealed an overrepresentation of genes involved in response to stress (table S9).

The fact that chromosome 3B exhibits a higher number of genes than its orthologs in other grasses and that at least 94% of the ancestral grass gene backbone is conserved indicates that

most insertions of nonsynthetic genes result from interchromosomal duplication with retention of the parental copy. To test this hypothesis, we used the sequences of the 18 bread wheat chromosomes nonhomeologous to group 3 chromosomes (19) to search for potential parental copies of chromosome 3B genes elsewhere in the genome (18) (table S10). A paralog was identified for 87% of the nonsynthetic genes (18), with no bias regarding the chromosomal origin of the interchromosomally duplicated genes (fig. S10). Duplications of DNA fragments to different locations in a genome have been shown to result from double-strand break (DSB) repair (in which a copy of the foreign DNA is used as filler to repair the break) or capture by active TEs (46, 47). We analyzed the composition of the regions flanking the syntenic versus nonsynthetic genes (20 kb on each side) and found a high association of nonsynthetic genes with a class II transposon superfamily: 41% more CACTAs were found around nonsynthetic genes than around syntenic genes (fig. S11). CACTA transposons are known to capture genes (31, 48) and may have contributed substantially to interchromosomal gene duplications in wheat.

We also investigated the time since duplication of nonsynthetic genes through the analysis of nucleotide substitution rates (K_s) (18). In total, 63% of these duplications were older than 10 My and, thus, are likely shared within other *Triticeae* species, whereas 37% are potentially wheat-specific. Comparison with the barley genome survey sequence data (34) showed that at least 29% of the 3B nonsynthetic genes (versus 51% of the syntenic genes) are orthologous with barley chromosome 3H, confirming that part of the nonsynthetic genes were relocated before the divergence of wheat and barley 10 to 14 Ma.

We next asked if the high gene duplication activity is also observed at the intrachromosomal level. We identified 809 gene families with two or more copies comprising 2216 genes on chromosome 3B, which is about three times as much as in rice, *Brachypodium*, and sorghum (table S11). This indicates that, in proportion, more than twice as many genes were duplicated or retained after intrachromosomal duplications in wheat (~37%) compared with the other three grasses (~15 to 18%). About 46% of the duplicated genes of chromosome 3B are found in tandem, whereas 54% are dispersed duplicates (18) (table S12). In other grass species, a majority of the duplicated genes are organized in tandem. Given the high interchromosomal duplication activity observed in our analyses (see above), it is possible that some dispersed duplicates on chromosome 3B originated through independent interchromosomal duplications rather than through intrachromosomal duplications, thereby leading to overestimates of the latter. However, even when considering syntenic dispersed duplicates—i.e., those genes that have remained at their ancestral locus and have undergone intrachromosomal duplication—23% of the whole gene set appears to have originated from recent intrachromosomal duplications, which is still higher than in other grass species. Thus, we conclude that both inter- and intrachromosomal

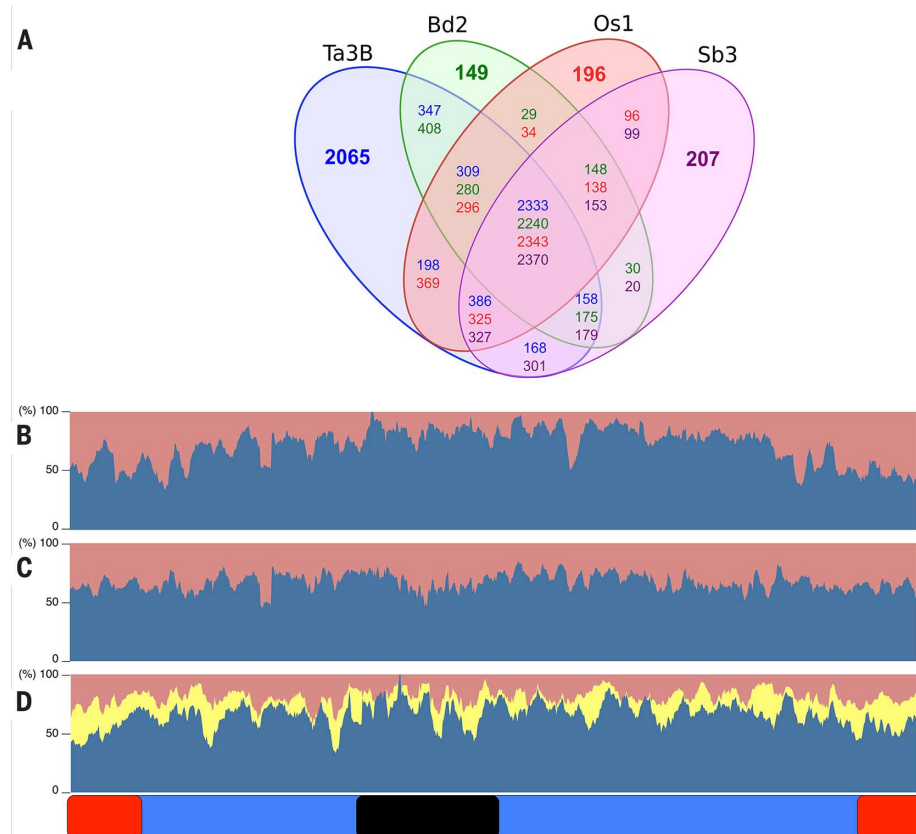


Fig. 2. Inter- and intraspecific comparative analyses of the gene content of wheat chromosome 3B. (A) Venn diagram displaying the number of genes conserved between wheat chromosome 3B (Ta3b, blue) and orthologous chromosomes in rice (Os1, red), *Brachypodium* (Bd2, green), and sorghum (Sb3, purple). The number of nonsynthetic genes is indicated in bold for each species. Distribution along the 774 Mb of the chromosome 3B pseudomolecule of the relative proportion of (B) syntenic (blue) versus nonsynthetic (red) genes, (C) interchromosomal duplications (duplicates in red, group 3-specific genes in blue), and (D) tandem (yellow) and dispersed (red) intrachromosomal duplications and singletons (blue). Chromosome 3B is represented at the bottom with distal regions in red and the centromeric/pericentromeric region in black.

rates of duplication are higher in wheat than in the other grass species analyzed so far. Interestingly, interchromosomal duplicates were distributed uniformly along the chromosome, whereas the proportion of tandem duplicates slightly increased in the distal regions (Fig. 2D). This suggests that long-distance and tandem duplications likely arose through different mechanisms. Finally, expression analysis of intra-chromosomal duplicated genes indicated that 49% of the families show expression of all copies in at least one condition. Similar to what was observed for the interchromosomal duplicated genes, the intrachromosomal duplicated genes tend to be expressed in fewer conditions as compared with nonduplicated genes (fig. S9 and table S8), suggesting that they may be undergoing subfunctionalization.

QTL mining

As exemplified in rice and other crops, a reference genome sequence provides a resource for gene discovery, marker development, and allele mining in support of crop improvement (49). We identified 153,190 insertion site-based polymorphism (ISBP) markers (50) and 35,579 micro-satellite markers along the 3B chromosome. We also located 121 quantitative trait loci (QTLs) for 50 different traits on chromosome 3B (table S13). Using these data, we conducted a meta-analysis that integrates QTLs defined in independent studies (51) and identified 18 metaQTLs with confidence intervals covering between 1.5 and 620 Mb of the chromosome 3B sequence. The largest one encompasses the centromeric region, where recombination is suppressed. Five metaQTLs with small intervals (<10 Mb) that contain between 23 and 266 protein-coding genes and between 511 and 4049 markers are suitable for fine mapping (table S14).

Discussion

We present a reference sequence of chromosome 3B that can be used to precisely delineate structural and functional features along a chromosome and establish correlations between recombination intensity, gene density, gene expression, and evolution rate. Our results indicate that during evolution, regions with distinct features become delineated along chromosome 3B, including relatively small distal regions that are preferential targets for recombination, adaptation, and genomic plasticity. Whether our observations reflect a general pattern for the wheat genome will need to be confirmed by the analysis of other chromosome reference sequences. Already, some of the features—such as the CACTA distribution, the high rate of intrachromosomal duplication, the absence of major gene loss since polyploidization, and the gradient of gene density—have been confirmed at the whole-genome level (19, 39). Moreover, the ordered chromosome 3B sequence allowed us to distinguish duplicated genes and provided evidence for superimposed mechanisms of gene duplications. The high level of gene duplication (allopolyploidy and inter- and intrachromosomal duplications) provides the

wheat genome with a vast reservoir of functional genes that likely contribute to wheat adaptation.

On the basis of this work, the IWGSC has already defined an adapted BAC pooling strategy to reach the same sequence quality while reducing sequencing costs for the remaining chromosomes. Although progress in sequencing technologies and cost reduction allows for more cost-efficient sequence production, the challenge of bioinformatics and limitations of current sequencing technologies remain (12). Solving these issues and improving methods to efficiently anchor and orient scaffolds within pseudomolecules should make the assembly of high-quality reference sequences of complex genomes routine work in the future. There is no doubt that, as witnessed after the release of the rice genome sequence (49), the number of genes cloned from wheat will grow exponentially in the near future, thereby enabling wheat researchers and breeders to cope with the urgent need to improve wheat yield in the face of climate change and food-security challenges (52).

Materials and methods

Sequencing, assembly, and scaffolding

A total of 8452 BACs representing the MTP of wheat chromosome 3B were pooled into 922 BAC pools. Each pool was used to create a bar-coded Roche/454 8-kb long paired-end library. In total, 150 sequencing runs were performed, leading to an average of 36-fold sequence coverage. After assembly with Newbler (Roche), we integrated 42,551 BAC end sequences to validate and improve scaffolding. Illumina reads generated from sorted DNA of chromosome 3B were used to fill gaps within scaffolds and correct potential sequencing errors remaining in the consensus sequence (18).

Anchoring scaffolds

SNP discovery was performed through sequence capture for 52,265 loci flanking TE junctions representing an average density of one locus per 16.2 kb (18). Out of 39,077 SNPs distributed along the chromosome, a subset of 3075 evenly distributed (38.2 ± 9.4 SNPs per 10 Mb) SNPs was selected to genotype 1025 lines from recombinant inbred and association panels. An anchor genetic map was built first by linkage analysis and integration of linkage disequilibrium data. A consensus map comprising 5318 markers was also built using 40 different genetic maps to anchor additional scaffolds (18). Finally, a position in the pseudomolecule was inferred for scaffolds without marker information but belonging to an anchored physical BAC contig.

Sequence annotation

Gene modeling was performed using an improved version of the TriAnnot pipeline (53). Noncoding RNA genes were predicted using three different programs (18), and predictions were manually curated. Predictions of TEs and reconstruction of the pattern of nested insertions were performed through the development of a specific program (18) that automatically curates similarity-search results obtained with a dedi-

cated databank comprising 4929 known wheat TEs classified into 521 families.

Gene expression analyses

Thirty RNA samples, corresponding to RNAs extracted in duplicates from five organs (root, leaf, stem, spike, and grain) at three developmental stages each from hexaploid wheat cultivar Chinese Spring (4), were used for gene expression analyses. RNA-Seq libraries were constructed using the IlluminaTruSeq (Illumina, CA, USA) RNA sample preparation kit and sequenced. An average of 50 ± 11 million paired-end reads per sample were mapped on the chromosome 3B scaffolds and used to reconstruct transcripts and estimate transcript abundance in units of fragments per kb of exon per million mapped reads (FPKM). Regions with FPKM values higher than zero were considered as expressed.

Distribution and segmentation analyses

Distributions of recombination rate, gene and TE densities, and expression breadth were calculated within a sliding window of 10 Mb (and 1 Mb for the recombination rate), with a step of 1 Mb along the chromosome sequence using a homemade Perl script. Segmentation analyses of these distributions were performed using the R package change-point v1.0.6 (54), with Segment Neighborhoods method and Bayesian information criterion penalty on the mean change.

Comparative genomics, gene duplications, and molecular evolution

We performed an all-by-all Basic Local Alignment Search Tool for Proteins (BLASTP) [cutoff expected value (Evalue), 1×10^{-5}] comparison between the amino-acid sequences of predicted genes of wheat chromosome 3B, rice (Michigan State University version 7.0), *Brachypodium* (*Brachypodium* Sequencing Initiative, 2.0), and sorghum (phytozome, version 1.4). We filtered out genes with no homology with at least one other gene in a compared species (cutoff 35% amino acid identity and 35% sequence overlap). Syntenic genes were defined as genes with a best BLAST hit on an orthologous chromosome in at least one other species. Nonsyntenic genes were defined as genes for which the best BLAST hit was on a nonorthologous chromosome in the other species. Clustering of orthologous and paralogous genes was performed using OrthoMCL [Evalue cutoff, 1×10^{-5} ; percentage match cutoff, 35% (55)]. All 3B genes clustered into the same family were considered intrachromosomal duplicates. 3B genes clustered in a family with wheat gene models annotated on another chromosome (19), not including genes from group 3, were considered as interchromosomal duplicates. Tandem duplicates were defined as genes in the same family with five or fewer spacer genes separating them on the pseudomolecule, and dispersed duplicates were defined as having more than five spacer genes. Synonymous (K_s) and non-synonymous (K_a) substitution rates were calculated based on ClustalW 2.1 (56) coding sequence alignments by the Nei and Gojobori method

using codeml [part of the PAML package (57)]. Age of gene divergence was estimated by the equation $K_2/2r$; where $r = 6.5 \times 10^{-9}$.

REFERENCES AND NOTES

- J. Dubcovsky, J. Dvorak, Genome plasticity a key factor in the success of polyploid wheat under domestication. *Science* **316**, 1862–1866 (2007). doi: [10.1126/science.1143986](https://doi.org/10.1126/science.1143986); pmid: [17600208](https://pubmed.ncbi.nlm.nih.gov/17600208/)
- F. Choulet *et al.*, Megabase level sequencing reveals contrasted organization and evolution patterns of the wheat gene and transposable element spaces. *Plant Cell* **22**, 1686–1701 (2010). doi: [10.1105/tpc.110.074187](https://doi.org/10.1105/tpc.110.074187); pmid: [20581307](https://pubmed.ncbi.nlm.nih.gov/20581307/)
- J. Zhang, Evolution by gene duplication: An update. *Trends Ecol. Evol.* **18**, 292–298 (2003). doi: [10.1016/S0169-5347\(03\)00033-8](https://doi.org/10.1016/S0169-5347(03)00033-8)
- C. Rustenholz *et al.*, A 3,000-loci transcription map of chromosome 3B unravels the structural and functional features of gene islands in hexaploid wheat. *Plant Physiol.* **157**, 1596–1608 (2011). doi: [10.1104/pp.111.183921](https://doi.org/10.1104/pp.111.183921); pmid: [22034626](https://pubmed.ncbi.nlm.nih.gov/22034626/)
- I. D. Wilson *et al.*, A transcriptomics resource for wheat functional genomics. *Plant Biotechnol. J.* **2**, 495–506 (2004). doi: [10.1111/j.1467-7652.2004.00096.x](https://doi.org/10.1111/j.1467-7652.2004.00096.x); pmid: [17147622](https://pubmed.ncbi.nlm.nih.gov/17147622/)
- P. R. Bhat *et al.*, Mapping translocation breakpoints using a wheat microarray. *Nucleic Acids Res.* **35**, 2936–2943 (2007). doi: [10.1093/nar/gkm148](https://doi.org/10.1093/nar/gkm148); pmid: [17439961](https://pubmed.ncbi.nlm.nih.gov/17439961/)
- L. L. Qi *et al.*, 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 (2004). doi: [10.1534/genetics.104.034868](https://doi.org/10.1534/genetics.104.034868); pmid: [15514046](https://pubmed.ncbi.nlm.nih.gov/15514046/)
- R. Brenchley *et al.*, Analysis of the bread wheat genome using whole-genome shotgun sequencing. *Nature* **491**, 705–710 (2012). doi: [10.1038/nature11650](https://doi.org/10.1038/nature11650); pmid: [23192148](https://pubmed.ncbi.nlm.nih.gov/23192148/)
- J. Jia *et al.*, Aegilops tauschii draft genome sequence reveals a gene repertoire for wheat adaptation. *Nature* **496**, 91–95 (2013). doi: [10.1038/nature12028](https://doi.org/10.1038/nature12028); pmid: [23535592](https://pubmed.ncbi.nlm.nih.gov/23535592/)
- H. Q. Ling *et al.*, Draft genome of the wheat A-genome progenitor *Triticum urartu*. *Nature* **496**, 87–90 (2013). doi: [10.1038/nature11997](https://doi.org/10.1038/nature11997); pmid: [23535596](https://pubmed.ncbi.nlm.nih.gov/23535596/)
- M. C. Schatz, A. L. Delcher, S. L. Salzberg, Assembly of large genomes using second-generation sequencing. *Genome Res.* **20**, 1165–1173 (2010). doi: [10.1101/gr.101360.109](https://doi.org/10.1101/gr.101360.109); pmid: [20508146](https://pubmed.ncbi.nlm.nih.gov/20508146/)
- V. Marx, Next-generation sequencing: The genome jigsaw. *Nature* **501**, 263–268 (2013). doi: [10.1038/501261a](https://doi.org/10.1038/501261a); pmid: [24025842](https://pubmed.ncbi.nlm.nih.gov/24025842/)
- P. S. Schnable *et al.*, The B73 maize genome: Complexity, diversity, and dynamics. *Science* **326**, 1112–1115 (2009). doi: [10.1126/science.1178534](https://doi.org/10.1126/science.1178534); pmid: [19965430](https://pubmed.ncbi.nlm.nih.gov/19965430/)
- X. Xu *et al.*, Genome sequence and analysis of the tuber crop potato. *Nature* **475**, 189–195 (2011). doi: [10.1038/nature10158](https://doi.org/10.1038/nature10158); pmid: [21743474](https://pubmed.ncbi.nlm.nih.gov/21743474/)
- J. Doležel, M. Kubaláková, E. Paux, J. Bartos, C. Feuillet, Chromosome-based genomics in the cereals. *Chromosome Res.* **15**, 51–66 (2007). doi: [10.1007/s10577-006-1106-x](https://doi.org/10.1007/s10577-006-1106-x); pmid: [17295126](https://pubmed.ncbi.nlm.nih.gov/17295126/)
- J. Safár *et al.*, Dissecting large and complex genomes: Flow sorting and BAC cloning of individual chromosomes from bread wheat. *Plant J.* **39**, 960–968 (2004). doi: [10.1111/j.1365-3113.2004.02179.x](https://doi.org/10.1111/j.1365-3113.2004.02179.x); pmid: [15341637](https://pubmed.ncbi.nlm.nih.gov/15341637/)
- E. Paux *et al.*, A physical map of the 1-gigabase bread wheat chromosome 3B. *Science* **322**, 101–104 (2008). doi: [10.1126/science.1161847](https://doi.org/10.1126/science.1161847); pmid: [18832645](https://pubmed.ncbi.nlm.nih.gov/18832645/)
- Supplementary materials are available on Science Online.
- International Wheat Genome Sequencing Consortium, A chromosome-based draft sequence of the hexaploid bread wheat (*Triticum aestivum*) genome. *Science* **345**, 1251788 (2014).
- B. S. Gill, B. Friebe, T. Endo, Standard karyotype and nomenclature system for description of chromosome bands and structural aberrations in wheat (*Triticum aestivum*). *Genome* **34**, 830–839 (1991). doi: [10.1139/g91-128](https://doi.org/10.1139/g91-128)
- P. S. Chain *et al.*, Genomics. Genome project standards in a new era of sequencing. *Science* **326**, 236–237 (2009). doi: [10.1126/science.1180614](https://doi.org/10.1126/science.1180614); pmid: [19815760](https://pubmed.ncbi.nlm.nih.gov/19815760/)
- B. Li *et al.*, Wheat centromeric retrotransposons: The new ones take a major role in centromeric structure. *Plant J.* **73**, 952–965 (2013). doi: [10.1111/tpj.12086](https://doi.org/10.1111/tpj.12086); pmid: [23253213](https://pubmed.ncbi.nlm.nih.gov/23253213/)
- H. Yan *et al.*, Intergenic locations of rice centromeric chromatin. *PLoS Biol.* **6**, e286 (2008). doi: [10.1371/journal.pbio.0060286](https://doi.org/10.1371/journal.pbio.0060286); pmid: [19067486](https://pubmed.ncbi.nlm.nih.gov/19067486/)
- H. Zhang, R. K. Dawe, Total centromere size and genome size are strongly correlated in ten grass species. *Chromosome Res.* **20**, 403–412 (2012). doi: [10.1007/s10577-012-9284-1](https://doi.org/10.1007/s10577-012-9284-1); pmid: [22552915](https://pubmed.ncbi.nlm.nih.gov/22552915/)
- T. Sakuno, K. Tada, Y. Watanabe, Kinetochore geometry defined by cohesion within the centromere. *Nature* **458**, 852–858 (2009). doi: [10.1038/nature07876](https://doi.org/10.1038/nature07876); pmid: [19370027](https://pubmed.ncbi.nlm.nih.gov/19370027/)
- Q. Pan, F. Ali, X. Yang, J. Li, J. Yan, Exploring the genetic characteristics of two recombinant inbred line populations via high-density SNP markers in maize. *PLoS ONE* **7**, e52777 (2012). doi: [10.1371/journal.pone.0052777](https://doi.org/10.1371/journal.pone.0052777); pmid: [23300772](https://pubmed.ncbi.nlm.nih.gov/23300772/)
- A. J. Lukaszewski, C. A. Curtis, Physical distribution of recombination in B-genome chromosomes of tetraploid wheat. *Theor. Appl. Genet.* **86**, 121–127 (1993). doi: [10.1007/BF00223816](https://doi.org/10.1007/BF00223816); pmid: [24193391](https://pubmed.ncbi.nlm.nih.gov/24193391/)
- C. Saintenac *et al.*, Detailed recombination studies along chromosome 3B provide new insights on crossover distribution in wheat (*Triticum aestivum* L.). *Genetics* **181**, 393–403 (2009). doi: [10.1534/genetics.108.097469](https://doi.org/10.1534/genetics.108.097469); pmid: [19064706](https://pubmed.ncbi.nlm.nih.gov/19064706/)
- J. Evans *et al.*, Extensive variation in the density and distribution of DNA polymorphism in sorghum genomes. *PLoS ONE* **8**, e79192 (2013). doi: [10.1371/journal.pone.0079192](https://doi.org/10.1371/journal.pone.0079192); pmid: [24265758](https://pubmed.ncbi.nlm.nih.gov/24265758/)
- International Rice Genome Sequencing Project, The map-based sequence of the rice genome. *Nature* **436**, 793–800 (2005). doi: [10.1038/nature03895](https://doi.org/10.1038/nature03895); pmid: [16100779](https://pubmed.ncbi.nlm.nih.gov/16100779/)
- A. H. Paterson *et al.*, The Sorghum bicolor genome and the diversification of grasses. *Nature* **457**, 551–556 (2009). doi: [10.1038/nature07723](https://doi.org/10.1038/nature07723); pmid: [19189423](https://pubmed.ncbi.nlm.nih.gov/19189423/)
- A. Gottlieb *et al.*, Insular organization of gene space in grass genomes. *PLoS ONE* **8**, e54101 (2013). doi: [10.1371/journal.pone.0054101](https://doi.org/10.1371/journal.pone.0054101); pmid: [23326580](https://pubmed.ncbi.nlm.nih.gov/23326580/)
- M. W. Ganal *et al.*, A large maize (*Zea mays* L.) SNP genotyping array: Development and germplasm genotyping, and genetic mapping to compare with the B73 reference genome. *PLoS ONE* **6**, e28334 (2011). doi: [10.1371/journal.pone.0028334](https://doi.org/10.1371/journal.pone.0028334); pmid: [22174790](https://pubmed.ncbi.nlm.nih.gov/22174790/)
- K. F. Mayer *et al.*, A physical, genetic and functional sequence assembly of the barley genome. *Nature* **491**, 711–716 (2012). pmid: [23075845](https://pubmed.ncbi.nlm.nih.gov/23075845/)
- R. S. Sekhon *et al.*, Maize gene atlas developed by RNA sequencing and comparative evaluation of transcriptomes based on RNA sequencing and microarrays. *PLoS ONE* **8**, e61005 (2013). doi: [10.1371/journal.pone.0061005](https://doi.org/10.1371/journal.pone.0061005); pmid: [23637782](https://pubmed.ncbi.nlm.nih.gov/23637782/)
- R. S. Baucum *et al.*, Exceptional diversity, non-random distribution, and rapid evolution of retroelements in the B73 maize genome. *PLoS Genet.* **5**, e1000732 (2009). doi: [10.1371/journal.pgen.1000732](https://doi.org/10.1371/journal.pgen.1000732); pmid: [19936065](https://pubmed.ncbi.nlm.nih.gov/19936065/)
- R. S. Baucum, J. C. Estill, J. Leebens-Mack, J. L. Bennetzen, Natural selection on gene function drives the evolution of LTR retrotransposon families in the rice genome. *Genome Res.* **19**, 243–254 (2009). doi: [10.1101/gr.083360.108](https://doi.org/10.1101/gr.083360.108); pmid: [19029538](https://pubmed.ncbi.nlm.nih.gov/19029538/)
- M. Charles *et al.*, Dynamics and differential proliferation of transposable elements during the evolution of the B and A genomes of wheat. *Genetics* **180**, 1071–1086 (2008). doi: [10.1534/genetics.108.092304](https://doi.org/10.1534/genetics.108.092304); pmid: [18780739](https://pubmed.ncbi.nlm.nih.gov/18780739/)
- E. M. Sergeeva, E. A. Salina, I. G. Adonina, B. Chalhou, Evolutionary analysis of the CACTA DNA-transposon Caspar across wheat species using sequence comparison and in situ hybridization. *Mol. Genet. Genomics* **284**, 11–23 (2010). doi: [10.1007/s00438-010-0544-5](https://doi.org/10.1007/s00438-010-0544-5); pmid: [2012353](https://pubmed.ncbi.nlm.nih.gov/2012353/)
- C. Lu *et al.*, Miniature inverted-repeat transposable elements (MITEs) have been accumulated through amplification bursts and play important roles in gene expression and species diversity in *Oryza sativa*. *Mol. Biol. Evol.* **29**, 1005–1017 (2012). doi: [10.1093/molbev/msr282](https://doi.org/10.1093/molbev/msr282); pmid: [22096216](https://pubmed.ncbi.nlm.nih.gov/22096216/)
- M. D. Gale, K. M. Devos, Comparative genetics in the grasses. *Proc. Natl. Acad. Sci. U.S.A.* **95**, 1971–1974 (1998). doi: [10.1073/pnas.95.5.1971](https://doi.org/10.1073/pnas.95.5.1971); pmid: [9482816](https://pubmed.ncbi.nlm.nih.gov/9482816/)
- K. M. Devos, M. D. Gale, Genome relationships: The grass model in current research. *Plant Cell* **12**, 637–646 (2000). doi: [10.1105/tpc.12.5.637](https://doi.org/10.1105/tpc.12.5.637); pmid: [10810140](https://pubmed.ncbi.nlm.nih.gov/10810140/)
- F. Murat *et al.*, Ancestral grass karyotype reconstruction unravels new mechanisms of genome shuffling as a source of plant evolution. *Genome Res.* **20**, 1545–1557 (2010). doi: [10.1101/gr.109744.110](https://doi.org/10.1101/gr.109744.110); pmid: [20876790](https://pubmed.ncbi.nlm.nih.gov/20876790/)
- International Brachypodium Initiative, Genome sequencing and analysis of the model grass *Brachypodium distachyon*. *Nature* **463**, 763–768 (2010). doi: [10.1038/nature08747](https://doi.org/10.1038/nature08747); pmid: [20148030](https://pubmed.ncbi.nlm.nih.gov/20148030/)
- E. D. Akhunov *et al.*, Synteny perturbations between wheat homoeologous chromosomes caused by locus duplications and deletions correlate with recombination rates. *Proc. Natl. Acad. Sci. U.S.A.* **100**, 10836–10841 (2003). doi: [10.1073/pnas.1934431100](https://doi.org/10.1073/pnas.1934431100); pmid: [12960374](https://pubmed.ncbi.nlm.nih.gov/12960374/)
- T. Wicker *et al.*, Frequent gene movement and pseudogene evolution is common to the large and complex genomes of wheat, barley, and their relatives. *Plant Cell* **23**, 1706–1718 (2011). doi: [10.1105/tpc.111.086629](https://doi.org/10.1105/tpc.111.086629); pmid: [21622801](https://pubmed.ncbi.nlm.nih.gov/21622801/)
- T. Wicker, J. P. Buchmann, B. Keller, Patching gaps in plant genomes results in gene movement and erosion of colinearity. *Genome Res.* **20**, 1229–1237 (2010). doi: [10.1101/gr.107284.110](https://doi.org/10.1101/gr.107284.110); pmid: [20530251](https://pubmed.ncbi.nlm.nih.gov/20530251/)
- M. Morgante *et al.*, Gene duplication and exon shuffling by helitron-like transposons generate intraspecific diversity in maize. *Nat. Genet.* **37**, 997–1002 (2005). doi: [10.1038/ng1615](https://doi.org/10.1038/ng1615); pmid: [16056225](https://pubmed.ncbi.nlm.nih.gov/16056225/)
- C. Feuillet, J. E. Leach, J. Rogers, P. S. Schnable, K. Eversole, Crop genome sequencing: Lessons and rationales. *Trends Plant Sci.* **16**, 77–88 (2011). doi: [10.1016/j.tplants.2010.10.005](https://doi.org/10.1016/j.tplants.2010.10.005); pmid: [21081278](https://pubmed.ncbi.nlm.nih.gov/21081278/)
- E. Paux *et al.*, Insertion site-based polymorphism markers open new perspectives for genome saturation and marker-assisted selection in wheat. *Plant Biotechnol. J.* **8**, 196–210 (2010). doi: [10.1111/j.1467-7652.2009.00477.x](https://doi.org/10.1111/j.1467-7652.2009.00477.x); pmid: [20078842](https://pubmed.ncbi.nlm.nih.gov/20078842/)
- B. Goffinet, S. Gerber, Quantitative trait loci: A meta-analysis. *Genetics* **155**, 463–473 (2000). pmid: [10790417](https://pubmed.ncbi.nlm.nih.gov/10790417/)
- J. A. Foley *et al.*, Solutions for a cultivated planet. *Nature* **478**, 337–342 (2011). doi: [10.1038/nature10452](https://doi.org/10.1038/nature10452); pmid: [21993620](https://pubmed.ncbi.nlm.nih.gov/21993620/)
- P. Leroy *et al.*, TriAnnot: A versatile and high performance pipeline for the automated annotation of plant genomes. *Front. Plant Sci.* **3**, 5 (2012). doi: [10.3389/fpls.2012.00005](https://doi.org/10.3389/fpls.2012.00005); pmid: [22645565](https://pubmed.ncbi.nlm.nih.gov/22645565/)
- J. Chen, A. K. Gupta, *Parametric Statistical Change Point Analysis: With Applications to Genetics, Medicine, and Finance* (Birkhäuser, Basel, 2012).
- L. Li, C. J. Stoekert Jr., D. S. Roos, OrthoMCL: Identification of ortholog groups for eukaryotic genomes. *Genome Res.* **13**, 2178–2189 (2003). doi: [10.1101/gr.1224503](https://doi.org/10.1101/gr.1224503); pmid: [12952885](https://pubmed.ncbi.nlm.nih.gov/12952885/)
- J. D. Thompson, D. G. Higgins, T. J. Gibson, CLUSTAL W: Improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**, 4673–4680 (1994). doi: [10.1093/nar/22.22.4673](https://doi.org/10.1093/nar/22.22.4673); pmid: [7984417](https://pubmed.ncbi.nlm.nih.gov/7984417/)
- Z. Yang, PAML: A program package for phylogenetic analysis by maximum likelihood. *CABIOS* **13**, 555–556 (1997). pmid: [9367129](https://pubmed.ncbi.nlm.nih.gov/9367129/)

ACKNOWLEDGMENTS

The authors thank the scientific advisory board (P. Schnable, S. Rounsley, D. Ware, J. Rogers, and K. Eversole) of the 3BSEQ project for fruitful discussions; K. Eversole for critical reading and editing of the manuscript; H. Rimbart, N. Cubizolles, and E. Rey for SNP marker discovery and genotyping; M. Kubaláková and J. Vrána for assistance with the preparation of DNA amplified from flow-sorted chromosome 3B; L. Couderc, A. Keliet, and S. Reboux for their support in database and system administration; and C. Poncet and the "Plateforme GENTYANE" for SNP genotyping. This work was supported by a grant from the French National Research Agency (ANR-09-GENM-025 3BSEQ), a grant from France Agrimer, and a grant (project DL-7E) from the INRA BAP Biologie et Amélioration des Plantes division. N.G. is funded by a grant from the European Commission research training program Marie-Curie Actions (FP7-MC-IF-NoncollinearGenes). J. Daron is funded by a grant from the French Ministry of Research. L.P. is funded by a grant from the Region Auvergne. K.V. is supported by the Ghent University Multidisciplinary Research Partnership ["Bioinformatics: From nucleotides to networks" (Project 01MR0310W)]. J. Dolezel is supported by the Czech Science Foundation (award no. P501/12/G090). The chromosome 3B BAC library and the pools of the MTP BAC clones are available upon request under a materials transfer agreement at the French Plant Genomic Center, INRA-Centre National de Ressources Génétiques Végétales. Annotation data and browser are available at https://urgi.versailles.inra.fr/gb2/gbrowse/wheat_annot_3B. Sequences and annotations of the reference pseudomolecule and unassigned scaffolds have been deposited in ENA (project PRJEB4376) under accession numbers HG670306 and CBUC010000001 to CBUC010001450, respectively. RNA-Seq data were deposited under accession number ERP004714.

SUPPLEMENTARY MATERIALS

www.sciencemag.org/content/345/6194/1249721/suppl/DC1
Materials and Methods
Figs. S1 to S11
Tables S1 to S14
References (58–125)

13 December 2013; accepted 30 May 2014
10.1126/science.1249721

Structural and functional partitioning of bread wheat chromosome 3B

Frédéric Choulet, Adriana Alberti, Sébastien Theil, Natasha Glover, Valérie Barbe, Josquin Daron, Lise Pingault, Pierre Sourdille, Arnaud Couloux, Etienne Paux, Philippe Leroy, Sophie Mangenot, Nicolas Guilhot, Jacques Le Gouis, Francois Balfourier, Michael Alaux, Véronique Jamilloux, Julie Poulain, Céline Durand, Arnaud Bellec, Christine Gaspin, Jan Safar, Jaroslav Dolezel, Jane Rogers, Klaas Vandepoele, Jean-Marc Aury, Klaus Mayer, Hélène Berges, Hadi Quesneville, Patrick Wincker, and Catherine Feuillet

Science, 345 (6194), 1249721. • DOI: 10.1126/science.1249721

View the article online

<https://www.science.org/doi/10.1126/science.1249721>

Permissions

<https://www.science.org/help/reprints-and-permissions>