# Long-read and chromosome-scale assembly of the hexaploid wheat genome achieves high resolution for research and breeding

Jean-Marc Aury, Stefan Engelen, Benjamin Istace, Cécile Monat, Pauline Lasserre-Zuber, Caroline Belser, Corinne Cruaud, Hélène Rimbert, Philippe Leroy, Sandrine Arribat, et al.

HAL Id: hal-03658639
https://hal.inrae.fr/hal-03658639v1

Submitted on 5 May 2022

# Long-read and chromosome-scale assembly of the hexaploid wheat genome achieves high resolution for research and breeding

Jean-Marc Aury [ID]1,*, Stefan Engelen [ID]1, Benjamin Istace [ID]1, Cécile Monat [ID]2, Pauline Lasserre-Zuber2, Caroline Belser [ID]1, Corinne Cruaud [ID]3, Hélène Rimbert [ID]2, Philippe Leroy [ID]2, Sandrine Arribat [ID]4, Isabelle Dufau4, Arnaud Bellec [ID]4, David Grimbichler [ID]5, Nathan Papon2, Etienne Paux [ID]2, Marion Ranoux2, Adriana Alberti [ID]1,6, Patrick Wincker [ID]1 and Frédéric Choulet [ID]2,*

1Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université Paris-Saclay, 91057 Evry, France
2GDEC, Université Clermont Auvergne, INRAE, UMR1095, 63000 Clermont-Ferrand, France
3Commissariat à l'Energie Atomique (CEA), Institut François Jacob, Genoscope, F-91057 Evry, France
4INRAE, CNRGV French Plant Genomic Resource Center, F-31320, Castanet Tolosan, France
5Mésocentre Clermont Auvergne, DOSI / Bâtiment Turing, 7 avenue Blaise Pascal, 63178 Aubière, France
6Present address: Université Paris-Saclay, CEA, CNRS, Institute for Integrative Biology of the Cell (I2BC), 91198, Gif-sur-Yvette, France
*Correspondence address. Jean-Marc Aury, Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université Paris-Saclay, 91057 Evry, France. E-mail: jmaury@genoscope.cns.fr; Frédéric Choulet, GDEC, Université Clermont Auvergne, INRAE, UMR1095, 63000 Clermont-Ferrand, France. E-mail: frederic.choulet@inrae.fr

## Abstract

**Background:** The sequencing of the wheat (*Triticum aestivum*) genome has been a methodological challenge for many years owing to its large size (15.5 Gb), repeat content, and hexaploidy. Many initiatives aiming at obtaining a reference genome of cultivar Chinese Spring have been launched in the past years and it was achieved in 2018 as the result of a huge effort to combine short-read sequencing with many other resources. Reference-quality genome assemblies were then produced for other accessions, but the rapid evolution of sequencing technologies offers opportunities to reach high-quality standards at lower cost. **Results:** Here, we report on an optimized procedure based on long reads produced on the Oxford Nanopore Technology PromethION device to assemble the genome of the French bread wheat cultivar Renan. **Conclusions:** We provide the most contiguous chromosome-scale assembly of a bread wheat genome to date. Coupled with an annotation based on RNA-sequencing data, this resource will be valuable for the crop community and will facilitate the rapid selection of agronomically important traits. We also provide a framework to generate high-quality assemblies of complex genomes using ONT.

**Keywords:** wheat, hexaploid genome, long-reads, nanopore sequencing, genome assembly, haplotype characterization, introgressions

## Introduction

Bread wheat (*Triticum aestivum*) is among the most important cereal crops, and a better knowledge in the area of wheat genomics is needed to face the main challenge of ensuring food security to a growing population in the context of climate change. Improving productivity requires both that local producers adapt their practices to increase their climate resilience and a better understanding of the wheat production systems. In this context, a better knowledge of the wheat genome and its gene content, but also the sequencing of numerous accessions, are essential.

However, the genome of bread wheat is particularly characterized by its complexity. Indeed this hexaploid genome is the result of 2 interspecific hybridization events. The earliest cultivated wheat was diploid, but humans have intensified the cultivation of polyploid species. Recent studies show that these polyploid species appear to be advantaged by their genomic plasticity [1]. Indeed, modifications of the gene space and related elements are buffered by the polyploid nature of wheat and open a wider field

to selection. Bread wheat is composed of 3 subgenomes A, B, and D derived from 3 ancestral diploid species that diverged between 2.5 and 6 million years ago [2].

The wheat genome is one of the largest among sequenced plant genomes (15.5 Gb), mainly composed of repetitive sequences (>85%), and contains many homoeologous regions between the 3 genomes (A, B, and D). Repetitive sequences and polyploidy pose serious challenges in the generation of genome assemblies. The adventure of sequencing the hexaploid wheat genome began in 2005 with the creation of the International Wheat Genome Sequencing Consortium (IWGSC) [3]. With the advent of sequencing technologies, the wheat genome has been competitively sequenced several times [4–6]. The first reference-quality genome sequence with a comprehensive annotation was published by the IWGSC in August 2018 [7] for the accession Chinese Spring (CS), referred to hereinafter as CS RefSeq v2.1 or simply CS. This assembly represents a tremendous resource for the scientific community and offers the promise of facilitating and accelerating breeding efforts.

More recently, 15 genomes of hexaploid wheat have been published [8], which represents a new step in the knowledge of the wheat model. Ten of these new wheat genomes have been assembled at the chromosome level, allowing for comparative analysis on a scale that was previously impossible. While being a valuable and highly validated resource using multiple technologies, these assemblies were produced using short-read technologies and therefore may contain a higher number of gaps compared to genomes assembled with long reads [9–13]. In 2017, an assembly of the CS genome using long reads was produced [5], although not annotated, highlighting the added value of long reads in such complex genomes. By accumulating long-read assemblies, the scientific community is now aware of the flaw in short-read strategies. Indeed they underestimate the repetitive content of the genome and more importantly can lack tandemly duplicated genes [14, 15]. Several years ago, Pacific Biosciences (PacBio) and Oxford Nanopore (ONT) sequencing technologies were commercialized with the promise to sequence long DNA fragments and revolutionize complex genome assemblies.

Here, we report the first hexaploid wheat genome based on ONT long reads. We sequenced the genome of the French variety Renan, one of the most used varieties in organic farming. The Renan genome carries multiple resistance genes against fungal pathogens (leaf rust, stem rust, yellow rust, eyespot) originating from introgression of DNA regions coming from the wild species *Aegilops ventricosa*. We used the PromethION device and organized the assembled contigs at the chromosome scale using optical maps (BioNano Genomics [BNG]) and Hi-C libraries (Arima Genomics [AG]). This assembly has a contig N50 of 2.2 Mb, which is a 30-fold improvement over existing chromosome-scale assemblies.

## Results

### Genome sequencing and optical maps

We sequenced genomic DNA using 20 ONT flow cells (2 MinION and 18 PromethION), which produced 12M reads representing 1.1 Tb. All the reads were originally base called using the guppy 2.0 software, but given the improvement of guppy software during our project, we decided to call bases using a newer version of the guppy software (version 3.6 with High Accuracy setting). This dataset represented a coverage of 63× of the hexaploid wheat genome, and the read N50 was of 24.6 kb. More importantly, we got 3.1M reads larger than 50 kb, representing a 14× genome coverage (Supplementary Table S1). In addition, we generated Illumina short reads and long-range data for, respectively, polishing and organizing nanopore contigs. We produced an optical map using the Saphyr instrument commercialized by BNG. High molecular weight (HMW) DNA was extracted and labelled using the Direct Label and Stain Chemistry (DLS) with the DLE-1 enzyme. The DLE-1 optical map was assembled using proprietary tools provided by BNG and had a cumulative size of 14.9 Gb with an N50 of 37.5 Mb (Supplementary Table S2). Four Hi-C libraries from 2 biological replicates were prepared using the AG protocol and sequenced on an Illumina sequencer to reach 537 Gb, i.e., a depth of 35×. We used a sample of 240 million read pairs (72 Gb, 5×) to build a Hi-C map.

### Genome assembly

Because the dataset was too large for many long-read assemblers, we sampled a 30× coverage by selecting the longest reads (Supplementary Table S1). This subset was assembled using multiple assembly tools dedicated to processing this large amount of data (Redbean [16], SMARTdenovo [17], and Flye [18]). SMARTdenovo is not among the fastest algorithms and has not been updated for several years, but because it can be easily parallelized, it remains an interesting choice for assembling large genomes. The overlap and consensus calculations were split into 60 chunks, and each were run on a 32-core server and took ∼2 days and 10 hours, respectively. In comparison, Redbean was able to generate an assembly after just 7s on a 64-core server with 3 TB of memory while Flye needed 43 days on the same computer server. Surprisingly, the redbean assembly had a cumulative size 2 times higher than the expected genome size (29.6 vs 14.5 Gb), a low contiguity and, contained a large amount of short contigs. The SMARTdenovo and Flye assemblies were highly comparable, but Flye was the most contiguous (contigs N50 of 1.8 vs 1.1 Mb) and SMARTdenovo had a cumulative size closer to the expected one (14.1 vs 13.0 Gb, Supplementary Table S3). Additionally, even though the assemblies were polished later, the raw SMARTdenovo assembly contained a higher number of complete BUSCO genes (83.0% vs 49.5%), which indicates that its consensus module is more efficient.

The SMARTdenovo and Flye assemblies were successively polished using Racon [19] and Medaka [20] with long reads and Hapo-G [21] with short reads. Polished contigs were validated and organized into scaffolds using the DLE-1 optical map and proprietary tools provided by BNG. As expected, owing to its lower cumulative size, Flye scaffolds contained a larger proportion of unknown bases (851 and 262 Mb). On the basis of these results (proportion of gaps and gene completion), the assembly produced by SMARTdenovo [17] was selected (Supplementary Table S4). Local contig duplications (negative gaps) were resolved using BiSCoT [44], which improved the contigs N50 from 1.2 Mb up to 2.1 Mb. Finally, the resulting assembly was polished 1 last time using Hapo-G [21] with short reads. This led to 2,904 scaffolds (larger than 30 kb) representing 14.26 Gb with an N50 of 48 Mb (79 scaffolds) and a maximum scaffold size of 254 Mb. Thus, the genome size is in the same range as all other available reference quality assemblies of *T. aestivum*: e.g., 14.29 Gb for *cv.* LongReach Lancer, 14.55 Gb for *cv.* Chinese Spring, and 14.96 Gb for *cv.* SY Mattis.

### Construction and validation of pseudomolecules

We then guided the construction of the 21 chromosome sequences (i.e., pseudomolecules) based on collinearity with the Chinese Spring (CS) RefSeq Assembly v2.1 [22]. Given the complexity of this hexaploid genome, we established a dedicated approach to anchor each Renan scaffold based on similarity search against CS. To avoid problems due to multiple mappings, we selected a dataset of uniquely mappable sequences. Genes are not uniquely mappable because most of them are repeated as 3 homoeologous copies sharing on average 97% nucleotide identity. In addition, the gene density (1 gene every 130 kb on average) is too low to anchor small Renan scaffolds that do not carry genes. Thus, we used 150-bp tags corresponding to the 5′ and 3′ junctions between a transposable element (TE) and its insertion site (75 bp on each side), which are called ISBP (insertion site–based polymorphism) markers and are highly abundant and uniquely mappable in the wheat genome [23]. We designed a dataset of 5.76 million ISBPs from the CS assembly, which represent 1 ISBP every 2.5 kb. Their mapping enabled the anchoring of 2,566 scaffolds on 21 pseudomolecules representing 14.20 Gb (99% of the assembly). We then used Hi-C data to validate the assembly and to correct the mis-ordered and mis-oriented scaffolds. The Hi-C map revealed only a few inconsistencies, demonstrating that the

collinearity between CS and Renan was strong enough to guide the anchoring in a very accurate manner. The Hi-C map-based curation led to the detection of 18 chimeric scaffolds that were split into 2 or 3 pieces and to the correction of the location and/or orientation of 198 scaffolds. The final assembly was composed of 21 pseudomolecules (Fig. 1) with 338 unanchored scaffolds representing 61 Mb only.

## Quality assessment of the assembly

First, we calculated the overall quality of the sequence using Merqury and Illumina reads. We obtained an average quality value (QV) of 32.8, a lower QV than that obtained with short-read assemblies but consistent with QV already reported for plant genomes sequenced by ONT [24]. Indeed, using Illumina reads and the CS RefSeq v2.1 assembly, Merqury computed a QV of 44.5 (Table 1). This shows that per-base quality is still an issue, at least with the version of the technology used in this study. However, this could be tempered by the fact that coding regions have higher precision. Indeed, exons longer than 150 bases have a QV of 35.9 (i.e., 1 error every 3.9 kb against 1 error every 2.0 kb in the whole genome).

The completeness and quality of the assembly were estimated by searching for the presence of known genes, i.e., the 107,891 high-confidence genes predicted in CS RefSeq v1.1. We used BLAST [25] to search for the presence of each of the 461,476 exons larger than 30 bp in the Renan scaffolds, and we considered only matches showing ≥95% identity over ≥95% query length. We found hits for 96.2% of the query exons with on average 99.3% identity, suggesting that the gene space is assembled at a high-quality level. The missing genes/exons would correspond, in most of the cases, to real presence/absence variations between CS and Renan, while the nucleotide divergence between exons is 0.7%. It was the first evidence that homoeologous gene copies, sharing on average 97% identity [7], were not collapsed in the Renan assembly. We confirmed this by showing that 62% of the CS exons are strictly identical in Renan (and carried by the same chromosome). Such a level of nucleotide divergence between CS and Renan is similar to what has been shown through whole-genome alignments [26].

We then assessed the assembly quality of the TE space by aligning the complete dataset of ISBP markers of CS onto the Renan assembly. We found that 94% markers were conserved (≥90% identity over 90% query length), i.e., present in the assembly, revealing that the TE space is extremely close to completeness. Indeed, 6% of missing markers is similar to the proportion of expected presence-absence variations affecting TEs [27].

Additionally, we searched for telomeric repeats (TTTAGGG) in the 21 chromosomes and found telomeric repeats at both ends of chromosome 7A, which is generally an indicator of the completion of the chromosome sequence. Both ends of chromosome 7A were also validated by the optical map (Supplementary Fig. S1).

## Impact of the polishing

Based on BUSCO and the alignment of the ISBP markers from the CS assembly, we monitored the evolution of the consensus quality through successive polishing iterations. As previously described, the SMARTdenovo consensus allowed the recovery of a greater number of complete BUSCO genes compared to that of Flye, which may be an indicator of its greater accuracy. However, the BUSCO score was still low (83%) especially for a hexaploid genome, underlining the importance of polishing raw assemblies. Likewise, we were able to find 80.4% of the ISBP markers but only 7% were

aligned without mismatch between the 2 genotypes (Supplementary Table S5). When polished with long reads, the BUSCO score reached 96.7% and 92.9% of the ISBP markers were retrieved (including 28.0% with perfect matches). The subsequent polishing step with short reads weakly decreased the BUSCO score (from 96.7% to 96.6%), but the proportion of duplicated genes increased from 83.1% to 87.0%, which is here wanted because in the case of a hexaploid genome most of the genes are in 3 copies. Moreover, the proportion of perfectly aligned ISBP markers drastically increased from 28.0% to 58.9%. Although the polishing with short reads weakly affects the BUSCO conserved genes, the ISBP markers underline its importance in the case of long-read assemblies. Because ISBPs are unique tags sampling the whole genome, this analysis revealed that nucleotide errors were frequent before polishing, affecting half of the sample loci. Thus, we showed that the polishing steps were successful, even in this large and polyploid genome, and drastically improved the quality of the consensus.
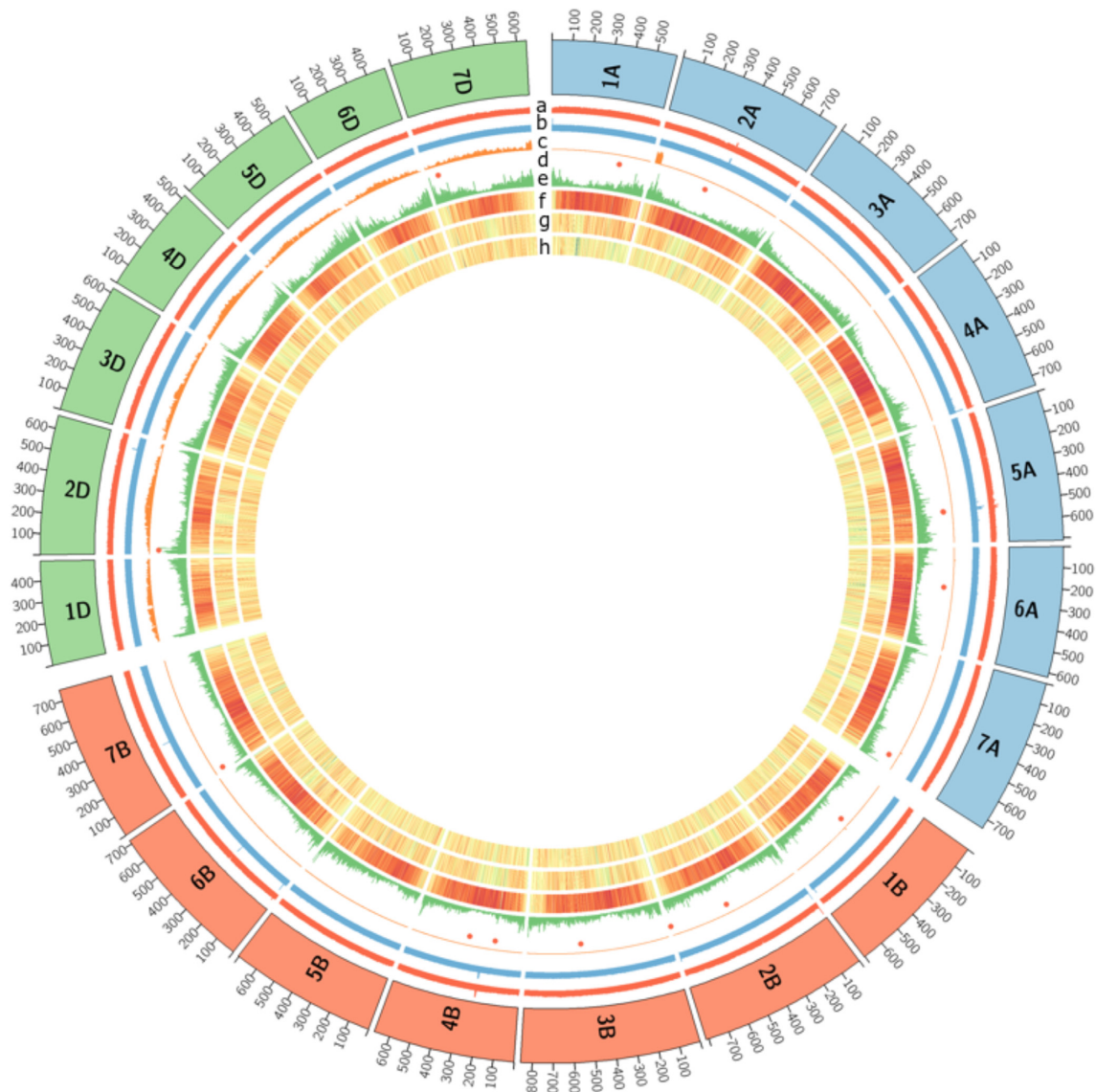
## Recent improvement of the ONT technology

The ONT technology is evolving rapidly, and improvements to the base-calling software are frequent, allowing old data to be analysed with the aim of improving read accuracy and subsequent analysis. To measure the gain brought by each new version during this project, we analysed a subset of ultra-long reads (longer than 100 kb) with different basecallers or versions of the same basecaller: guppy 2.0, guppy 3.0.3 (High Accuracy mode), guppy 3.6 (High Accuracy mode), and the recent bonito v0.3.1. We observed a strong difference in accuracy, of ~7%, between guppy 2.0 and the newer basecaller (bonito v0.3.1), representing the gain over the past 2 years (Supplementary Fig. S2a). This significant improvement could lead nanopore users to reanalyse their old sequencing data to improve the quality of their assemblies. As an example, the accuracy of raw nanopore reads gained ~2% on average using guppy 3.6 (Supplementary Table S6). We observed a reduction of the number of contigs of 19% and an improvement of the contig N50 of 26%. Likewise, the cumulative size is slightly higher in the guppy 3.6 assembly, which may underline a smaller amount of collapsed repetitive regions (Supplementary Table S7).

More importantly, the identity percentage obtained when aligning ONT reads on the wheat assembly is lower than what was obtained on yeast and human samples (Supplementary Fig. S2b). This difference can be explained by the fact that, first, the consensus of the wheat genome is not perfect and second, that basecallers are trained on a mixture that contains yeast and human data. Indeed, DNA modification patterns can differ between taxa, and read accuracy seems better when the model was trained on native DNA from the same species [28]. This huge difference between the read accuracy of yeast and wheat samples should motivate nanopore users to train basecaller models to their targeted species.

## Annotation of transposable elements and protein-coding genes

We annotated TEs on the basis of similarity search against our wheat-specific TE library ClariTeRep [29], and raw results were then refined using CLARITE, a homemade program able to resolve prediction conflicts, merge adjacent features into a single complete element, and identify nested insertion patterns. We detected 3.9 million copies of TEs in the Renan genome assembly, representing 12.0 Gb, i.e., 84% of the assembly size. The proportions of each superfamily were similar to what has been described for CS [30] (Table 2). Gene annotation was achieved by, first, transferring

**Figure 1:** Genome overview of the 21 chromosomes of hexaploid *T. aestivum* Renan (the 7A chromosomes are in blue, the 7B chromosomes in orange, and the 7D chromosomes in green). From inner to outer track: (a) Coverage with short reads, (b) Coverage with long reads, (c) coverage with *A. ventricosa* short reads, (d) red dots represent large deletions (>500 kb), (e) gene density, (f) density of CACTA (DNA transposon) elements, (g) density of Copia elements, (h) density of Gypsy elements. All densities and coverage are calculated in 1-Mb windows; yellow and red colours in density plots indicate lower and higher values, respectively.

genes predicted in CS RefSeq v2.1 by homology using the MAG-ATT pipeline [22]. This allowed us to accurately transfer 105,243 (of 106,801 [98.5%]) high-confidence genes and 155,021 (of 159,846 [97%]) low-confidence genes. Such a transfer of genes predicted in another genotype (here CS) avoided genome-wide *de novo* gene prediction that may artificially lead to many differences between the annotations. We thus focused *de novo* predictions using Tri-Annot [31] only on the unannotated part of the genome, representing 8.5% of the 14.2 Gb, after having masked transferred genes and predicted TEs. For that purpose, we produced RNASeq data for Renan from 28 samples corresponding to 14 different organs/conditions in replicates: grains at 4 developmental stages (100, 250, 500, and 700 degree days) under heat stress and control conditions, stems at 2 developmental stages, leaves at 3 stages, and roots at 1 stage, representing on average 78.8 million read-pairs per sample, i.e., 2.2 billion read-pairs in total. This method allowed us to predict 4,440 genes specific to Renan compared to CS,

i.e., 4% of the gene complement. This is consistent with the extent of structural variations affecting genomes of *Triticeae* [27]. Transfer of known genes, novel predictions, and manual curation (limited to storage protein encoding genes) led us to annotate 109,552 protein-coding genes on the Renan pseudomolecules.

## Comparison with existing hexaploid genome assemblies

We compared our long-read assembly with 10 other available chromosome-scale assemblies of wheat genomes. Although the gene content was similar between the different assemblies, as expected, the assemblies based on short reads had a lower contiguity (contig N50 values <100 kb compared to the 2 Mb of the assembly of the Renan genome, Fig. 2A and B). Logically, they also contained more gaps (~40 times, Fig. 2C). Interestingly, we found more gaps per Mb in the D subgenome compared to the A and B subgenomes in Renan (Supplementary Fig. S3). This indicates

**Table 1:** Comparison of *Triticum aestivum L.* genome assemblies

| Parameter | Renan (this study) | Chinese Spring RefSeq_v2.1 (Zhu et al. [22]) |
|---|---|---|
| No. of contigs | 12,982 | 306,746 |
| Cumulative size (bp) | 14,001,122,256 | 14,317,423,665 |
| N50 (bp) | 2,159,703 | 341,062 |
| L50 | 1,958 | 12,223 |
| N90 (bp) | 598,285 | 32,302 |
| L90 | 6,645 | 59,261 |
| NG50* (bp) | 1,973,000 | 322,161 |
| LG50 | 2,202 | 13,254 |
| NG90* (bp) | 264,272 | 16,550 |
| LG90 | 8,816 | 85,688 |
| Longest contig (bp) | 15,116,687 | 3,528,546 |
| No. of chromosomes | 21 | 21 |
| Cumulative size (bp) | 14,195,643,615 | 14,225,829,371 |
| N50[1] (bp) | 703,299,328 | 713,360,512 |
| L50 | 10 | 10 |
| N90[1] (bp) | 520,815,552 | 518,332,608 |
| L90 | 19 | 19 |
| Longest (bp) | 854,463,248 | 851,934,019 |
| % of N | 1.78 | 1.52 |
| BUSCO on assemblies (%) (N = 4,896) | | |
| Complete | 99.1 | 99.3 |
| Duplicated | 94.7 | 96.1 |
| Fragmented | 0.1 | 0.1 |
| Missing | 0.8 | 0.6 |
| Base accuracy—quality value (*k*-mer) | 32.8 | 44.5 |
| No. of genes | 109,552 | 107,891 |
| Mean No. of exons | 5.10 | 5.33 |
| BUSCO on gene predictions (%) (N = 4,896) | | |
| Complete | 99.1 | 99.5 |
| Duplicated | 94.6 | 98.2 |
| Fragmented | 0.2 | 0.1 |
| Missing | 0.7 | 0.4 |

[1]Calculated using a genome size of 15 Gb.

that the D subgenome is more difficult to assemble even though it has a smaller genome size and contains fewer repetitive elements. The same trend was already observed in another polyploid genome, the rapeseed and its 2 subgenomes A and C [11]. Chromosomes from the different assemblies had similar length except for the Arina*LrFor* and the SY_Mattis variety, in which a translocation has been previously described between chromosomes 5B and 7B [8] (Fig. 2D).

In addition, we generated dot plots between CS and Renan homeologous chromosomes and confirmed the strong collinearity between the 2 genomes (Fig. 3). Whole-chromosome alignments highlighted 16 large-scale inversions (>5 Mb; up to 118 Mb) on 10 chromosomes and 1 translocation of a ~45 Mb segment on chromosome 4A. We performed the same comparisons with the 10 other available genomes of related varieties assembled at the pseudomolecule level (Supplementary Data S1). It showed that only 2 of these inversions are specific to Renan while the others are shared between several accessions. They correspond to

regions of 23 Mb on chr6B (position 398–421 Mb) and 10 Mb on chr7B (position 267–277 Mb).

## Haplotype characterization

Crop breeding involves the selection of desired traits and their combination to generate improved genotypes. Generally, these traits correspond to genomic regions carrying genetic variations or genes [26]. These regions of interest are inherited from their parents in the form of large genomic blocks. The availability of several assemblies of the wheat genome now allows the detection of these haplotypic blocks. Using the 11 chromosome-scale wheat assemblies and an approach based on coloured de Bruijn graphs, we investigated these haplotypic blocks and applied our method to the 21 chromosomes of wheat. First, a coloured de Bruijn graph was built for each chromosome, where each colour represents a different cultivar. Short (1 kb) and evenly distributed (every 20 kb) markers were extracted from each chromosome and compared to the coloured de Bruijn graph to extract their pres-

**Table 2:** TE class proportions in Chinese Spring and Renan genome assemblies

| Parameter | Chinese Spring assembly from Zhu et al. [22] | | Renan RefSeq_v2.0 |
| --- | --- | --- | --- |
| | RefSeq_v1.0 | RefSeq_v2.1 | |
| Genome size (bp) | 14,066,280,851 | 14,225,829,371 | 14,195,643,615 |
| TE (bp) | 11,921,309,743 | 12,092,094,168 | 11,967,447,100 |
| TE (%) | | | |
| All | 84.7 | 85.0 | 84.3 |
| Class I (retrotransposons) | 67.6 | 66.9 | 66.6 |
| Gypsy (RLG) | 46.7 | 46.1 | 45.8 |
| Copia (RLC) | 16.7 | 16.5 | 16.5 |
| Unclassified LTR retrotransposons (RLX) | 3.24 | 3.3 | 3.2 |
| LINE (RIX) | 0.9 | 1.1 | 1.1 |
| SINE (SIX) | 0.01 | 0.01 | 0.01 |
| Class II (DNA transposons) Subclass 1 | 16.5 | 17.0 | 16.9 |
| CACTA (DTC) | 15.5 | 15.9 | 15.8 |
| Mutator (DTM) | 0.38 | 0.44 | 0.44 |
| Unclassified DNA transposons with TIR (DTX) | 0.21 | 0.24 | 0.24 |
| Harbinger (DTH) | 0.16 | 0.18 | 0.18 |
| Mariner (DTT) | 0.16 | 0.17 | 0.17 |
| Unclassified DNA transposons (DXX) | 0.06 | 0.06 | 0.06 |
| hAT (DTA) | 0.006 | 0.009 | 0.009 |
| Helitrons (DHH) | 0.004 | 0.01 | 0.01 |
| Unclassified TE (XXX) | 0.68 | 0.95 | 0.82 |

LINE: long interspersed nuclear element; LTR: long terminal repeat; SINE: short interspersed nuclear element; TIR: terminal inverted repeat.



**Figure 2:** Comparison of existing hexaploid genome assemblies. **A.** contig N50 values in Mb. **B.** Proportion of complete BUSCO genes found in each assembly (N = 4,896). **C.** Number of gaps in each chromosome. **D.** chromosome length in Mb.

ence/absence in each wheat cultivar. On each chromosome, the 15 most abundant presence/absence profiles were selected and used to characterize haplotypic blocks. The haplotype blocks of chromosome 6A, which is associated with productivity traits (e.g., yield, grain size, and height), have already been expertized using a different method [26]. We obtained similar results (Fig. 4), except for the CS chromosome 6A. Previous results have assigned a unique haplotype to this wheat line. But in our case CS ex-

hibits the same haplotype as SY Mattis, Jagger, LongReach Lancer, and Norin61, which had previously been described as sharing the same haplotype. These differences may be explained by the stringency of the comparison, which perhaps should be adjusted separately for each chromosome. Concerning the Renan cultivar, the chromosome 6A has haplotype blocks similar to those of the ArinaLrFor line. Additionally, we used this method to investigate haplotypic blocks that are specific to 1 or a subset of wheat cultivars.

**Figure 3:** Dot plot comparisons of the 21 chromosomes of Renan (y axis) with the Chinese Spring RefSeq v2.1 assembly (x axis).

## Identification of introgressions

Introgression is an important source of genetic variation that is generally the signature of breeding programmes, especially in wheat [32]. Several introgressions have already been reported [8], notably in chromosomes 2B and 3D in LongReach Lancer and in chromosome 2A in Jagger, Mace, SY Mattis, and CDC Stanley. Using our approach, we were able to clearly identify the 2 introgressions in LongReach Lancer (Fig. 5A) and the *A. ventricosa* introgression in chromosome 2A (Fig. 5B). In addition, we found that this introgression of *A. ventricosa* is also present in the Renan cultivar

**Figure 4:** Representation of haplotype blocks in chromosome 6A for the 11 chromosome-scale cultivars (based on 1-Mb blocks). Regions with the same colour represent common regions in wheat lines, except white regions, which are not contained in haplotype blocks. The grey and black regions represent haplotypes respectively shared by ≥10 cultivars or specific to a given cultivar.

(Fig. 5B). The optical map was aligned with this 34-Mb region of Renan and validated the correct structure of this important region carrying multiple resistance genes (Yr17, Lr37, Sr38, Cre5). More importantly, the 34 Mb consisted of 22 contigs in Renan and 2,339 in Jagger. A comparison of the fragmentation near the introgression point is presented in Fig. 5D and shows a large difference between the long- and short-read assemblies. Additionally, we also identified several candidate introgressions, which had already been spotted through retrotransposon profiles [8]: (i) a 45-Mb region on chromosome 2D that is shared between the lines Julius, ArinaLrFor, SY Mattis, Jagger, and also Renan (Fig. 6A); (ii) a 53-Mb region at the end of chromosome 3D in LongReach Lancer (Fig. 6B); (iii) a 48-Mb region at the beginning of chromosome 3D in SY Mattis (Fig. 6B); and (iv) the *A. ventricosa* introgression of 30 Mb in chromosome 7D, which carries Pch1 resistance gene (Fig. 6C).

Moreover, a known large-scale structural variation in chromosomes 5B and 7B of ArinaLrFor and SY Mattis cultivars was also easily identifiable using haplotypic blocks of individual chromosomes (Supplementary Fig. S4).
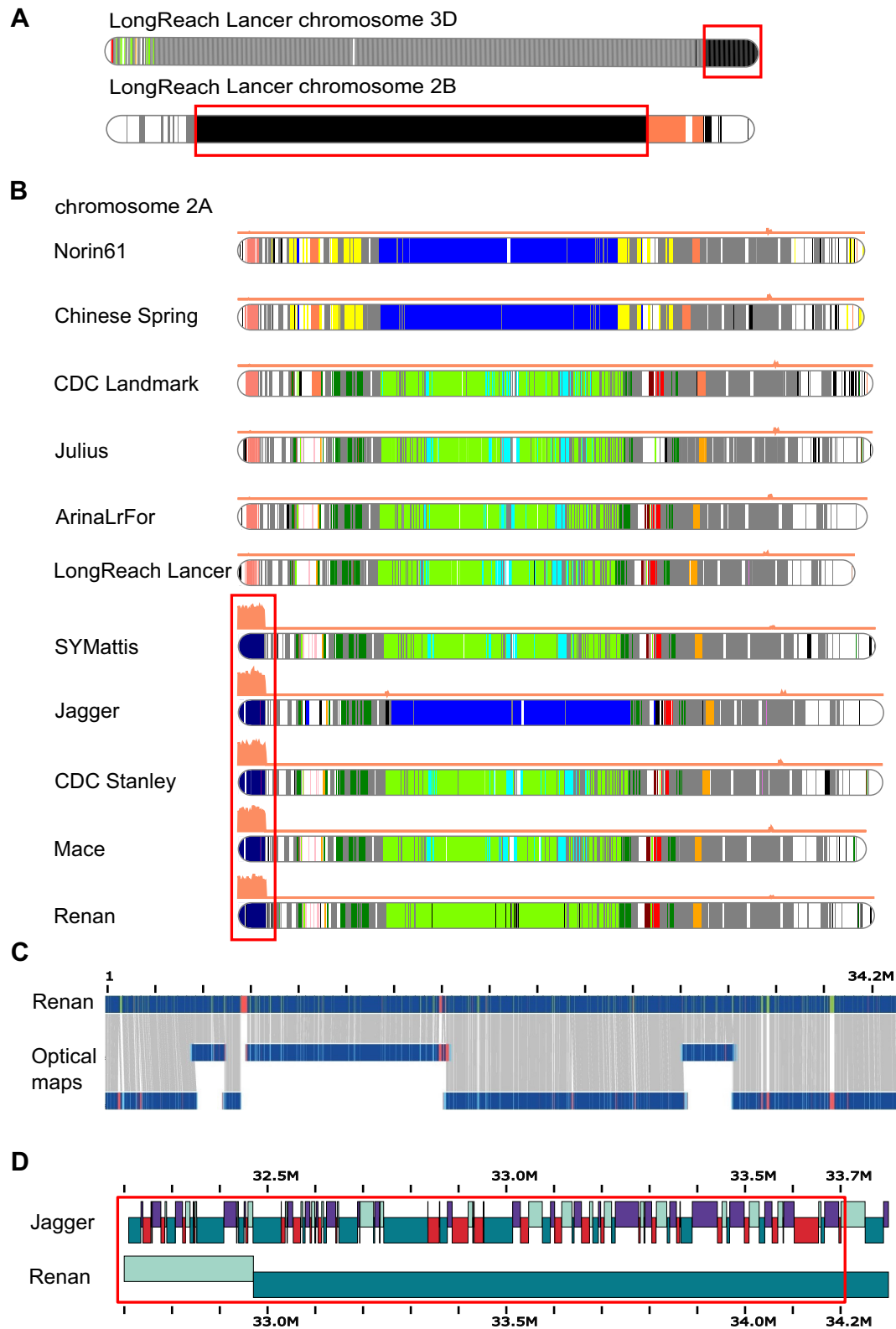
## Comparative analysis of a storage protein–coding gene cluster in *T. aestivum*

Tandem duplications are an important mechanism in plant genome evolution and adaptation [33, 34], but the assembly of tandemly duplicated gene clusters is difficult, especially with short reads. To illustrate the gain brought by this optimized assembly process, we focused on an important locus on chromosome 1B known to carry multiple copies of storage protein and disease resistance genes [35, 36]. Among them, the genes encoding $\omega$-gliadins are not only duplicated in tandem but also composed of microsatellite DNA in their coding part, making them particularly hard to assemble properly from short reads. We compared orthologous regions harbouring these genes between CS and Renan, spanning 1.58 and 2.32 Mb, respectively. The CS region was more fragmented, with 101 gaps versus only 3 in Renan (Fig. 7A). The number of copies of $\omega$-gliadin–encoding genes was quite similar: 9 in CS and 10 in Renan. The most striking difference came from the completeness of the microsatellite motifs: 8 copies out of 9 contain N stretches in CS RefSeq v2.1, revealing that the microsatellite is usually too large to be fully assembled with short reads (Fig. 7B). In contrast, all 10 copies predicted in Renan were assembled completely. More generally, we mapped the $\omega$-gliadin–encoding genes, annotated on CS in a previous study [35], back to the locus and showed that it was better reconstructed in the Renan assembly, with a mean protein alignment length of 99% compared to 58% in CS (Fig. 7C). In addition, the optical map and long reads were used to validate the structure of this region in Renan, which turns out to be consistent (Fig. 7D and Supplementary Fig. S5).

**Figure 5:** Haplotypic blocks in wheat chromosomes. Colours represent common regions in wheat cultivars. The grey and black regions represent haplotypes respectively shared by ≥10 cultivars or specific to a given cultivar. The orange curve, when present, represents coverage with *A. ventricosa* short reads. The red boxes frame the introgressions. **A.** Known introgressions in chromosomes 3D and 2B in LongReach Lancer. Regions in black represent genomic regions that are specific to LongReach Lancer and are respectively *Triticum ponticum* and *Triticum timopheevii* introgressions as described previously [8]. **B.** *A. ventricosa* introgression on chromosome 3D in CDC Stanley, Mace, SY Mattis, and Jagger. This known introgression is also present in Renan. The dark blue block represents the region shared across the 5 cultivars. **C.** Validation of the introgression in Renan (chromosome 2A from 1 to 34.2 Mb) using Bionano maps. **D.** Comparison of the contig composition of the first megabases from the introgression point in Jagger and Renan cultivars.
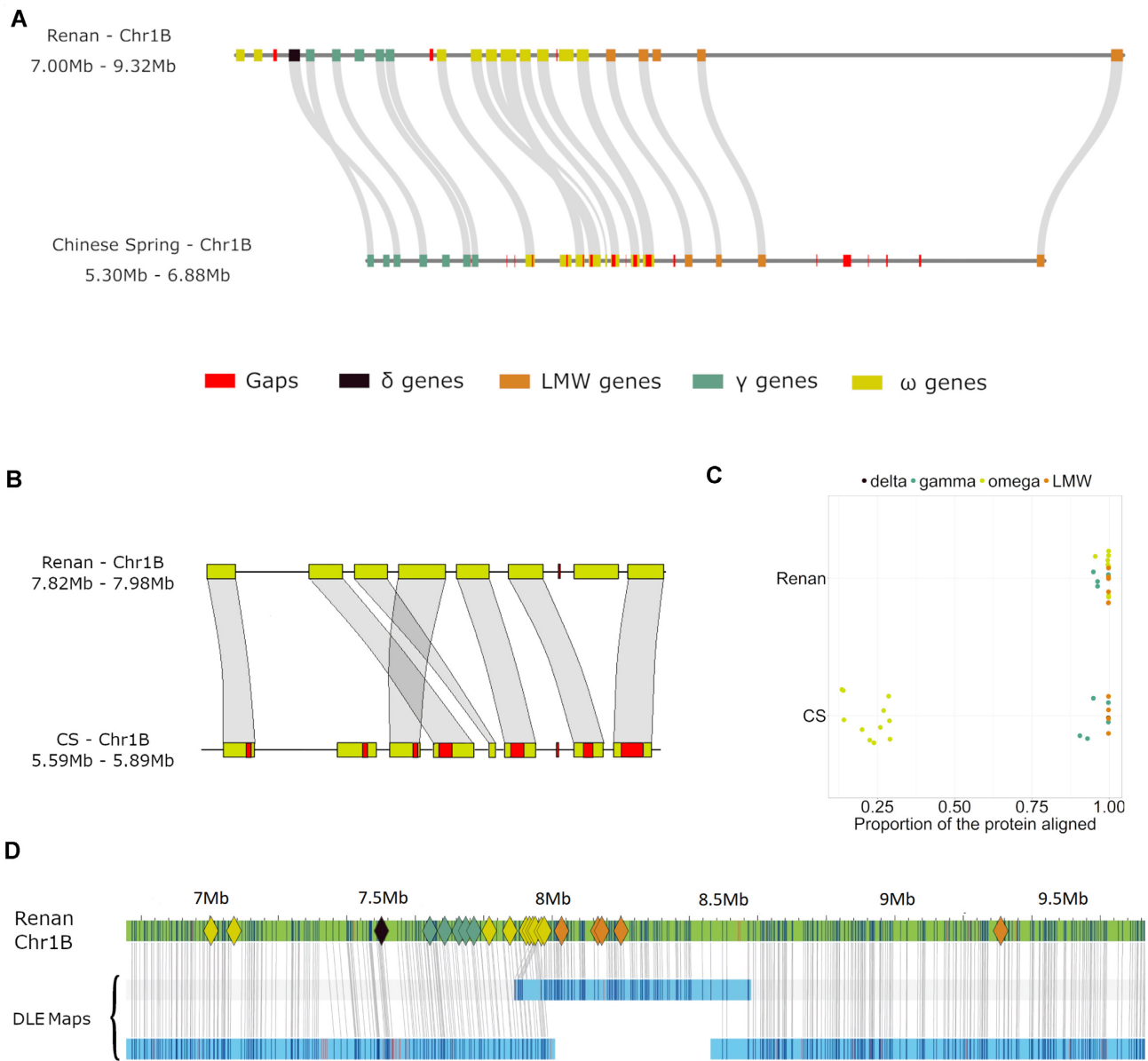
**Figure 6:** Haplotypic blocks in wheat chromosomes. Colours represent common regions in wheat cultivars. The grey and black regions represent haplotypes respectively shared by ≥10 cultivars or specific to a given cultivar. The orange curve represents coverage with *A. ventricosa* short reads. Non-zero coverage of the D subgenome is expected because this subgenome is evolutionarily close to the Dv genome of *A. ventricosa*. The red boxes frame the introgressions. **A.** Candidate introgression (green block) on chromosomes 2D in Julius, ArinaLrFor, SY Mattis, Jagger, and Renan. **B.** Candidate introgressions (black blocks) on chromosome 3D in LongReach Lancer and SY Mattis. **C.** *A. ventricosa* introgression (black block) on chromosome 7D in Renan.

## Comparative analysis of the locus that provides resistance to the orange wheat blossom midge

Like a few other wheat cultivars, Renan is resistant to the orange wheat blossom midge (OWBM). The *Sm1* gene is known to confer resistance to wheat, and a previous study has shown that CDC Landmark is also resistant to the OWBM and carries a 7.3-Mb haplotype within the *Sm1* locus on chromosome 2B [8]. We extracted and aligned the corresponding region of CDC Landmark on each cultivar to precisely locate the corresponding region on each chromosome 2B. From these 11 regions of 1–2 Mb, we computed the haplotypic blocks using a higher resolution than previously (1 kb marker every 5 kb). This analysis revealed a strong similarity of the *Sm1* locus between CDC Landmark and Renan (Fig. 8A), the presence of the *Sm1* gene in blocks shared between the 2 cultivars. In addition, a comparison of the fragmentation of these 2 regions underlines the higher contiguity of the Renan assembly, with 4 contigs in the Renan *Sm1* locus compared to 62 in CDC Landmark (Fig. 8B). The *Sm1* locus of Renan is in agreement with the optical map and shows clearly the 3 remaining gaps that may correspond to smaller and unanchored contigs.

## Discussion

In this study, we showed that the recent improvement of the ONT, in terms of error rate and throughput, has opened up new perspectives in the age of long-read technologies. Indeed, the sequencing and assembly of complex genomes, like hexaploid wheat, is now accessible to sequencing facilities. Additionally, the ability to sequence ultra-long reads using ONT devices is a real advantage over the other long-read technology, namely, PacBio. In this study, we were able to generate a coverage of 14× with reads longer than 50 kb, whereas PacBio libraries, used to generate HiFi (High-Fidelity) reads, are generally sized ~15 kb [37, 38]. Several studies have already underlined the positive impact of these ONT ultra-long reads on the assembly contiguity [9, 37, 39]. In contrast, the error rate that was previously a thorn in their side has been drastically reduced over the past year. Herein we reported a quality score near Q10–Q15 for individual ONT reads, as already shown [28], which is still far from what HiFi reads can provide, generally near Q30 [37]. The high accuracy of HiFi reads might be sufficient to distinguish copies from repeat regions if they present few variations. The impact of ultra-long reads will lie mainly in the case
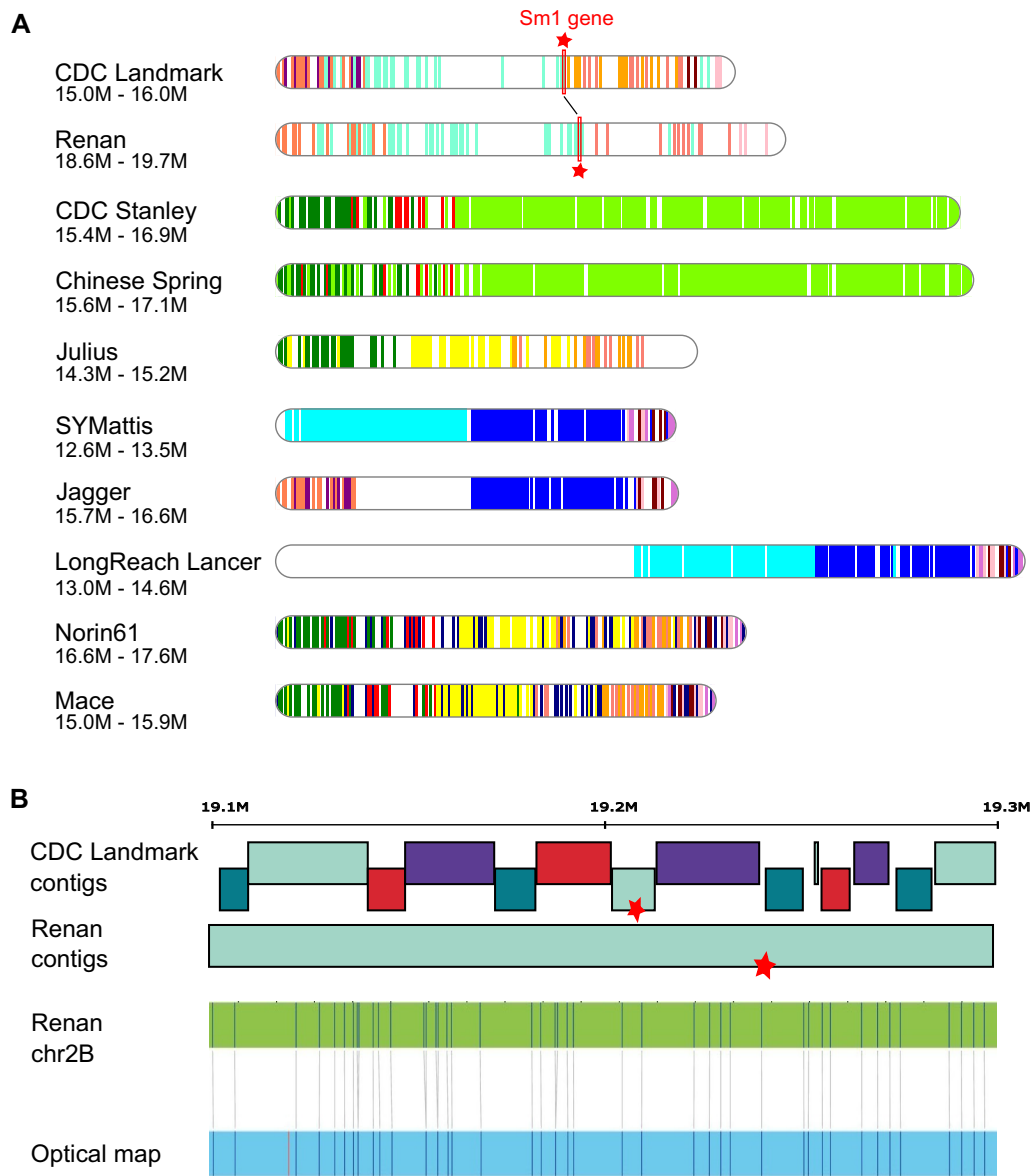
**Figure 7:** Comparative view of an important locus on chromosome 1B containing prolamin and resistance genes, tandemly duplicated. **A.** Representation of the region with gaps and genes on the 2 assemblies of Renan and CS. **B.** Zoomed view on the ω-gliadin gene cluster. **C.** Proportion of the length of the proteins that were aligned in the genomic region of Renan and CS. Aligned protein sequences were annotated in CS by Huo et al. [35]. **D.** Alignment view of Bionano maps on the Renan cluster; coloured diamond shapes represent genes belonging to the ω-gliadin gene cluster. The optical maps are in blue and the chromosome sequence in green. Restriction sites are represented by vertical lines and are joined between the sequence and the map when properly aligned.

of identical repeats, and obviously, the presence of these particular cases will depend on the evolutionary history of the studied genomes. In addition, this high error rate has an impact on the consensus quality, and at the moment, a combination of ONT and Illumina reads is still needed to achieve a decent per-base accuracy.

By following basecaller evolution, we noticed that the gain when using a recent basecaller is high and we guess that this observation will encourage users to reprocess older data. However, this is not trivial and it requires sufficient computing resources. Interestingly, we observed that the error rate of ONT data is organism dependent and that basecaller training has a significant impact on the overall quality of the reads [28]. This is, in our opinion, an important fact because a large proportion of *de novo* as-

semblies now concern non-model organisms and users will have to address this limitation of current software. There are existing methods to train the basecaller on non-model species [40, 41], but this can still be a big barrier, depending on the size of the dataset, for many end users. However, as highlighted in this study, the combination of long- and short-read sequencing with polishing methods greatly improves the consensus sequence of a given genome assembly, and these algorithms seem sufficient at least in coding regions.

Even though there are now several chromosome-scale assemblies of the hexaploid wheat genome, this assembly of the Renan variety based on long reads will benefit biologists and geneticists because it offers a high resolution. We show that our chromosome-scale assembly of Renan based on long reads can

**Figure 8:** Comparison of the *Sm1* loci. **A.** Representation of haplotype blocks (5-kb bins) of the region surrounding the *Sm1* gene on chromosome 2B. Colours represent common regions in wheat cultivars. The genomic region of CDC Landmark (15–16 Mb) was aligned against other cultivars to localize the *Sm1* loci. The *Sm1* gene in CDC Landmark and Renan, the 2 *Sm1* carrier cultivars, is represented by a red star. **B.** Comparison of the contig composition in the *Sm1* region of CDC Landmark and Renan, and validation of the assembly structure in Renan using Bionano optical maps. The optical map is in blue and the chromosome sequence in green. Restriction sites are represented by vertical lines and are joined between the sequence and the map when properly aligned.

bring new insight into genomic regions of interest. In particular, this can happen in regions that carry multiple resistance genes, such as a large *A. ventricosa* introgression shared with other cultivars on chromosome 2A and a unique *A. ventricosa* introgression on chromosome 7D. The lower number of gaps in these regions will help to localize genes of interest and to provide a better understanding of the impact of these introgressions. Additionally, we demonstrated, by examining 2 important loci containing prolamin and resistance genes, that such regions are truly enhanced and contain very few gaps compared to assemblies based on short reads.

Moreover, unlike recent chromosome-scale assemblies, Renan's gene prediction is not only a projection of CS gene models but also includes *de novo* annotation with RNA-Seq data, which is of real benefit for the construction of pan genome (or pan annota-

tion) or when cultivar-specific genes are examined. For all of these reasons, we believe that this high-resolution assembly will benefit the wheat community and help breeding programs dedicated to the bread wheat genome.

## Methods

### Plant material and DNA extraction

*Triticum aestivum* cv. Renan seeds were provided by the INRAE Biological Resource Center on small-grain cereals and grown for 2 weeks, after which a dark treatment was applied on the seedlings for 2 days before collecting leaf tissue samples.

For the sequencing experiments, DNA was isolated from frozen leaves using Qiagen Genomic-tips 100/G kit (Cat No./ID: 10243) and following the tissue extraction protocol. Briefly, 1 g of leaves

were ground in liquid nitrogen with mortar and pestle. After 3 hours of lysis and 1 centrifugation step, the DNA was immobilized on the column. After several washing steps, DNA was eluted from the column, then desalted and concentrated by alcohol precipitation. The DNA was then resuspended in the TE buffer.

To generate the optical map, ultra-HMW DNA was purified from 0.5 g of very young fresh leaves according to the Bionano Prep Plant tissue DNA Isolation Base Protocol (30068–Bionano Genomics, San Diego, CA, USA) with the following specifications and modifications. Briefly, the leaves were fixed using a fixing solution (BNG) containing formaldehyde (Sigma-Aldrich, Saint-Louis, MO, USA) and then ground in a homogenization buffer (BNG) using a Tissue Ruptor grinder (Qiagen, MD, USA). Nuclei were washed and embedded in agarose plugs. After overnight proteinase K digestion in Lysis Buffer (BNG) and 1-hour treatment with RNAse A (Qiagen, MD, USA), plugs were washed 4 times in 1× Wash Buffer (BNG) and 5 times in 1× TE Buffer (ThermoFisher Scientific, Waltham, MA). Then, plugs were melted for 2 minutes at 70°C and solubilized with 2 μL of 0.5 U/μL AGARase enzyme (ThermoFisher Scientific, Waltham, MA) for 45 minutes at 43°C. A dialysis step was performed in 1× TE Buffer (ThermoFisher Scientific, Waltham, MA) for 45 minutes to purify DNA from any residues. The DNA samples were quantified by using the Qubit dsDNA BR Assay (Invitrogen, Carlsbad, CA, USA). Quality of megabase-size DNA was validated by pulsed-field gel electrophoresis.

## Illumina Sequencing

DNA (1.5 μg) was sonicated using a Covaris E220 sonicator (Covaris, Woburn, MA, USA). Fragments (1 μg) were end-repaired, 3′-adenylated, and Illumina adapters (Bioo Scientific, Austin, TX, USA) were then added using the Kapa Hyper Prep Kit (Kapa Biosystems, Wilmington, MA, USA). Ligation products were purified with AMPure XP beads (Beckman Coulter Genomics, Danvers, MA, USA). Libraries were then quantified by qPCR using the KAPA Library Quantification Kit for Illumina Libraries (Kapa Biosystems, Wilmington, MA, USA), and library profiles were assessed using a DNA High Sensitivity LabChip kit on an Agilent Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA). The library was sequenced on an Illumina NovaSeq instrument (Illumina, San Diego, CA, USA) using 150 base-length read chemistry in a paired-end mode. After the Illumina sequencing, an in-house quality control process was applied to the reads that passed the Illumina quality filters [42]. These trimming and removal steps were achieved using Fastxtend tools [43].

## Nanopore Sequencing

Libraries were prepared according to the protocol Genomic DNA by ligation (SQK-LSK109 kit). Genomic DNA fragments (1.5 μg) were repaired and 3′-adenylated with the NEBNext FFPE DNA Repair Mix and the NEBNext® Ultra™ II End Repair/dA-Tailing Module (New England Biolabs, Ipswich, MA, USA). Sequencing adapters provided by ONT (Oxford Nanopore Technologies Ltd, Oxford, UK) were then ligated using the NEBNext Quick Ligation Module (NEB). After purification with AMPure XP beads (Beckmann Coulter, Brea, CA, USA), the library was mixed with the Sequencing Buffer (ONT) and the Loading Bead (ONT) and loaded on MinION or PromethION R9.4.1 flow cells. One PromethION run was performed with Genomic DNA purified with Short Read Eliminator kit (Circulomics, Baltimore, MD, USA) before the library preparation.

## Optical maps

Labelling and staining of the ultra-HMW DNA were performed according to the Bionano Prep Direct Label and Stain (DLS) protocol (30206–Bionano Genomics, San Diego, CA, USA). Briefly, labelling was performed by incubating 750 ng genomic DNA with 1× DLE-1 Enzyme (BNG) for 2 hours in the presence of 1× DL-Green (BNG) and 1× DLE-1 Buffer (BNG). Following proteinase K digestion and DL-Green clean-up, the DNA backbone was stained by mixing the labelled DNA with DNA Stain solution (BNG) in the presence of 1× Flow Buffer (BNG) and 1× DTT (BNG), and incubating overnight at room temperature. The DLS DNA concentration was measured with the Qubit dsDNA HS Assay (Invitrogen, Carlsbad, CA, USA).

Labelled and stained DNA was loaded on Saphyr chips. Loading of the chips and running of the BNG Saphyr System were all performed according to the Saphyr System User Guide (30247–Bionano Genomics, San Diego, CA, USA). Data processing was performed using the BNG Access software.

A total of 4,541 Gb data were generated. From these data, molecules of size >150 kb were filtered, generating 1,931 Gb of data. These filtered data, corresponding to 128× coverage of the *Triticum aestivum* cv. Renan, consist of 7,810,298 molecules with an N50 of 237.5 kb and an average label density of 14.3/100 kb. The filtered molecules were aligned using RefAligner with default parameters. It produced 1,053 genome maps with an N50 of 37.5 Mb, for a total genome map length of 14,946.8 Mb.

## RNA extraction

Samples of several tissues (stem, leaf, root, or grain) were collected on plants with different growth conditions and of different ages. Each of these 28 tissues was subjected to RNA extraction with the following protocol: 200 mg to 1 g of fine powder was put in a 50-mL falcon tube with 4.5 mL of NTES buffer (0.1 M NaCl, 1% SDS, 10 mM Tris-HCl [pH 7.4], 1 mM EDTA [pH 8]). After vortexing the tube, 3 mL of phenol-chloroforme-IAA were added. The tube was mixed for 10 minutes and centrifuged for 20 minutes at 5,000 rpm (15°C). The aqueous phase was collected and placed in a new 15-mL tube. Then 3 mL of phenol-chloroforme-IAA were added. The tube was mixed for 10 minutes and centrifuged for 20 minutes at 5,000 rpm (15°C). The aqueous phase was collected and placed in a new 50-mL tube. Then 1/10 of AcNa 3M (pH 5.2) and 2 volumes of 100% ethanol were added. The tube was mixed gently by turning and centrifuged for 20 minutes at 5,000 rpm (4°C). The supernatant was removed. The precipitate was dried and resuspended in 20 μL RNAse-free water. A treatment with DNase was realized and the RNA were purified on a MinElute column (Qiagen, MD, USA). A second treatment with DNAse was realized by adding DNAse directly on the filter. After ethanol clean-up, the column was eluted with 14 μL of RNAse-free water. The quality of the RNA was evaluated using RNA 6000 Nano Assay chip for size and RIN estimation and spectrophotometry (A260/A280 and A260/A230 ratios) for purity estimation. The RNA were quantified using Qubit RNA high sensitivity Assay kit (Invitrogen, Carlsbad, CA, USA).

## RNA sequencing

RNA-Seq library preparations were carried out from 500–2,000 ng of total RNA using the TruSeq Stranded messenger RNA (mRNA) kit (Illumina, San Diego, CA, USA), which allows mRNA strand orientation (sequence reads occur in the same orientation as antisense RNA). Briefly, polyadenylated RNA was selected with oligo(dT) beads, chemically fragmented, and converted into single-stranded cDNA using random hexamer priming. Then, the

second strand was generated to create double-stranded complementary DNA (cDNA). The cDNA was then 3′-adenylated, and Illumina adapters were added. Ligation products were PCR-amplified. Ready-to-sequence Illumina libraries were then quantified by qPCR using the KAPA Library Quantification Kit for Illumina Libraries (Kapa Biosystems, Wilmington, MA, USA), and library profiles were evaluated with an Agilent 2100 Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA). Each library was sequenced using 151-bp paired-end read chemistry on an Illumina NovaSeq 6000 sequencer (Illumina, San Diego, CA, USA).

### Long-read genome assembly

The 20 ONT runs were basecalled using 2 versions of guppy: 3.3 HAC and 3.6 HAC (Supplementary Table S6). We monitored the gain of each guppy basecaller release and evaluated 3 different assemblers in the context of large genomes: Redbean [16] v2.5 (git commit 3d51d7e), SMARTdenovo (SMARTdenovo, RRID:SCR_017 622) [17] (git commit 5cc1356), and Flye (Flye, RRID:SCR_017016) [18] v2.7 (git commit 5c12b69). All assemblers were launched using a subset of reads consisting of 30× of the longest reads (Supplementary Table S3). Then, we selected 1 of the assemblies based not only on contiguity metrics such as N50 but also cumulative size and proportion of unknown bases. The Flye (longest reads) and SMARTdenovo (all reads) assemblies were very similar in terms of contiguity, but we decided to keep the SMARTdenovo assembly because its cumulative size was higher. The SMARTdenovo assembler using the longest reads resulted in a contig N50 of 1.1 Mb and a cumulative size of 14.07 Gb. Because nanopore reads contain systematic error in homopolymeric regions, we polished the consensus of the selected assembly with nanopore reads as input to the Racon (v1.3.2, git commit 5e2ecb7) and Medaka software packages. In addition, we polished the assembly 2 additional times using Illumina reads as input to the Hapo-G tool (v1.0).

### Long-range genome assembly

The BNG scaffolding workflow (Bionano Solve version 3.5.1) was launched with the nanopore contigs and the Bionano map. We found in several cases that the nanopore contigs were overlapping (based on the optical map), and these overlaps were corrected using the BiSCoT software [44] with default parameters. Finally, the consensus sequence was polished once more using Hapo-G and short reads, to ensure correction of duplicate regions that were collapsed (Supplementary Table S4).

### Validation of the *Triticum aestivum* cv Renan assembly

The quality value (QV) of the Renan and CS assemblies was obtained using Merqury [45]. First, 31-mers were extracted from the Renan and CS Illumina sequencing reads (accessions SRR5893651, SRR5893652, SRR5893653, and SRR5893654), and then the QV of each genome assembly was computed using Merqury (version 1.3, git commit 6b5405e).

We used BLAST [25] to search for the presence of 107,891 high-confidence genes from CS RefSeq v1.1 in the Renan genome sequence. We extracted the 461,476 individual exons larger than 30 bp and without Ns from this dataset and computed exon-by-exon BLAST to avoid spurious sliced alignments. An exon was considered present if it matched the Renan scaffolds with ≥95% identity over ≥95% of its length. To estimate the proportion of identical exons between CS and Renan and the average nucleotide identity, we used the same BLAST-based procedure but while restricting the dataset to 454,008 CS exons that are on pseudomolecules (ex-

cluding chrUn) and considering Renan pseudomolecules instead of scaffolds; i.e., only exons carried by the same chromosome in CS and Renan were considered. We extracted all available ISBPs (150 bp each) from the CS RefSeq v1.1 and filtered out ISBPs containing Ns and those that do not map uniquely on the CS genome. This led to the design of a dataset containing 5,394,172 ISBPs that were aligned on the Renan scaffolds using BLAST. We considered an ISBP to be conserved in Renan if it matched with ≥90% identity over 90% of its length. We used the same ISBP dataset to study the effect of polishing on error rate in the assembly while using BLAST and considering ≥90% identity over ≥145 aligned nucleotides.

### Anchoring of the *Triticum aestivum* cv Renan assembly

We guided the construction of 21 Renan pseudomolecules based on collinearity with the CS RefSeq Assembly v2.1. For this, we used the positions of conserved ISBPs as anchors (5,087,711 ISBPs matching with ≥80% identity over ≥90% query overlap). This represented 357 ISBPs/Mb, meaning that even the smallest scaffolds (30 kb) generally carried >10 potential anchors. However, some ISBPs match at non-orthologous positions, which create noise to precisely determine the order and orientation of some scaffolds. To overcome this issue, we considered ISBPs by pairs. Only pairs of adjacent ISBPs (i.e., separated by <50 kb on both CS and Renan genomes) were kept as valid anchors, allowing isolated mismapped ISBPs to be filtered out. Only scaffolds harbouring ≥50% of valid ISBP pairs on a single chromosome were kept. The others were considered unanchored and they comprised the "chrUn." We calculated the median position of matching ISBP pairs along each CS chromosome for defining the order of the Renan scaffolds relative to each other. Their orientation was retrieved from the orientation of all matching ISBP pairs in CS following the majority rule. We thus built 21 pseudomolecules that were then corrected according to the Hi-C map as explained hereafter.

Two Hi-C biological replicates were prepared from 10-days plantlets of *T. aestivum* cv. Renan following the AG Hi-C protocol (AG Hi-C User Guide for Plant Tissues DOC A160106 v01). For each replicate, 2 libraries were constructed using the Kapa Hyper Prep kit (Roche) according to AG's recommendation (Library Preparation using KAPA Hyper Prep Kit DOC A160108 v01). The technical replicates were then pooled and sent to Genewiz for sequencing on an Illumina HiSeq4000 (4 lanes in total), reaching a 35× coverage. We mapped a sample of 240 million read pairs with BWA-MEM [46]) to the formerly built 21 pseudomolecules, filtered out for low quality, sorted, and deduplicated using the Juicer pipeline (Juicer, RRID:SCR_017226) [47]. We produced a Hi-C map from the Juicer output by the candidate assembly visualizer mode of 3D-DNA pipeline [48] and visualized it with the Juicebox Assembly Tools software. Based on abnormal frequency contacts' signals revealing a lack of contiguity, scaffold-level modifications of order, orientation, and/or chimeric scaffolds were identified to improve the assembly. In case of chimeric scaffolds, coordinates of resulting fragments were retrieved from the Juicebox Assembly Tools application but then recalculated to correspond precisely to the closest gap in the scaffold. Pseudomolecules were eventually rebuilt from initial scaffolds and new fragments while adding 100N gaps between neighbour scaffolds. A final Hi-C map was built to validate the accuracy of the final assembly.

### Calculation of chromosome coverage

Short (*T. aestivum* cv Renan and *A. ventricosa*) and long reads (*T. aestivum* cv Renan) were aligned using minimap2 (with the following

parameters: "-I 17G -2 –sam-hit-only -a -x sr" and "-I 17G -2 –sam-hit-only –secondary = no -a -x map-ont", respectively). Coverage of individual chromosomes was calculated in 1-Mb windows using mosdepth [49] (version 0.3.1) and the following parameters: "–by 1 000 000 -n -i 2 -Q 10 -m." Note that the "-i 2" and "-Q 10" parameters were used to keep only alignments of reads that mapped in a proper pair and with a minimal quality value of 10. Coverage of individual chromosomes is plotted in Fig. 1. In addition, large deletions and duplications were detected using CNVnator (CNVnator, RRID:SCR_010821) [50] with the Illumina bam file and a window of 100 bp. We focused on large events (>500 kb) and detected only 15 deletions and no duplication (Fig. 1).

## Transposable element annotation

TEs were annotated using CLARITE [29]. Briefly, TEs were identified through a similarity search approach based on the ClariTeRep curated databank of repeated elements using RepeatMasker (RepeatMasker, RRID:SCR_012954) and modelled with the CLARITE program that was developed to resolve overlapping predictions, merge adjacent fragments into a single element when necessary, and identify patterns of nested insertions [29].

## Gene prediction

We used the MAGATT pipeline [51] to map the full set of 106,801 High Confidence and 159,848 Low Confidence genes predicted in Chinese Spring IWGSC RefSeq v2.1. The workflow implemented in this pipeline was described in Zhu et al. [22]. Briefly, it uses gene-flanking ISBP markers to determine an interval that is predicted to contain the gene before homology-based annotation transfer, limiting problems due to multiple mapping. When the interval is identified, MAGATT uses BLAT [52] to align the gene (untranslated regions, exons, and introns) sequence and recalculate all sub-feature coordinates if the alignment is full-length and without indels. If the alignment is partial or contains indels, it runs GMAP [53] to perform spliced alignment of the candidate CDS inside the interval. If no ISBP-flanked interval was determined or if both BLAT and GMAP failed to transfer the gene, MAGATT runs GMAP against the whole genome, including the unanchored fraction of the Renan assembly. We kept the best hit considering a minimum identity of 70% and a minimum coverage of 70%, with "cross_species" parameter enabled.

We then masked the genome sequence based on mapped genes and predicted TE coordinates using BEDTools (BEDTools, RRID:SCR_006646) thfthy [54], mergeBed, and maskfasta v2.27.1. Hence, we computed a *de novo* gene prediction on the unannotated part of the genome. We used TriAnnot [31] to call genes based on a combination of evidence: RNA-Seq data, *de novo* predictions of gene finders (FGeneSH, Augustus), and similarity with known proteins in Poaceae, as described previously [7]. For that purpose, we mapped RNA-Seq reads with hisat2 [55] v2.0.5, called 277,505 transcripts with StringTie [56] v2.0.3, extracted their sequences with Cufflink (Cufflinks, RRID:SCR_014597) [57] gffread v2.2.1, and provided this resource as input to TriAnnot. We optimized TriAnnot workflow to ensure a flawless use on a cloud-based high-performance compute cluster (10 nodes with 32 CPUs/128 GB RAM each and shared file system) using the IaaS Openstack infrastructure from the UCA Mesocentre. Gene models were then filtered as follows: we discarded gene models that shared strong identity (≥92% identity, ≥95% query coverage) with an unannotated region of the Chinese Spring RefSeq v2.1, considered as doubtful predictions. We then kept all predictions that matched RNASeq-derived transcripts (≥99% identity, ≥70% query and subject coverage). For

those that did not show evidence of transcription, we kept gene models sharing protein similarity (≥40% identity, ≥50% query and subject coverage) with a Poaceae protein having a putative function (filtering out based on terms "unknown," "uncharacterized," and "predicted protein").

## Comparison of genome assemblies

Genome assemblies were downloaded from https://webblast.ipk-gatersleben.de/downloads. Contigs were extracted by splitting input sequences at each N and standard metrics were computed. Gene completion metrics were calculated using BUSCO v5.0 and version 10 of the poales geneset, which contains 4,896 genes.

We built dot plots between Renan, CS, and 10 other reference quality genomes (ArinaLrFor, CDC Landmark, CDC Stanley, Jagger, Julius, LongReach Lancer, Mace, Norin61, SY Mattis, spelta PI190962) by using orthologous positions of conserved ISBPs (1 ISBP every 2.5 kb on average) identified by mapping them with BWA-MEM (maximum 2 mismatches, 100% coverage, and minimal mapping quality of 30).

## Characterization of haplotypic blocks

First a coloured de Bruijn graph was built for each chromosome from the 11 available chromosome-scale assemblies of wheat (Renan, CS, ArinaLrFor, CDC Landmark, CDC Stanley, Jagger, Julius, LongReach Lancer, Mace, Norin61, and SY Mattis). The coloured de Bruijn graph was created using Bifrost [58] with 31-mers and a unique colour for each wheat cultivar. In a second step, we extracted short markers (1 kb) evenly spaced (20 or 5 kb) on each chromosome and queried the coloured de Bruijn graph using Bifrost and the

parameter "-e 0.95" (for the comparison of each chromosome) and "-e 0.97" (for the comparison of the *Sm1* locus). This parameter is the ratio of *k*-mers from queries that must occur in the graph to be reported as present. For whole-chromosome analyses, the 20-kb blocks were merged into 1-Mb blocks (the most abundant colour in the 50 20-kb blocks was retained for the 1-Mb block). Individual blocks and *A. ventricosa* coverage were displayed using RIdeogram [59].

## Comparison of a storage protein–coding gene cluster

We performed manual curation of the gene models encoding storage proteins predicted in Renan. Protein sequences of prolamin and resistance genes [35], annotated from a 1B chromosome locus in the PacBio-based assembly of Chinese Spring [5], were downloaded and aligned to the IWGSC RefSeq v2.1 and Renan assemblies using BLAT [52] with default parameters. Draft alignments were refined by aligning the given protein sequence and the genomic region defined by the blat alignment using GeneWise (GeneWise, RRID:SCR_015054) with default parameters. Resulting alignments were filtered to conserve only the best match for each position by keeping only the highest-scoring alignment, and the genomic region containing the gene cluster was extracted. Then, we used the jcvi suite [60] with the mcscan pipeline to find synteny blocks between both genomes. First, we used the "jcvi.compara.catalog" command to find orthologs and then the "jcvi.compara.synteny mcscan" with "–iter = 1" command to extract synteny blocks. Finally, we generated the figure with the "jcvi.graphics.synteny" command and manually edited the generated svg file to improve the resulting image by changing gene colours, incorporating gaps, and renaming genes. Moreover, to make the figure clearer, we artificially reduced the intergenic

space by 95% so that gene structures appear bigger. The $\omega$ gene cluster representation figure was generated by using DnaFeaturesViewer [61] with coordinates of features generated by the mcscan pipeline used previously.

## Data Availability

The Illumina and PromethION sequencing data and the Bionano optical map are available in the European Nucleotide Archive under the project PRJEB49351. The genome assembly and gene predictions are freely available from the Genoscope website [62]. Additionally, all data and scripts used to produce the main figures as well as the haplotypic blocks and graphical visualizations are available on a GitHub repository [63]. Supporting data are also available via the GigaScience database GigaDB [64].

## Additional Files

Supplementary Table S1. Statistics of the ONT sequencing data.

Supplementary Table S2. BNG Optical map.

Supplementary Table S3. Raw long-read assemblies.

Supplementary Table S4. Hybrid assemblies obtained using ONT and BNG data.

Supplementary Table S5. Impact of the polishing .

Supplementary Table S6. Statistics of the ONT reads obtained with two different versions of the guppy basecaller .

Supplementary Table S7. Long-read assemblies with ONT reads obtained with two different versions of the guppy basecaller.

Supplementary Table S8. RNASeq data for Renan from 28 samples corresponding to 14 different organs/conditions in replicates.

Supplementary Figure S1. Validation of both ends of the chromosome 7A .

Supplementary Figure S2. Comparison of the accuracy of different ONT basecallers.

Supplementary Figure S3. Number of gaps per Mbp in Chinese Spring and Renan genome assemblies.

Supplementary Figure S4. Large-scale structural variation in chromosomes 5B and 7B.

Supplementary Figure S5. Validation of the omega gliadin gene cluster.

Supplementary File S1. Dot plots of the 21 chromosomes of Renan with other wheat genome assemblies.

## Abbreviations

AG: Arima Genomics; BLAST: Basic Local Alignment Search Tool; BNG: BioNano Genomics; bp: base pairs; BUSCO: Benchmarking Universal Single-Copy Orthologs; BWA: Burrows-Wheeler Aligner; CPU: central processing unit; CS: Chinese Spring; Gb: gigabase pairs; ISBP: insertion site–based polymorphism; IWGSC: International Wheat Genome Sequencing Consortium; kb: kilobase pairs; MAGATT: Marker Assisted Gene Annotation Transfer for Triticeae; Mb: megabase pairs; mRNA: messenger RNA; ONT: Oxford Nanopore Technologies; OWBM: orange wheat blossom midge; PacBio: Pacific Biosciences; QV: quality value; RAM: random access memory; TE: transposable element.

## Competing Interests

J.M.A. received travel and accommodation expenses to speak at Oxford Nanopore Technologies conferences. J.M.A. and C.B. received accommodation expenses to speak at Bionano Genomics user meetings. The authors declare that they have no other competing interests.

## Authors' Contributions

S.A., I.D., and A.B. extracted the sequenced DNA and generated the optical map. C. C. and A.A. optimized and performed the nanopore and Illumina sequencing. N.P., E.P., and M.R. generated the Hi-C libraries and sequences. J.M.A., S.E., B.I., C.M., P.L.Z., C.B., H.R., P.L., D.G., and F.C. performed the bioinformatic analyses. J.M.A., S.E., B.I., C.M., P.L.Z., C.B., C.C., H.R., P.L., and F.C. wrote the manuscript. J.M.A., P.W., and F.C. supervised the study.

## References

1. Dubcovsky, J, Dvorak, J. Genome plasticity a key factor in the success of polyploid wheat under domestication. *Science* 2007;**316**(5833):1862–6.
2. Marcussen, T, Sandve, SR, Heier, L, *et al*. Ancient hybridizations among the ancestral genomes of bread wheat. *Science* 2014;**345**(6194):1250092.
3. Guan, J, Garcia, DF, Zhou, Y, *et al*. The battle to sequence the bread wheat genome: a tale of the three kingdoms. *Genomics Proteomics Bioinformatics* 2020;**18**(3):221–9.
4. Chapman, JA, Mascher, M, Buluç, A, *et al*. A whole-genome shotgun approach for assembling and anchoring the hexaploid bread wheat genome. *Genome Biol* 2015;**16**(1):26.
5. Zimin, AV, Puiu, D, Hall, R, *et al*. The first near-complete assembly of the hexaploid bread wheat genome, *Triticum aestivum*. *Gigascience* 2017;**6**(11):doi:10.1093/gigascience/gix097.
6. Clavijo, BJ, Venturini, L, Schudoma, C, *et al*. An improved assembly and annotation of the allohexaploid wheat genome identifies complete families of agronomic genes and provides genomic evidence for chromosomal translocations. *Genome Res* 2017;**27**(5):885–96.
7. International Wheat Genome, Sequencing Consortium (IWGSC), Appels, R, Eversole, K, *et al*. Shifting the limits in wheat research and breeding using a fully annotated reference genome. *Science* 2018;**361**(6403):eaar7191.

8. Walkowiak, S, Gao, L, Monat, C, *et al.* Multiple wheat genomes reveal global variation in modern breeding. *Nature* 2020;**588**(7837):277–83.

9. Miga, KH, Koren, S, Rhie, A, *et al.* Telomere-to-telomere assembly of a complete human X chromosome. *Nature* 2020;**585**(7823):79–84.

10. Belser, C, Istace, B, Denis, E, *et al.* Chromosome-scale assemblies of plant genomes using nanopore long reads and optical maps. *Nat Plants* 2018;**4**(11):879–887.

11. Rousseau-Gueutin, M, Belser, C, Da Silva, C, *et al.* Long-read assembly of the *Brassica napus* reference genome Darmor-bzh. *Gigascience* 2020;**9**(12):giaa137.

12. Li, G, Wang, L, Yang, J, *et al.* A high-quality genome assembly highlights rye genomic characteristics and agronomically important genes. *Nat Genet* 2021;**53**(4):574–84.

13. Liu, J, Seetharam, AS, Chougule, K, *et al.* Gapless assembly of maize chromosomes using long-read technologies. *Genome Biol* 2020;**21**(1):121.

14. Tørresen, OK, Star, B, Mier, P, *et al.* Tandem repeats lead to sequence assembly errors and impose multi-level challenges for genome and protein databases. *Nucleic Acids Res* 2019;**47**(21):10994–11006.

15. Li, C, Xiang, X, Huang, Y, *et al.* Long-read sequencing reveals genomic structural variations that underlie creation of quality protein maize. *Nat Commun* 2020;**11**(1):17.

16. Ruan, J, Li, H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods* 2020;**17**(2):155–8.

17. Liu, H, Wu, S, Li, A, *et al.* SMARTdenovo: a de novo assembler using long noisy reads. *Gigabyte* 2021;doi:10.46471/gigabyte.15.

18. Kolmogorov, M, Yuan, J, Lin, Y, *et al.* Assembly of long, error-prone reads using repeat graphs. *Nat Biotechnol* 2019;**37**(5):540–6.

19. Vaser, R, Sović, I, Nagarajan, N, *et al.* Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Res* 2017;**27**(5):737–46.

20. https://github.com/nanoporetech/medaka.

21. Aury, J-M, Istace, B. Hapo-G, haplotype-aware polishing of genome assemblies with accurate reads. *NAR Genom Bioinform* 2021;**3**(2):lqab034.10.1093/nargab/lqab034

22. Zhu, T, Wang, L, Rimbert, H, *et al.* Optical maps refine the bread wheat *Triticum aestivum* cv. Chinese Spring genome assembly. *Plant J* 2021;**107**(1):303–14.

23. Rimbert, H, Darrier, B, Navarro, J, *et al.* High throughput SNP discovery and genotyping in hexaploid wheat. *PLoS One* 2018;**13**(1):e0186329.

24. Istace, B, Belser, C, Falentin, C, *et al.* Sequencing and chromosome-scale assembly of plant genomes, *Brassica rapa* as a use case. *Biology (Basel)* 2021;**10**(8):732.

25. Altschul, SF, Gish, W, Miller, W, *et al.* Basic Local Alignment Search Tool. *J Mol Biol* 1990;**215**(3):403–10.

26. Brinton, J, Ramirez-Gonzalez, RH, Simmonds, J, *et al.* A haplotype-led approach to increase the precision of wheat breeding. *Commun Biol* 2020;**3**(1):712.

27. De Oliveira, R, Rimbert, H, Balfourier, F, *et al.* Structural variations affecting genes and transposable elements of Chromosome 3B in wheats. *Front Genet* 2020;**11**:891.

28. Wick, RR, Judd, LM, Holt, KE. Performance of neural network basecalling tools for Oxford Nanopore sequencing. *Genome Biol* 2019;**20**(1):129.

29. Daron, J, Glover, N, Pingault, L, *et al.* Organization and evolution of transposable elements along the bread wheat chromosome 3B. *Genome Biol* 2014;**15**(12):546.

30. Wicker, T, Gundlach, H, Spannagl, M, *et al.* Impact of transposable elements on genome structure and evolution in bread wheat. *Genome Biol* 2018;**19**(1):103.

31. Leroy, P, Guilhot, N, Sakai, H, *et al.* TriAnnot: a versatile and high performance pipeline for the automated annotation of plant genomes. *Front Plant Sci* 2012;**3**:5.

32. Hao, M, Zhang, L, Ning, S, *et al.* The resurgence of introgression breeding, as exemplified in wheat improvement. *Front Plant Sci* 2020;11:252.

33. Kondrashov, FA. Gene duplication as a mechanism of genomic adaptation to a changing environment. *Proc Biol Sci* 2012;**279**(1749):5048–57.

34. Panchy, N, Lehti-Shiu, M, Shiu, S-H. Evolution of gene duplication in plants. *Plant Physiol* 2016;**171**(4):2294–316.

35. Huo, N, Zhang, S, Zhu, T, *et al.* Gene duplication and evolution dynamics in the homeologous regions harboring multiple prolamin and resistance gene families in hexaploid wheat. *Front Plant Sci* 2018;**9**:673.

36. Xu, J-H, Messing, J. Organization of the prolamin gene family provides insight into the evolution of the maize genome and gene duplications in grass species. *Proc Natl Acad Sci U S A* 2008;**105**(38):14330–5.

37. Lang, D, Zhang, S, Ren, P, *et al.* Comparison of the two up-to-date sequencing technologies for genome assembly: HiFi reads of Pacific Biosciences Sequel II system and ultralong reads of Oxford Nanopore. *Gigascience* 2020;**9**(12):giaa123.

38. Hon, T, Mars, K, Young, G, *et al.* Highly accurate long-read HiFi sequencing data for five complex genomes. *Sci Data* 2020;**7**(1):399.

39. Belser, C, Baurens, F-C, Noel, B, *et al.* Telomere-to-telomere gapless chromosomes of banana using nanopore sequencing. *Commun Biol* 2021;**4**(1):1047.

40. Lv, X, Chen, Z, Lu, Y, *et al.* An end-to-end Oxford Nanopore basecaller using convolution-augmented transformer. In: *2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. 2020:337–42.

41. Huang, N, Nie, F, Ni, P, *et al.* An attention-based neural network basecaller for Oxford Nanopore sequencing data. In: *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. 2019:390–4.

42. Alberti, A, Poulain, J, Engelen, S, *et al.* Viral to metazoan marine plankton nucleotide sequences from the Tara Oceans expedition. *Sci Data* 2017;**4**:170093.

43. fastxtend. https://www.genoscope.cns.fr/fastxtend/.

44. Istace, B, Belser, C, Aury, J-M. BiSCoT: improving large eukaryotic genome assemblies with optical maps. *PeerJ* 2020;**8**:e10150.

45. Rhie, A, Walenz, BP, Koren, S, *et al.* Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies. *Genome Biol* 2020;**21**(1):245.

46. Li, Heng. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv* 2013;10.48550/ARXIV.1303.3997.https://arxiv.org/abs/1303.3997

47. Durand, NC, Shamim, MS, Machol, I, *et al.* Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Syst* 2016;**3**(1):95–8.

48. Dudchenko, O, Batra, SS, Omer, AD, *et al.* De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds. *Science* 2017;**356**(6333):92–5.

49. Pedersen, BS, Quinlan, AR. Mosdepth: quick coverage calculation for genomes and exomes. *Bioinformatics* 2018;**34**(5):867–8.

50. Abyzov, A, Urban, AE, Snyder, M, *et al.* CNVnator: an approach to discover, genotype, and characterize typical and atypical CNVs from family and population genome sequencing. *Genome Res* 2011;**21**(6):974–84.

51. MAGATT pipeline, accessed date: 16 sept. 2021. https://forgemia.inra.fr/umr-gdec/magatt

52. Kent, WJ. BLAT—The BLAST-Like Alignment Tool. *Genome Res* 2002;**12**(4):656–64.

53. Wu, TD, Watanabe, CK. GMAP: a genomic mapping and alignment program for mRNA and EST sequences. *Bioinformatics* 2005;**21**(9):1859–75.

54. Quinlan, AR, Hall, IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 2010;**26**(6):841–2.

55. Kim, D, Paggi, JM, Park, C, *et al.* Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol* 2019;**37**(8):907–15.

56. Pertea, M, Pertea, GM, Antonescu, CM, *et al.* StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol* 2015;**33**(3):290–5.

57. Trapnell, C, Williams, BA, Pertea, G, *et al.* Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* 2010;**28**(5):511–5.

58. Holley, G, Melsted, P. Bifrost: highly parallel construction and indexing of colored and compacted de Bruijn graphs. *Genome Biol* 2020;**21**(1):249.

59. Hao, Z, Lv, D, Ge, Y, *et al.* RIdeogram: drawing SVG graphics to visualize and map genome-wide data on the idiograms. *PeerJ Comput Sci* 2020;**6**:e251.

60. Tang, H, Bowers, JE, Wang, X, *et al.* Synteny and collinearity in plant genomes. *Science* 2008;**320**(5875):486–8.

61. Zulkower, V, Rosser, S. DNA Features Viewer: a sequence annotation formatting and plotting library for Python. *Bioinformatics* 2020;**36**(15):4350–2.

62. Genoscope, accessed date: 16 sept. 2021. http://www.genoscope.cns.fr/plants/.

63. Renan-associated-data – GitHub repository, accessed date: 16 sept. 2021. https://github.com/institut-de-genomique/Renan-associated-data.

64. Aury, JM, Engelen, S, Istace, B, *et al.* Supporting data for "Long-read and chromosome-scale assembly of the hexaploid wheat genome achieves high resolution for research and breeding." *GigaScience Database* 2022. http://dx.doi.org/10.5524/102205.