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

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

* corresponding authors


[1] Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université Paris-Saclay, 91057 Evry, France
[2] GDEC, Université Clermont Auvergne, INRAE, UMR1095, 63000 Clermont-Ferrand, France
[3] Commissariat à l'Energie Atomique (CEA), Institut François Jacob, Genoscope, F-91057 Evry, France
[4] INRAE, CNRGV French Plant Genomic Resource Center, F-31320, Castanet Tolosan, France
[5] Mésocentre Clermont Auvergne, DOSI / Bâtiment Turing, 7 avenue Blaise Pascal, 63178 Aubière CEDEX
[6] Current address: Université Paris-Saclay, CEA, CNRS, Institute for Integrative Biology of the Cell (I2BC), 91198, Gif-sur-Yvette, France.

## Abstract

26

27 The sequencing of the wheat (*Triticum aestivum*) genome has been a methodological

28 challenge for many years due to its large size (15.5 Gb), repeat content, and hexaploidy.

29 Many initiatives aiming at obtaining a reference genome of cultivar Chinese Spring have

30 been launched in the past years and it was achieved in 2018 as the result of a huge effort to

31 combine short-read sequencing with many other resources. Reference-quality genome

32 assemblies were then produced for other accessions but the rapid evolution of sequencing

33 technologies offers opportunities to reach high-quality standards at lower cost. Here, we

34 report on an optimized procedure based on long-reads produced on the ONT (Oxford

35 Nanopore Technology) PromethION device to assemble the genome of the French bread

36 wheat cultivar Renan. We provide the most contiguous and complete chromosome-scale

37 assembly of a bread wheat genome to date. Coupled with an annotation based on RNA-Seq

38 data, this resource will be valuable for the crop community and will facilitate the rapid

39 selection of agronomically important traits. We also provide a framework to generate high-

40 quality assemblies of complex genomes using ONT.

41

## Introduction

43 Bread wheat (*Triticum aestivum*) is among the most important cereal crops and a better

44 knowledge in the area of wheat genomics is needed to face the main challenge of ensuring

45 food security to a growing population in the context of climate change. Improving productivity

46 requires both that local producers adapt their practices to increase their climate resilience

47 and a better understanding of the wheat production systems. In this context, a better

48 knowledge of the wheat genome and its gene content, but also the sequencing of numerous

49 accessions, are essential.

50 However, the genome of bread wheat is particularly characterized by its complexity. Indeed,

51 this hexaploid genome is the result of two interspecific hybridization events. The earliest

52 cultivated wheat was diploid, but humans have intensified the cultivation of polyploid

53 species. Recent studies show that these polyploid species appear to be advantaged by their

54 genomic plasticity[1]. Indeed, modifications of the gene space and related elements are

55 buffered by the polyploid nature of wheat and open a wider field to selection. Bread wheat is

56 composed of three subgenomes A, B and D derived from three ancestral diploid species that

57 diverged between 2.5 and 6 million years ago[2].

58 The wheat genome is one of the largest among sequenced plant genomes (15.5 Gb), mainly

59 composed of repetitive sequences (ca. >85%), and contains many homoeologous regions

60 between the three subgenomes (A, B and D). Repetitive sequences and polyploidy pose

61 serious challenges in the generation of genome assemblies. The adventure of sequencing

62   the hexaploid wheat genome began in 2005 with the creation of the International Wheat

63   Genome Sequencing Consortium (IWGSC)[3]. With the advent of sequencing technologies,

64   the wheat genome has been competitively sequenced several times[4–6]. The first

65   reference-quality genome sequence with a comprehensive annotation was published by the

66   IWGSC in August 2018[7] for the accession Chinese Spring (CS). This assembly represents

67   a tremendous resource for the scientific community and offers the promise of facilitating and

68   accelerating breeding efforts.

69   More recently, fifteen genomes of hexaploid wheat have been published[8] which represents

70   a new step in the knowledge of the wheat model. Ten of these new wheat genomes have

71   been assembled at the chromosome level, allowing for comparative analysis on a scale that

72   was previously impossible. While being a valuable and highly validated resource using

73   multiple technologies, these assemblies were produced using short-read technologies and

74   therefore may contain a higher number of gaps compared to genomes assembled with long

75   reads[9–13]. In 2017, an assembly of the CS genome using long-reads was produced[5],

76   although not annotated, highlighting the added-value of long-reads in such complex

77   genomes. By accumulating long-read assemblies, the scientific community is now aware of

78   the flaw in short-read strategies. Indeed they underestimate the repetitive content of the

79   genome and more importantly can lack tandemly duplicated genes[14,15]. Several years

80   ago, Pacific Biosciences (PACBIO) and Oxford Nanopore (ONT) sequencing technologies

81   were commercialized with the promise to sequence long DNA fragments and revolutionize

82   complex genome assemblies.

83   Here, we report the first hexaploid wheat genome based on ONT long-reads. We sequenced

84   the genome of the French variety Renan, one of the most used varieties in organic farming.

85   The Renan genome carries multiple resistance genes against fungal pathogens (leaf rust,

86   stem rust, yellow rust, eyespot) originating from introgression of DNA regions coming from

87   the wild species *Aegilops ventricosa*. We used the PromethION device and organized the

88   assembled contigs at the chromosome scale using optical maps (BioNano Genomics, BNG)

89   and Hi-C libraries (Arima Genomics, AG). This assembly has a contig N50 of 2.2 Mb, which

90   is a 30-fold improvement over existing chromosome-scale assemblies.

91

## Results

93   **Genome sequencing and optical maps**

94   We sequenced genomic DNA using 20 ONT flow cells (2 MinION and 18 PromethION)

95   which produced 12M reads representing 1.1 Tb. All the reads were originally base called

96   using the guppy 2.0 software, but given the improvement of guppy software during our

97   project, we decided to call bases using a newer version of the guppy software (version 3.6

98   with High Accuracy setting). This dataset represented a coverage of 63x of the hexaploid

99   wheat genome and the read N50 was of 24.6kb. More importantly, we got 3.1M reads larger

100  than 50kb representing a 14x genome coverage (Table S1). In addition, we generated

101  Illumina short-reads and long-range data for respectively polishing and organizing nanopore

102  contigs. We produced an optical map using the Saphyr instrument commercialized by

103  Bionano Genomics (BNG). High molecular weight DNA was extracted and labeled using the

104  Direct Label and Stain Chemistry (DLS) with the DLE-1 enzyme. The DLE-1 optical map was

105  assembled using proprietary tools provided by BNG and had a cumulative size of 14.9 Gb

106  with an N50 of 37.5 Mb (Table S2). Four Hi-C libraries from two biological replicates were

107  prepared using the Arima Genomics protocol and sequenced on an Illumina sequencer to

108  reach 537 Gb i.e., a depth of 35x. We used a sample of 240 million read pairs (72 Gb, 5x) to

109  build a Hi-C map.

## Genome assembly

111  Since the dataset was too large for many long-read assemblers, we sampled a 30x coverage

112  by selecting the longest reads (Table S1). This subset was assembled using multiple

113  assembly tools dedicated to processing this large amount of data (Redbean[16],

114  SMARTdenovo[17] and Flye[18]). SMARTdenovo is not among the fastest algorithms and

115  has not been updated for several years, but since it can be easily parallelized, it remains an

116  interesting choice for assembling large genomes. The overlap and consensus calculations

117  were split into 60 chunks and each were run on a 32-core server and took about two days

118  and ten hours respectively. In comparison, Redbean was able to generate an assembly after

119  just seven days on a 64-core server with 3TB of memory while Flye needed 43 days on the

120  same computer server. Surprisingly, the redbean assembly had a cumulative size two times

121  higher than the expected genome size (29.6Gb vs 14.5Gb), a low contiguity and contained a

122  large amount of short contigs. The SMARTdenovo and Flye assemblies were highly

123  comparable, but Flye was the most contiguous (contigs N50 of 1.8 Mb vs 1.1 Mb) and

124  SMARTdenovo had a cumulative size closer to the expected one (14.1 Gb vs 13.0 Gb, Table

125  S3). Additionally, even though the assemblies were polished later, the raw SMARTdenovo

126  assembly contained a higher number of complete BUSCO genes (83.0% vs 49.5%) which

127  indicates that its consensus module is more efficient.

128  The SMARTdenovo and Flye assemblies were successively polished using Racon[19] and

129  Medaka (https://github.com/nanoporetech/medaka) with long reads and Hapo-G[20] with

130  short reads. Polished contigs were validated and organized into scaffolds using the DLE-1

131  optical map and proprietary tools provided by BNG. As expected, due to its lower cumulative

132     size, Flye scaffolds contained a larger proportion of unknown bases (851 Mb and 262 Mb).

133     Based on these results (proportion of gaps and gene completion), the assembly produced by

134     SMARTdenovo[17] was selected (Table S4). Local contig duplications (negative gaps) were

135     resolved using BiSCoT[22], which improved the contigs N50 from 1.2 Mb up to 2.1 Mb. Finally,

136     the resulting assembly was polished one last time using Hapo-G[20] with short reads. This

137     led to 2,904 scaffolds (larger than 30kb) representing 14.26 Gb with a N50 of 48 Mb (79

138     scaffolds) and a maximum scaffold size of 254 Mb. Thus, the genome size is in the same

139     range as all other available reference quality assemblies of *T. aestivum*: e.g. 14.29 Gb for

140     *cv.* LongReach Lancer, 14.55 Gb for *cv.* Chinese Spring, and 14.96 Gb for *cv.* SY Mattis.

141     **Construction and validation of pseudomolecules**

142     We then guided the construction of the 21 chromosome sequences (i.e. pseudomolecules)

143     based on collinearity with the CS (Chinese Spring) RefSeq Assembly v2.1[22]. Given the

144     complexity of this hexaploid genome, we established a dedicated approach in order to

145     anchor each Renan scaffold based on similarity search against CS. To avoid problems due

146     to multiple mappings, we selected a dataset of uniquely mappable sequences. Genes are

147     not uniquely mappable since most of them are repeated as three homoeologous copies

148     sharing on average 97% nucleotide identity. In addition, the gene density (1 gene every

149     130kb on average) is too low to anchor small Renan scaffolds that do not carry genes. Thus,

150     we used 150 bp tags corresponding to the 5' and 3' junctions between a transposable

151     element (TE) and its insertion site (75 bps on each side) which are called ISBP (Insertion

152     Site-Based Polymorphism) markers and are highly abundant and uniquely mappable in the

153     wheat genome[23]. We designed a dataset of 5.76 million ISBPs from CS assembly which

154     represent 1 ISBP every 2.5kb. Their mapping enabled the anchoring of 2,566 scaffolds on

155     21 pseudomolecules representing 14.20 Gb (99% of the assembly). We then used Hi-C data

156     to validate the assembly and to correct the mis-ordered and mis-oriented scaffolds. The Hi-C

157     map revealed only a few inconsistencies, demonstrating that the collinearity between CS

158     and Renan was strong enough to guide the anchoring in a very accurate manner. The Hi-C

159     map-based curation led to the detection of 18 chimeric scaffolds that were split into 2 or 3

160     pieces and to the correction of the location and/or orientation of 198 scaffolds. The final

161     assembly was composed of 21 pseudomolecules (Figure 1) with 338 unanchored scaffolds

162     representing 61 Mb only.

163     **Quality assessment of the assembly**

164     First, we calculated the overall quality of the sequence using Merqury and Illumina reads.

165     We obtained an average quality value (QV) of 32.8, a lower QV than that obtained with

5

166  short-reads assemblies, but consistent with QV already reported for plant genomes
167  sequenced by ONT[24]. Indeed, using Illumina reads and the CS RefSeq v2.1 assembly,
168  Merqury computed a QV of 44.5 (Table 1). This shows that per-base quality is still an issue,
169  at least with the version of the technology used in this study. However, this could be
170  tempered by the fact that coding regions, due to lower repeating regions, may have higher
171  precision.

172  The completeness and quality of the assembly was estimated by searching for the presence
173  of known genes, i.e. the 107,891 High Confidence (HC) genes predicted in CS RefSeq v1.1.
174  We used BLAST[25] to search for the presence of each of the 461,476 exons larger than 30
175  bps in the Renan scaffolds, and we considered only matches showing at least 95% identity
176  over at least 95% query length. We found hits for 96.2% of the query exons with on average
177  99.3% identity, suggesting that the gene space is assembled at a high-quality level. The
178  missing genes/exons would correspond, in most of the cases, to real presence/absence
179  variations between CS and Renan while the nucleotide divergence between exons is 0.7%.
180  It was the first evidence that homoeologous gene copies, sharing on average 97%
181  identity[7], were not collapsed in the Renan assembly. We confirmed this by showing that
182  62% of the CS exons are strictly identical in Renan (and carried by the same chromosome).
183  Such level of nucleotide divergence between CS and Renan is similar to what has been
184  shown through whole genome alignments (Brinton et al. 2020).

185  We then assessed the assembly quality of the TE space by aligning the complete dataset of
186  ISBP markers of CS onto the Renan assembly. We found that 94% markers were conserved
187  (at least 90% identity over 90% query length) i.e., present in the assembly, revealing that the
188  TE space is extremely close to completeness. Indeed, 6% of missing markers is similar to
189  the proportion of expected Presence-Absence variations (PAVs) affecting TEs[26].

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

194  **Impact of the polishing**

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

201 were able to find 80.4% of the IBSP markers but only 7% were aligned without mismatch
202 between the two genotypes (Table S5). When polished with long-reads, the BUSCO score
203 reached 96.7% and 92.9% of the IBSP markers were retrieved (including 28.0% with perfect
204 matches). The subsequent polishing step with short reads weakly decreased the BUSCO
205 score (from 96.7% to 96.6%), but the proportion of duplicated genes increased from 83.1%
206 to 87.0% which is here wanted because in the case of a hexaploid genome most of the
207 genes are in three copies. Moreover, the proportion of perfectly aligned ISBP markers
208 drastically increased from 28.0% up to 58.9%. Although the polishing with short reads
209 weakly impacts the BUSCO conserved genes, the IBSP markers underline its importance in
210 the case of long reads assemblies. Since ISBPs are unique tags sampling the whole
211 genome, this analysis revealed that nucleotide errors were frequent before polishing,
212 affecting half of the sample loci. Thus, we showed that the polishing steps were successful,
213 even in this large and polyploid genome, and drastically improved the quality of the
214 consensus.

## Recent improvement of the ONT technology

216 Oxford Nanopore Technology is evolving rapidly, and improvements to the base calling
217 softwares are frequent, allowing old data to be analyzed with the aim of improving read
218 accuracy and subsequent analysis. To measure the gain brought by each new version
219 during this project, we analyzed a subset of ultra-long reads (longer than 100kb) with
220 different basecallers or versions of the same basecaller: guppy 2.0, guppy 3.0.3 (High
221 Accuracy mode), guppy 3.6 (High Accuracy mode) and the recent bonito v0.3.1. We
222 observed a strong difference in accuracy, of around 7%, between guppy 2.0 and the newer
223 basecaller (bonito v0.3.1), representing the gain over the last two years (Figure S2A). This
224 significant improvement could lead nanopore users to reanalyze their old sequencing data to
225 improve the quality of their assemblies. As an example, the accuracy of raw nanopore reads
226 gained about 2% on average using guppy 3.6 (Table S6). We observed a reduction of the
227 number of contigs of 19%, and an improvement of the contig N50 of 26%. Likewise, the
228 cumulative size is slightly higher in the guppy 3.6 assembly, which may underline a smaller
229 amount of collapsed repetitive regions (Table S7).
230 More importantly, the identity percentage obtained when aligning ONT reads on the wheat
231 assembly is lower than what was obtained on yeast and human samples (Figure S2B). This
232 difference can be explained by the fact that, first, the consensus of the wheat genome is not
233 perfect and secondly, that basecallers are trained on a mixture which contains yeast and
234 human data. Indeed, DNA modification patterns can differ between taxa, and read accuracy
235 seems better when the model was trained on native DNA from the same species[27]. This

236  huge difference between the read accuracy of yeast and wheat samples should motivate

237  nanopore users to train basecaller models to their targeted species.


238  **Annotation of transposable elements and protein-coding genes**

239  We annotated TEs based on similarity search against our wheat-specific TE library

240  ClariTeRep[28] and raw results were then refined using CLARITE, a homemade program

241  able to resolve prediction conflicts, merge adjacent features into a single complete element,

242  and identify nested insertion patterns. We detected 3.9 million copies of TEs in the Renan

243  genome assembly, representing 12.0 Gb i.e. 84% of the assembly size. The proportions of

244  each superfamily were similar to what has been described for CS[29] (Table 2).

245  Gene annotation was achieved by, first, transferring genes predicted in CS RefSeq v2.1 by

246  homology using the MAGATT pipeline[22]. This allowed us to accurately transfer 105,243

247  (out of 106,801; 98%) HC genes and 155,021 (out of 159,846; 97%) Low Confidence genes.

248  Such a transfer of genes predicted in another genotype (here CS) avoided genome-wide *de*

249  *novo* gene prediction that may artificially lead to many differences between the annotations.

250  We thus focused *de novo* predictions using TriAnnot[30] only on the unannotated part of the

251  genome, representing 8.5% of the 14.2 Gb, after having masked transferred genes and

252  predicted TEs. For that purpose, we produced RNASeq data for Renan from 28 samples

253  corresponding to 14 different organs/conditions in replicates: grains at four developmental

254  stages (100, 250, 500, and 700 degree days) under heat stress and control conditions,

255  stems at two developmental stages, leaves at three stages, and roots at one stage),

256  representing on average 78.8 million read-pairs per sample i.e 2.2 billion read-pairs in total.

257  This method allowed us to predict 4,440 genes specific to Renan compared to CS i.e., 4% of

258  the gene complement. This is consistent with the extent of structural variations affecting

259  genomes of *Triticeae*[26]. Transfer of known genes, novel predictions, and manual curation

260  (limited to storage protein encoding genes), led us to annotate 109,552 protein-coding genes

261  on the Renan pseudomolecules.


262  **Comparison with existing hexaploid genome assemblies**

263  We compared our long-read assembly with 10 other available chromosome-scale

264  assemblies of wheat genomes. Although the gene content was similar between the different

265  assemblies, as expected, the assemblies based on short reads had a lower contiguity

266  (contig N50 values lower than 100kb compared to the 2 Mb of the assembly of the Renan

267  genome, Figure 2A-B). Logically, they also contained more gaps (around 40 times, Figure

268  2C). Interestingly, we found more gaps per Mb in the D subgenome compared to the A and

269  B subgenomes in Renan (Figure S3). This indicates that the D subgenome is more difficult

270　to assemble even though it has a smaller genome size and contains less repetitive

271　elements. The same trend was already observed in another polyploid genome, the rapeseed

272　and its two subgenomes A and C[11]. Chromosomes from the different assemblies had

273　similar length except for the Arina*LrFor* and the SY_Mattis variety in which a translocation

274　has been previously described between chromosomes 5B and 7B[8] (Figure 2D).

275　In addition, we generated dotplots between CS and Renan homeologous chromosomes and

276　confirmed the strong collinearity between the two genomes (Figure 3). Whole chromosome

277　alignments highlighted 16 large-scale inversions (>5 Mb; up to 118 Mb) on 10 chromosomes

278　and 1 translocation of a ca. 45 Mb segment on chromosome 4A. We performed the same

279　comparisons with the 10 other available genomes of related varieties assembled at the

280　pseudomolecule level (Supplementary Data 1). It showed that only 2 of these inversions are

281　specific to Renan while the others are shared between several accessions. They correspond

282　to regions of 23 Mb on chr6B (position 398-421 Mb) and 10 Mb on chr7B (position 267-277

283　Mb).

## Haplotype characterisation

285　Crop breeding involves the selection of desired traits and their combination to generate

286　improved genotypes. Generally, these traits correspond to genomic regions carrying genetic

287　variations or genes[31]. These regions of interest are inherited from their parents in the form

288　of large genomic blocks. The availability of several assemblies of the wheat genome now

289　allows the detection of these haplotypic blocks. Using the 11 chromosome-scale wheat

290　assemblies and an approach based on colored de Bruijn graphs, we investigated these

291　haplotypic blocks and applied our method to the 21 chromosomes of wheat. First, a colored

292　de Bruijn graph was built for each chromosome, where each colour represents a different

293　cultivar. Short (1kb) and evenly distributed (every 20kb) markers were extracted from each

294　chromosome and compared to the colored de Bruijn graph to extract their presence/absence

295　in each wheat cultivar. On each chromosome, the 15 most abundant presence/absence

296　profiles were selected and used to characterise haplotypic blocks. The haplotype blocks of

297　chromosome 6A, which is associated with productivity traits (as for example yield, grain size

298　and height), have already been expertized using a different method[31]. We obtained similar

299　results (Figure 4), except for the Chinese Spring chromosome 6A. Previous results have

300　assigned a unique haplotype to this wheat line. But in our case Chinese Spring exhibits the

301　same haplotype as SY Mattis, Jagger, Lancer and Norin61, which had previously been

302　described as sharing the same haplotype. These differences may be explained by the

303　stringency of the comparison, which perhaps should be adjusted separately for each

304　chromosome. Concerning the Renan cultivar, the chromosome 6A has haplotype blocks

305  similar to those of the ArinaLrFor line. Additionally, we used this method to investigate
306  haplotypic blocks that are specific to one or a subset of wheat cultivars.

**Identification of introgressions**

308  Introgression is an important source of genetic variation which is generally the signature of
309  breeding programmes, especially in wheat[32]. Several introgressions have already been
310  reported[8], notably in chromosomes 2B and 3D in LongReach Lancer and in chromosome
311  2A in Jagger, Mace, SY Mattis and CDC Stanley. Using our approach, we were able to
312  clearly identify the two introgressions in LongRead Lancer (Figure 5a), and the *Ae.*
313  *ventricosa* introgression in chromosome 2A (Figure 5b). In addition, we found that this
314  introgression of *Ae. ventricosa* is also present in the Renan cultivar (Figure 5b). The optical
315  map was aligned with this 34 Mb region of Renan and validated the correct structure of this
316  important region carrying multiple resistance genes (Yr17, Lr37, Sr38, Cre5). More
317  importantly, the 34 Mb consisted of 22 contigs in Renan and 2,339 in Jagger. A comparison
318  of the fragmentation near the introgression point is presented in Figure 5d and shows a large
319  difference between the long- and short-reads assemblies. Additionally, we also identified
320  several candidate introgressions, which had already been spotted through retrotransposon
321  profiles[8]: i) a 45-Mb region on chromosome 2D which is shared between the lines Julius,
322  ArinaLrFor, SY Mattis, Jagger and also Renan (Figure 6a); ii) a 53-Mb region at the end of
323  chromosome 3D in Lancer (Figure 6b); iii) a 48-Mb region at the beginning of chromosome
324  3D in SY Mattis (Figure 6b) and iv) the *Ae. ventricosa* introgression of 30-Mb in chromosome
325  7D which carries Pch1 resistance gene (Figures 6c).
326  Moreover, a known large-scale structural variation in chromosomes 5B and 7B of ArinaLrFor
327  and SY Mattis cultivars was also easily identifiable using haplotypic blocks of individual
328  chromosomes (Figure S4).

**Comparative analysis of a storage protein coding gene cluster in *T. aestivum***

330  Tandem duplications are an important mechanism in plant genome evolution and
331  adaptation[33,34] but the assembly of tandemly duplicated gene clusters is difficult,
332  especially with short reads. In order to illustrate the gain brought by this optimized assembly
333  process, we focused on an important locus on chromosome 1B known to carry multiple
334  copies of storage protein and disease resistance genes[35,36]. Among them, the genes
335  encoding omega-gliadins are not only duplicated in tandem, but are also composed of
336  microsatellite DNA in their coding part, making them particularly hard to assemble properly
337  from short reads. We compared orthologous regions harboring these genes between CS and
338  Renan, spanning 1.58 Mb and 2.32 Mb, respectively. The CS region was more fragmented

10

339   with 101 gaps versus only 3 in Renan (Figures 5a). The number of copies of omega-gliadin

340   encoding genes was quite similar: 9 in CS and 10 in Renan. The most striking difference

341   came from the completeness of the microsatellite motifs: 8 copies out of 9 contain N

342   stretches in CS RefSeq v2.1, revealing that the microsatellite is usually too large to be fully

343   assembled with short reads (Figure 5b). In contrast, all 10 copies predicted in Renan were

344   assembled completely. More generally, we mapped the corresponding proteins back to the

345   locus and showed that it was better reconstructed in the Renan assembly, with a mean

346   protein alignment length of 99% compared to 58% in CS (Figure 5c). In addition, the optical

347   map was used to validate the structure of this region in Renan and the assembly was

348   consistent with the three maps of this loci (Figure 5d).

349

350   **Comparative analysis of the locus that provides resistance to the orange**

351   **wheat blossom midge**

352   Like a few other wheat cultivars, Renan is resistant to the orange wheat blossom midge

353   (OWBM). The *Sm1* gene is known to confer resistance to wheat and a previous study has

354   shown that CDC Landmark is also resistant to the OWBM, and carries a 7.3-Mb haplotype

355   within the *Sm1* locus on chromosome 2B[8]. We extracted and aligned the corresponding

356   region of CDC Landmark on each cultivar, to precisely locate the corresponding region on

357   each chromosome 2B. From these eleven regions of 1-2 Mb, we computed the haplotypic

358   blocks using a higher resolution than previously (1 kb marker every 5 kb). This analysis

359   revealed a strong similarity of the *Sm1* locus between CDC Landmark and Renan (Figure

360   8a), the presence of the *Sm1* gene in blocks shared between the two cultivars.  In addition, a

361   comparison of the fragmentation of these two regions underlines the higher contiguity of the

362   Renan assembly, with 4 contigs in the Renan *Sm1* locus compared to 62 in CDC Landmark

363   (Figure 8b). The *Sm1* locus of Renan is in agreement with the optical map and shows clearly

364   the three remaining gaps that may correspond to smaller and unanchored contigs.

365

# Discussion

367   In this study, we showed that the recent improvement of the Oxford Nanopore Technology,

368   in terms of error rate and throughput, has opened up new perspectives in the age of long-

369   read technologies. Indeed, the sequencing and assembly of complex genomes, like the

370   hexaploid wheat, is now accessible to sequencing facilities. Additionally, the ability to

371   sequence ultra-long reads using ONT devices is a real advantage over the other long-read

372   technology, namely PACBIO. In this study, we were able to generate a coverage of 14X with

373   reads longer than 50kb, whereas PACBIO libraries, used to generate HiFi (High-Fidelity)

374   reads, are generally sized around 15kb[37,38]. Several studies have already underlined the

375  positive impact of these ONT ultra-long reads on the assembly contiguity[9,37,39]. In
376  contrast, the error rate that was previously a thorn in their side has been drastically reduced
377  over the last year. Herein we reported a quality score near Q10-Q15 for individual ONT
378  reads, as already shown[27], which is still far from what HiFi reads can provide, generally
379  near Q30[37]. The high accuracy of HiFi reads might be sufficient to distinguish copies from
380  repeat regions if they present few variations. The impact of ultra-long reads will lie mainly in
381  the case of identical repeats, and obviously, the presence of these particular cases will
382  depend on the evolutionary history of the studied genomes. In addition, this high error rate
383  has an impact on the consensus quality, and at the moment, a combination of ONT and
384  Illumina reads is still needed to achieve a decent per-base accuracy.

385

386  By following basecallers evolution, we noticed that the gain when using recent basecaller is
387  high and we guess this observation will encourage users to reprocess older data. However,
388  this is not trivial and it requires sufficient computing resources. Interestingly, we observed
389  that the error rate of ONT data is organism dependent and that the training of basecaller has
390  a significant impact on the overall quality of the reads[27]. This is, in our opinion, an
391  important fact because a large proportion of *de novo* assemblies now concern non-model
392  organisms and users will have to address this limitation of current software. There are
393  existing methods to train the basecaller on non-model species[40,41], but this can still be a
394  big barrier, depending on the size of the dataset, for many end users. However, as
395  highlighted in this study, the combination of long- and short-reads sequencing with polishing
396  methods greatly improves the consensus sequence of a given genome assembly and these
397  algorithms seem sufficient at least in coding regions.

398

399  Even though there are now several chromosome-scale assemblies of the hexaploid wheat
400  genome, this assembly of the Renan variety based on long-reads will benefit biologists and
401  geneticists as it offers a high resolution. We show that our chromosome-scale assembly of
402  Renan based on long reads can bring new insight into genomic regions of interest. In
403  particular, in regions that carry multiple resistance genes, as a large *Ae. ventricosa*
404  introgression shared with other cultivars on chromosome 2A and a unique *Ae. ventricosa*
405  introgression on chromosome 7D. The lower number of gaps in these regions will help to
406  localize genes of interest and to have a better understanding of the impact of these
407  introgressions. Additionally, we demonstrated by examining two important locus, containing
408  prolamin and resistance genes that such regions are truly enhanced and contain very few
409  gaps compared to assemblies based on short reads.

410

411   Moreover, unlike recent chromosome-scale assemblies, Renan's gene prediction is not only
412   a projection of Chinese Spring gene models, but also includes *de novo* annotation with RNA-
413   Seq data which is of real benefit for the construction of pan genome (or pan annotation) or
414   when cultivar-specific genes are examined. For all of these reasons, we believe this high
415   resolution assembly will benefit the wheat community and help breeding programs dedicated
416   to the bread wheat genome.

417

# 418   Methods

### 419   **Plant material and DNA extraction**

420   *Triticum aestivum* cv. Renan seeds were provided by the INRAE Biological Resource Center
421   on small grain cereals and grown for two weeks and a dark treatment was applied on the
422   seedlings for two days before collecting leaf tissues.

423   For the sequencing experiments, DNA was isolated from frozen leaves using QIAGEN
424   Genomic-tips 100/G kit (Cat No./ID: 10243) and following the tissue protocol extraction.
425   Briefly, 1g of leaves were ground in liquid nitrogen with mortar and pestle. After 3h of lysis
426   and one centrifugation step, the DNA was immobilized on the column. After several washing
427   steps, DNA is eluted from the column, then desalted and concentrated by alcohol
428   precipitation. The DNA is resuspended in the TE buffer.

429   To generate the optical map, uHMW DNA were purified from 0.5 gram of very young fresh
430   leaves according to the Bionano Prep Plant tissue DNA Isolation Base Protocol (30068 -
431   Bionano Genomics) with the following specifications and modifications. Briefly, the leaves
432   were fixed using a fixing solution (Bionano Genomics) containing formaldehyde (Sigma-
433   Aldrich) and then grinded in a homogenization buffer (Bionano Genomics) using a Tissue
434   Ruptor grinder (Qiagen). Nuclei were washed and embedded in agarose plugs. After
435   overnight proteinase K digestion in Lysis Buffer (Bionano Genomics) and one hour treatment
436   with RNAse A (Qiagen), plugs were washed four times in 1x Wash Buffer (Bionano
437   Genomics) and five times in 1x TE Buffer (ThermoFisher Scientific). Then, plugs were
438   melted two minutes at 70°C and solubilized with 2 µL of 0.5 U/µL AGARase enzyme
439   (ThermoFisher Scientific) for 45 minutes at 43°C. A dialysis step was performed in 1x TE
440   Buffer (ThermoFisher Scientific) for 45 minutes to purify DNA from any residues. The DNA
441   samples were quantified by using the Qubit dsDNA BR Assay (Invitrogen). Quality of
442   megabase size DNA was validated by pulsed field gel electrophoresis (PFGE).

**Illumina Sequencing**

DNA (1.5μg) was sonicated using a Covaris E220 sonicator (Covaris, Woburn, MA, USA). Fragments (1μg) were end-repaired, 3′-adenylated and Illumina adapters (Bioo Scientific, Austin, TX, USA) were then added using the Kapa Hyper Prep Kit (KapaBiosystems, 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 (KapaBiosystems), 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 (https://www.genoscope.cns.fr/fastxtend/).

**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 Oxford Nanopore Technologies (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**

Labeling and staining of the uHMW DNA were performed according to the Bionano Prep Direct Label and Stain (DLS) protocol (30206 - Bionano Genomics). Briefly, labeling was performed by incubating 750 ng genomic DNA with 1× DLE-1 Enzyme (Bionano Genomics) for 2 hours in the presence of 1× DL-Green (Bionano Genomics) and 1× DLE-1 Buffer (Bionano Genomics). Following proteinase K digestion and DL-Green cleanup, the DNA backbone was stained by mixing the labeled DNA with DNA Stain solution (Bionano Genomics) in presence of 1× Flow Buffer (Bionano Genomics) and 1× DTT (Bionano Genomics), and incubating overnight at room temperature. The DLS DNA concentration was measured with the Qubit dsDNA HS Assay (Invitrogen).

477 Labelled and stained DNA was loaded on Saphyr chips. Loading of the chips and running of
478 the Bionano Genomics Saphyr System were all performed according to the Saphyr System
479 User Guide (30247 - Bionano Genomics). Data processing was performed using the
480 Bionano Genomics Access software.
481 A total of 4541 Gb data were generated. From this data, molecules with a size larger than
482 150kb were filtered generating 1931 Gb of data. These filtered data, corresponding to 128x
483 coverage of the *Triticum aestivum* cv. Renan consists of 7,810,298 molecules with an N50 of
484 237.5kb and an average label density of 14.3/100kb. The filtered molecules were aligned
485 using RefAligner with default parameters. It produced 1053 genome maps with a N50 of 37.5
486 Mbp for a total genome map length of 14946.8 Mbp.

## RNA extraction

488 Several tissues (stem, leaves, root or grain) were collected on plants with different growth
489 conditions and of different ages. Each of these 28 tissues was subjected to RNA extraction
490 with the following protocole: 200mg to 1g of fine powder was put in a 50ml falcon tube with
491 4.5 ml of NTES buffer [0.1 M NaCl, 1% SDS, 10 mM Tris-HCl (pH 7.4), 1 mM EDTA(pH 8)].
492 After vortexing the tube, 3ml of phenol-chloroforme-IAA were added. The tube was mixed for
493 10 minutes and centrifuged for 20 minutes at 5,000 rpm (15°C). The aqueous phase was
494 collected and placed in a new 15ml tube. 3ml of phenol-chloroforme-IAA were added. The
495 tube was mixed for 10 minutes and centrifuged for 20 minutes at 5,000 rpm (15°C). The
496 aqueous phase was collected and placed in a new 50ml tube. 1/10 of AcNa 3M (pH 5.2)
497 and 2 volumes of 100% ethanol were added. The tube was mixed gently by turning and
498 centrifuged 20 minutes at 5,000 rpm (4°C). The supernatant was removed. The precipitate
499 was dried and resuspended in 20 µl RNAse free water. A treatment with DNase was
500 realized and the RNA were purified on a MinElute column (Qiagen). A second treatment with
501 DNAse was realized by adding DNAse directly on the filter. After ethanol cleanup, the
502 column was eluted with 14 µl of RNAse free water. The quality of the RNA was evaluated
503 using RNA 6000 Nano Assay chip for size and RIN estimation and spectrophotometry
504 (A260/A280 and A260/A230 ratios) for purity estimation. The RNA were quantified using
505 Qubit RNA high sensitivity Assay kit (Invitrogen).

## RNA sequencing

507 RNA-Seq library preparations were carried out from 500ng to 2000ng of total RNA using the
508 TruSeq Stranded mRNA kit (Illumina, San Diego, CA, USA), which allows mRNA strand
509 orientation (sequence reads occur in the same orientation as antisense RNA). Briefly,
510 poly(A)+ RNA was selected with oligo(dT) beads, chemically fragmented and converted into

15

511    single-stranded cDNA using random hexamer priming. Then, the second strand was

512    generated to create double-stranded cDNA. cDNA were then 3'-adenylated, and Illumina

513    adapters were added. Ligation products were PCR-amplified. Ready-to-sequence Illumina

514    libraries were then quantified by qPCR using the KAPA Library Quantification Kit for Illumina

515    Libraries (KapaBiosystems, Wilmington, MA, USA), and libraries profiles evaluated with an

516    Agilent 2100 Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA). Each library was

517    sequenced using 151 bp paired end reads chemistry on an Illumina NovaSeq 6000

518    sequencer (Illumina, San Diego, CA, USA).

**Long reads genome assembly**

519

520    The 20 ONT runs were basecalled using two versions of guppy: 3.3 HAC and 3.6 HAC

521    (Table S6). We monitored the gain of each guppy basecaller release and evaluated three

522    different assemblers in the context of large genomes: Redbean[16] v2.5 (git commit

523    3d51d7e), SMARTdenovo[17] (git commit 5cc1356) and Flye[18] v2.7 (git commit 5c12b69).

524    All assemblers were launched using a subset of reads consisting of 30X of the longest reads

525    (Table S3). Then, we selected one of the assemblies based not only on contiguity metrics

526    such as N50 but also cumulative size, proportion of unknown bases. The Flye (longest

527    reads) and SMARTdenovo (all reads) assemblies were very similar in terms of contiguity but

528    we decided to keep the SMARTdenovo assembly as its cumulative size was higher. The

529    SMARTdenovo assembler using the longest reads resulted in a contig N50 of 1.1Mb and a

530    cumulative size of 14.07Gb. As nanopore reads contain systematic error in homopolymeric

531    regions, we polished the consensus of the selected assembly with nanopore reads as input

532    to the Racon (v1.3.2, git commit 5e2ecb7) and Medaka softwares. In addition, we polished

533    the assembly two additional times using Illumina reads as input to the Hapo-G tool (v1.0, git

534    commit ).

**Long range genome assembly**

535

536    The Bionano Genomics scaffolding workflow (Bionano Solve version 3.5.1) was launched

537    with the nanopore contigs and the Bionano map. We found in several cases that the

538    nanopore contigs were overlapping (based on the optical map) and these overlaps were

539    corrected using the BiSCoT software[21] with default parameters. Finally, the consensus

540    sequence was polished once more using Hapo-G and short reads, to ensure correction of

541    duplicate regions that were collapsed (Table S4).

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

542

543    The quality value (QV) of the Renan and CS assemblies was obtained using Merqury[43].

544    First, 31-mers were extracted from the Renan and CS Illumina sequencing reads

545    (accessions SRR5893651, SRR5893652, SRR5893653 and SRR5893654) and then the QV

546    of each genome assembly was computed using Merqury (version 1.3, git commit 6b5405e).

547    We used BLAST[25] to search for the presence of 107,891 HC genes from CS RefSeq v1.1

548    in the Renan genome sequence. We extracted the 461,476 individual exons larger than 30

549    bps and without Ns from this dataset and computed exon-by-exon BLAST in order to avoid

550    spurious sliced alignments. An exon was considered present if it matched the Renan

551    scaffolds with at least 95% identity over at least 95% of its length. To estimate the proportion

552    of identical exons between CS and Renan and the average nucleotide identity, we used the

553    same BLAST-based procedure but while restricting the dataset to 454,008 CS exons that

554    are on pseudomolecules (excluding chrUn) and considering Renan pseudomolecules

555    instead of scaffolds i.e., only exons carried by the same chromosome in CS and Renan were

556    considered. We extracted all available ISBPs (150 bps each) from the CS RefSeq v1.1 and

557    filtered out ISBPs containing Ns and those that do not map uniquely on the CS genome. This

558    led to the design of a dataset containing 5,394,172 ISBPs which were aligned on the Renan

559    scaffolds using BLAST. We considered an ISBP was conserved in Renan if it matched with

560    at least 90% identity over 90% of its length. We used the same ISBP dataset to study the

561    impact of polishing on error rate in the assembly while using BLAST and considering at least

562    90% identity over at least 145 aligned nucleotides.

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

563

564    We guided the construction of 21 Renan pseudomolecules based on collinearity with the CS

565    RefSeq Assembly v2.1. For this, we used the positions of conserved ISBPs as anchors

566    (5,087,711 ISBPs matching with >=80% identity over >=90% query overlap). This

567    represented 357 ISBPs/Mb, meaning that even the smallest scaffolds (30kb) carried

568    generally more than 10 potential anchors. However, some ISBPs match at non-orthologous

569    positions which create noise to precisely determine the order and orientation of some

570    scaffolds. To overcome this issue, we considered ISBPs by pairs. Only pairs of adjacent

571    ISBPs (i.e. separated by less than 50kb on both CS and Renan genomes) were kept as valid

572    anchors, allowing the filtering out of isolated mis-mapped ISBPs. Only scaffolds harboring at

573    least 50% of valid ISBP pairs on a single chromosome were kept. The others were

574    considered unanchored and they comprised the "chrUn". We calculated the median position

575    of matching ISBP pairs along each CS chromosome for defining the order of the Renan

576    scaffolds relative to each other. Their orientation was retrieved from the orientation of all

577    matching ISBP pairs in CS following the majority rule. We thus built 21 pseudomolecules

578    that were then corrected according to the HiC map as explained hereafter.

579    Two Hi-C biological replicates were prepared from ten-days plantlets of *Triticum aestivum*

580    cv. Renan following the Arima Hi-C protocol (Arima Hi-C User Guide for Plant Tissues DOC

581    A160106 v01). For each replicate, two libraries were constructed using the Kapa Hyper Prep

582    kit (Roche) according to Arima's recommendation (Library Preparation using KAPA Hyper

583    Prep Kit DOC A160108 v01). The technical replicates were then pooled and sent to Genewiz

584    for sequencing on an Illumina HiSeq4000 (four lanes in total), reaching a 35x coverage. We

585    mapped a sample of 240 million read pairs with BWA-MEM (Burrows-Wheeler Aligner, Heng

586    Li, 2013) to the formerly built 21 pseudomolecules, filtered out for low quality, sorted, and

587    deduplicated using the Juicer pipeline[44]. We produced a Hi-C map from the Juicer output

588    by the candidate assembly visualizer mode of 3D-DNA pipeline[45] and visualized it with the

589    Juicebox Assembly Tools software. Based on abnormal frequency contacts signals revealing

590    a lack of contiguity, scaffold-level modifications of order, orientation and/or chimeric scaffolds

591    were identified in order to improve the assembly. In case of chimeric scaffolds, coordinates

592    of resulting fragments were retrieved from the Juicebox Assembly Tools application but then

593    recalculated to correspond precisely to the closest gap in the scaffold. Pseudomolecules

594    were eventually rebuilt from initial scaffolds and new fragments while adding 100N gaps

595    between neighbor scaffolds. A final Hi-C map was built to validate the accuracy of the final

596    assembly.

597    **Calculation of chromosome coverage**

598    Short (*Triticum aestivum* cv Renan and *Ae. ventricosa*) and long-reads (*Triticum aestivum* cv

599    Renan) were aligned using minimap2 (with the following parameters '-I 17G -2 --sam-hit-only

600    -a -x sr' and '-I 17G -2 --sam-hit-only --secondary=no -a -x map-ont' respectively).

601    Coverage of individual chromosomes was calculated in 1 Mb windows using mosdepth[46]

602    (version 0.3.1) and the following parameters '--by 1000000 -n -i 2 -Q 10 -m'. Note that the '-i

603    2' and '-Q 10' parameters were used to keep only alignments of reads that mapped in a

604    proper pair and with a minimal quality value of 10. Coverage of individual chromosomes was

605    plotted in Figure 1. In addition, large deletions and duplications were detected using

606    CNVnator[47] with the Illumina bam file and a window of 100bp. We focused on large events

607    (>500kb) and detected only 15 deletions and no duplication (Figure 1).

608    **Transposable elements annotation**

609    Transposable elements were annotated using CLARITE[28]. Briefly, TEs were identified

610    through a similarity search approach based on the ClariTeRep curated databank of repeated

611     elements using RepeatMasker (www.repeatmasker.org) and modelled with the CLARITE

612     program that was developed to resolve overlapping predictions, merge adjacent fragments

613     into a single element when necessary, and identify patterns of nested insertions[28].


614     **Gene prediction**

615     We used MAGATT pipeline (Marker Assisted Gene Annotation Transfer for Triticeae,

616     https://forgemia.inra.fr/umr-gdec/magatt) to map the full set of 106,801 High Confidence and

617     159,848 Low Confidence genes predicted in Chinese Spring IWGSC RefSeq v2.1. The

618     workflow implemented in this pipeline was described in Zhu et al.[22]. Briefly, it uses gene

619     flanking ISBP markers in order to determine an interval that is predicted to contain the gene

620     before homology-based annotation transfer, limiting problems due to multiple mapping.

621     When the interval is identified, MAGATT uses BLAT[48] to align the gene (UTRs, exons, and

622     introns) sequence and recalculate all sub-features coordinates if the alignment is full-length

623     and without indels. If the alignment is partial or contains indels, it runs GMAP[49] to perform

624     spliced alignment of the candidate CDS inside the interval. If no ISBP-flanked interval was

625     determined or if both BLAT and GMAP failed to transfer the gene, MAGATT runs GMAP

626     against the whole genome, including the unanchored fraction of the Renan assembly. We

627     kept the best hit considering a minimum identity of 70% and a minimum coverage of 70%,

628     with *cross_species* parameter enabled.

629     We then masked the genome sequence based on mapped genes and predicted

630     transposable elements coordinates using BEDTools[50] mergeBed and maskfasta v2.27.1.

631     Hence, we computed a *de novo* gene prediction on the unannotated part of the genome. We

632     used TriAnnot[30] to call genes based on a combination of evidence: RNA-Seq data, *de*

633     *novo* predictions of gene finders (FGeneSH, Augustus), similarity with known proteins in

634     *Poaceae*, as described previously[7]. For that purpose, we mapped RNA-Seq reads with

635     hisat2[51] v2.0.5, called 277,505 transcripts with StringTie[52] v2.0.3, extracted their

636     sequences with Cufflink[53] gffread v2.2.1, and provided this resource as input to TriAnnot.

637     We optimized TriAnnot workflow to ensure a flawless use on a cloud-based hpc cluster (10

638     nodes with 32 CPUs/128GB RAM each and shared file system) using the IaaS Openstack

639     infrastructure from the UCA Mesocentre. Gene models were then filtered as follows: we

640     discarded gene models that shared strong identity (>=92% identity, >=95% query coverage)

641     with an unannotated region of the Chinese Spring RefSeq v2.1, considered as doubtful

642     predictions. We then kept all predictions that matched RNASeq-derived transcripts (>=99%

643     identity, >=70% query and subject coverage). For those that did not show evidence of

644     transcription, we kept gene models sharing protein similarity (>=40% identity, >=50% query

645 and subject coverage) with a *Poaceae* protein having a putative function (filtering out based
646 on terms "unknown", "uncharacterized", and "predicted protein").

## Comparison of genome assemblies

648 Genome assemblies were downloaded from https://webblast.ipk-gatersleben.de/downloads.
649 Contigs were extracted by splitting input sequences at each N and standard metrics were
650 computed. Gene completion metrics were calculated using BUSCO v5.0 and version 10 of
651 the poales geneset which contains 4896 genes.
652 We built dotplots between Renan, CS and 10 other reference quality genomes (Arina*LrFor*,
653 CDC Landmark, CDC Stanley, Jagger, Julius, LongReach Lancer, Mace, Norin61, SY
654 Mattis, spelta PI190962) by using orthologous positions of conserved ISBPs (1 ISBP every
655 2.5kb on average) identified by mapping them with BWA-MEM (maximum 2 mismatches,
656 100% coverage and minimal mapping quality of 30).

## Characterisation of haplotypic blocks

658 First a colored de Bruijn graph was built for each chromosome from the eleven available
659 chromosome-scale assemblies of wheat (Renan, CS, Arina*LrFor*, CDC Landmark, CDC
660 Stanley, Jagger, Julius, LongReach Lancer, Mace, Norin61 and SY Mattis). The colored de
661 Bruijn graph was created using Bifrost[54] with 31-mers and a unique color for each wheat
662 cultivar. In a second step, we extracted short markers (1kb) evenly spaced (20kb or 5kb) on
663 each chromosome and queried the colored de Bruijn graph using Bifrost and the following
664 parameter '-e 0.95' (for the comparison of each chromosome) and '-e 0.97' (for the
665 comparison of the *Sm1* locus). This parameter is the ratio of k-mers from queries that must
666 occur in the graph to be reported as present. For whole chromosome analyses, the 20kb
667 blocks were merged into 1-Mb blocks (the most abundant colour in the 50 20kb blocks was
668 retained for the 1Mb block). Individual blocks and *Ae. ventricosa* coverage were displayed
669 using RIdeogram[55].

## Comparison of a storage protein coding gene cluster

671 We performed manual curation of the gene models encoding storage proteins predicted in
672 Renan. Protein sequences of prolamin and resistance genes[35] from a 1B chromosome
673 locus were downloaded and aligned to the CS and Renan genomes using BLAT[48] with
674 default parameters. Draft alignments were refined by aligning the given protein sequence
675 and the genomic region defined by the blat alignment using Genewise with default
676 parameters. Resulting alignments were filtered in order to conserve only the best match for
677 each position by keeping only the highest-scoring alignment and the genomic region

678    containing the gene cluster was extracted. Then, we used the jcvi suite[56] with the mcscan

679    pipeline to find synteny blocks between both genomes. First, we used the

680    "jcvi.compara.catalog" command to find orthologs and then the "jcvi.compara.synteny

681    mcscan" with "--iter=1" command to extract synteny blocks. Finally, we generated the figure

682    with the "jcvi.graphics.synteny" command and manually edited the generated svg file to

683    improve the quality of the resulting image by changing gene colors, incorporating gaps and

684    renaming genes. Moreover, to make the figure clearer, we artificially reduced the intergenic

685    space by 95% so that gene structures appear bigger. The omega gene cluster

686    representation figure was generated by using DnaFeaturesViewer[57] with coordinates of

687    features generated by the mcscan pipeline used previously.

688

## 689    Additional files

690    All the supporting data are included in two additional files: (a) A supplementary file which

691    contains Supplementary Tables 1-7 and Supplementary Figures 1-3; (b) A supplementary

692    file which contains dotplots of the 21 chromosomes of Renan with other wheat genome

693    assemblies.

694

## 695    Acknowledgements

## 706    Availability of supporting data

707    The Illumina and PromethION sequencing data and the Bionano optical map are available in

708    the European Nucleotide Archive under the following project PRJEB49351. The genome

709    assembly and gene predictions are freely available from the Genoscope website

710    http://www.genoscope.cns.fr/plants/.

711 Additionally, all the data and scripts used to produce the main figures are available on a

712 github repository https://github.com/institut-de-genomique/Renan-associated-data

# Competing interests

714 The authors declare that they have no competing interests. JMA received travel and

715 accommodation expenses to speak at Oxford Nanopore Technologies conferences. JMA

716 and CB received accommodation expenses to speak at Bionano Genomics user meetings.

# Funding

# Author's contributions

722 SA, ID and AB extracted the sequenced DNA and generated the optical map. KL and AA

723 optimized and performed the nanopore and Illumina sequencing. NP, EP and MR generated

724 the Hi-C libraries and sequences. JMA, SE, BI, CM, PLZ, CB, HR, PL, DG and FC

725 performed the bioinformatic analyses. JMA, SE, BI, CM, PLZ, CB, CC, HR, PL and FC wrote

726 the article. JMA, PW and FC supervised the study.

727 **Table 1:** Comparison of *Triticum aestivum L.* genome assemblies. *NG50 and NG90 were
728 calculated using a genome size of 15Gb.

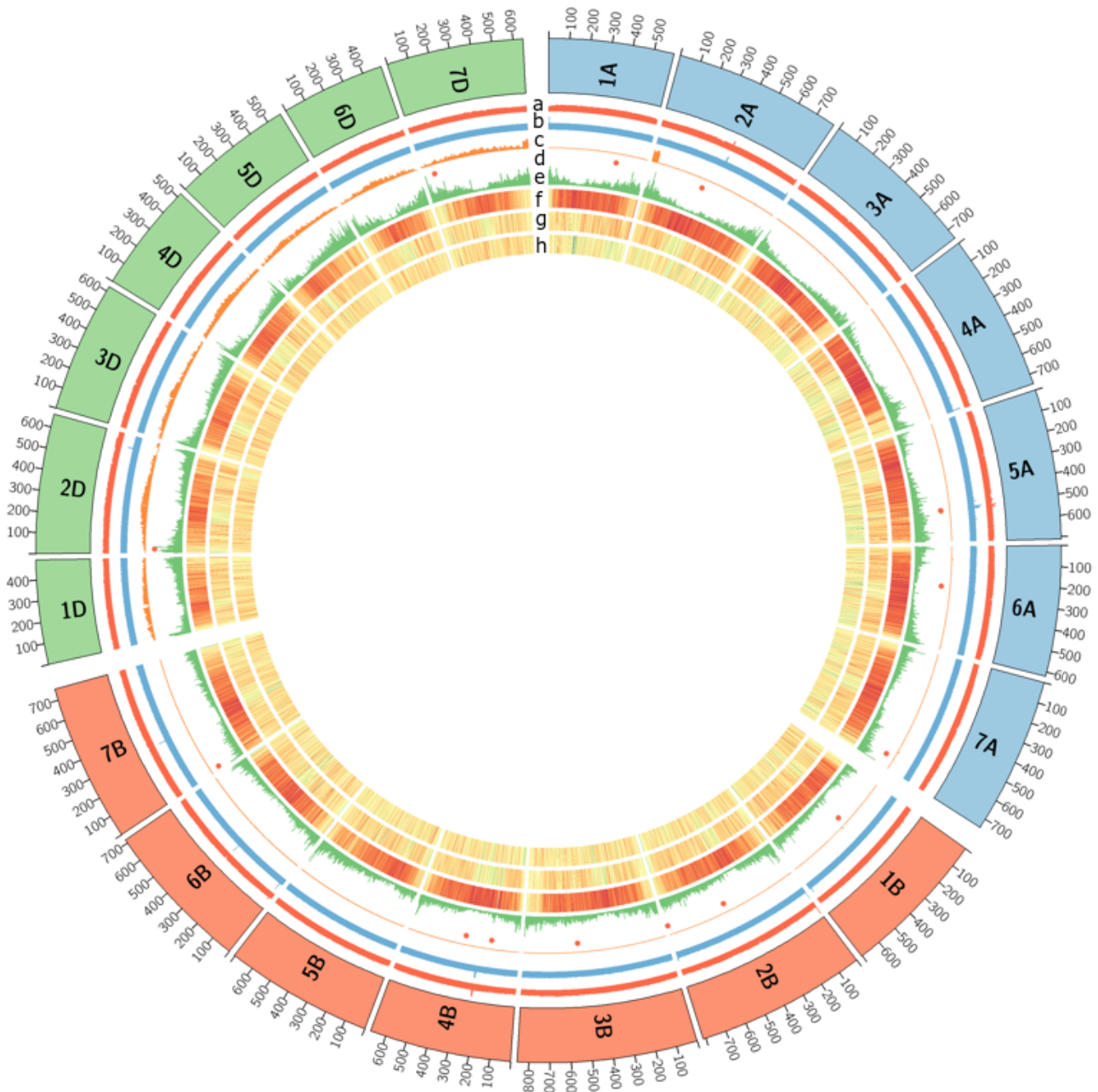|  |  | Renan<br>This study | Chinese Spring<br>RefSeq_v2.1<br>from Zhu et al.[22] |
|---|---|---|---|
| Number of contigs |  | 12,982 | 306,746 |
| Cumulative size (bp) |  | 14,001,122,256 | 14,317,423,665 |
| N50 (bp)<br>L50 |  | 2,159,703<br>1,958 | 341,062<br>12,223 |
| N90 (bp)<br>L90 |  | 598,285<br>6,645 | 32,302<br>59,261 |
| NG50* (bp)<br>LG50 |  | 1,973,000<br>2,202 | 322,161<br>13,254 |
| NG90* (bp)<br>LG90 |  | 264,272<br>8,816 | 16,550<br>85,688 |
| Longest contig (bp) |  | 15,116,687 | 3,528,546 |
| Number of chromosomes |  | 21 | 21 |
| Cumulative size (bp) |  | 14,195,643,615 | 14,225,829,371 |
| N50 (bp)<br>L50 |  | 703,299,328<br>10 | 713,360,512<br>10 |
| N90 (bp)<br>L90 |  | 520,815,552<br>19 | 518,332,608<br>19 |
| Longest (bp) |  | 854,463,248 | 851,934,019 |
| % of N |  | 1.78% | 1.52% |
| BUSCO on assemblies<br>(N=4,896) | Complete | 99.1% | 99.3% |
|  | Duplicated | 94.7% | 96.1% |
|  | Fragmented | 0.1% | 0.1% |
|  | Missing | 0.8% | 0.6% |
| Base accuray - Quality Value (kmer) |  | 32.8 | 44.5 |
| Number of genes |  | 109,552 | 107,891 |
| Average number of exons |  | 5.10 | 5.33 |
| BUSCO on gene predictions<br>(N=4,896) | Complete | 99.1% | 99.5% |
|  | Duplicated | 94.6% | 98.2% |
|  | Fragmented | 0.2% | 0.1% |
|  | Missing | 0.7% | 0.4% |

729 **Table 2:** TE classes proportions in Chinese Spring and Renan genome assemblies.
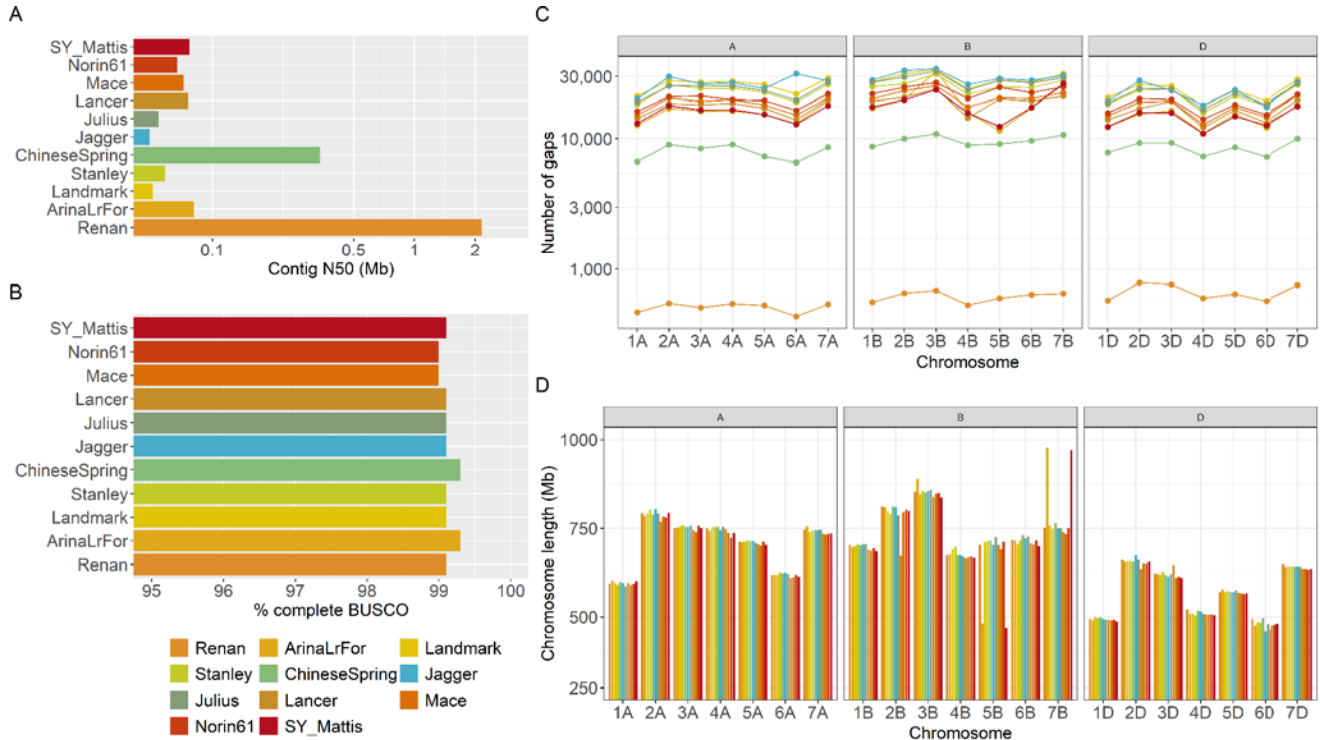
730

| | Chinese Spring RefSeq_v1.0 from Zhu et al. [22] | Chinese Spring RefSeq_v2.1 from Zhu et al.[22] | Renan RefSeq_v2.0 |
|---|---|---|---|
| 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 (%) | 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 |

731

732 **Figure 1.** Genome overview of the 21 chromosomes of hexaploid *T. aestivum* Renan (the 7
733 A chromosomes are in blue, the 7 B chromosomes in orange and the 7 D chromosomes in
734 green). From inner to outer track: (a) Coverage with short reads, (b) Coverage with long
735 reads, (c) coverage with *Ae. ventricosa* short reads, (d) Red dots represent large deletions
736 (>500Kb), (e) Gene density, (f) Density of CACTA (DNA transposon) elements, (g) Density
737 of Copia elements,  (h) Density of Gypsy elements. All densities and coverage are calculated
738 in 1-Mb windows; yellow and red colors in density plots indicate lower and higher values,
739 respectively.



740

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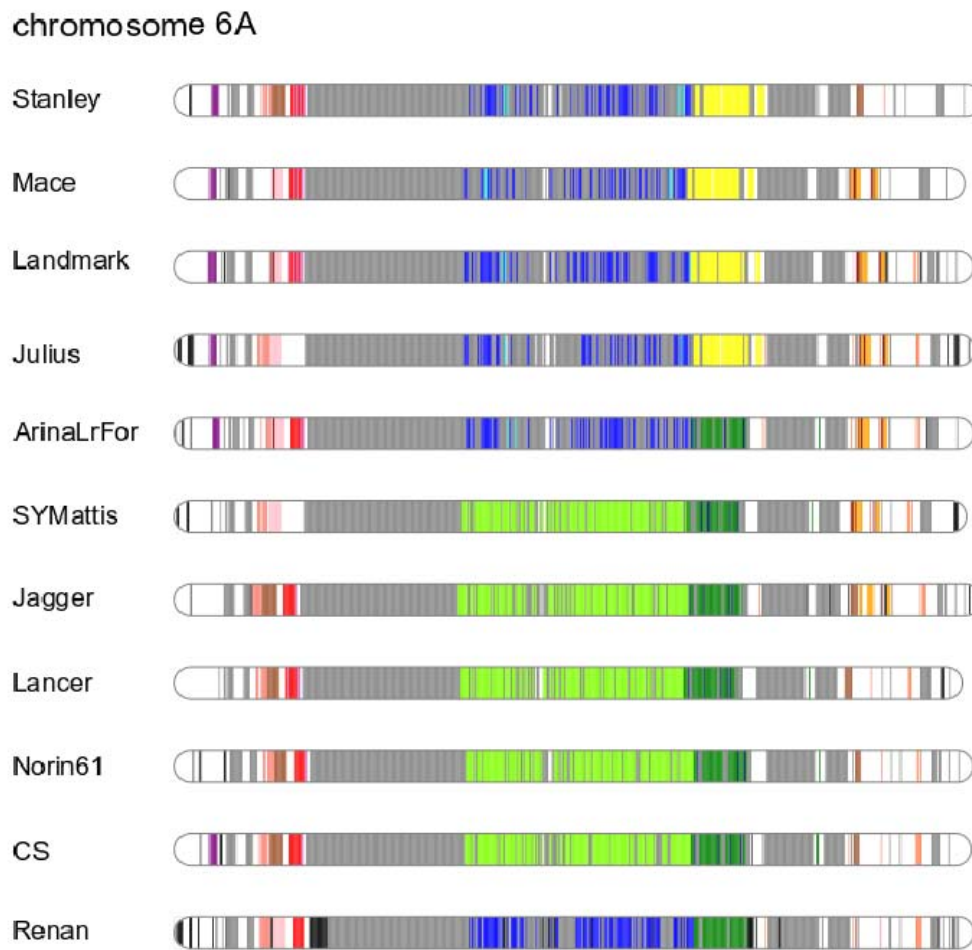


744

745 **Figure 3.** Dotplot comparisons of the 21 chromosomes of Renan (y axis) with the Chinese
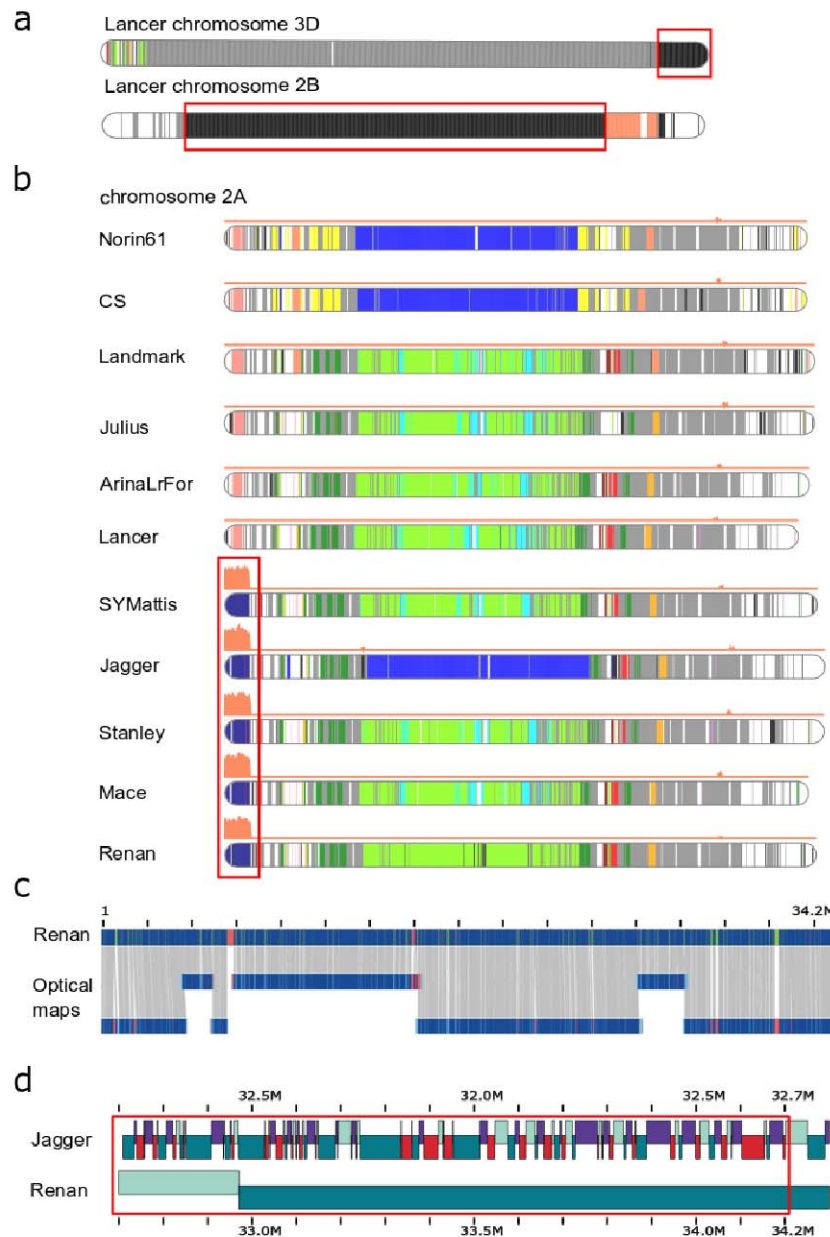746 Spring RefSeq v2.1 assembly (x axis).

747

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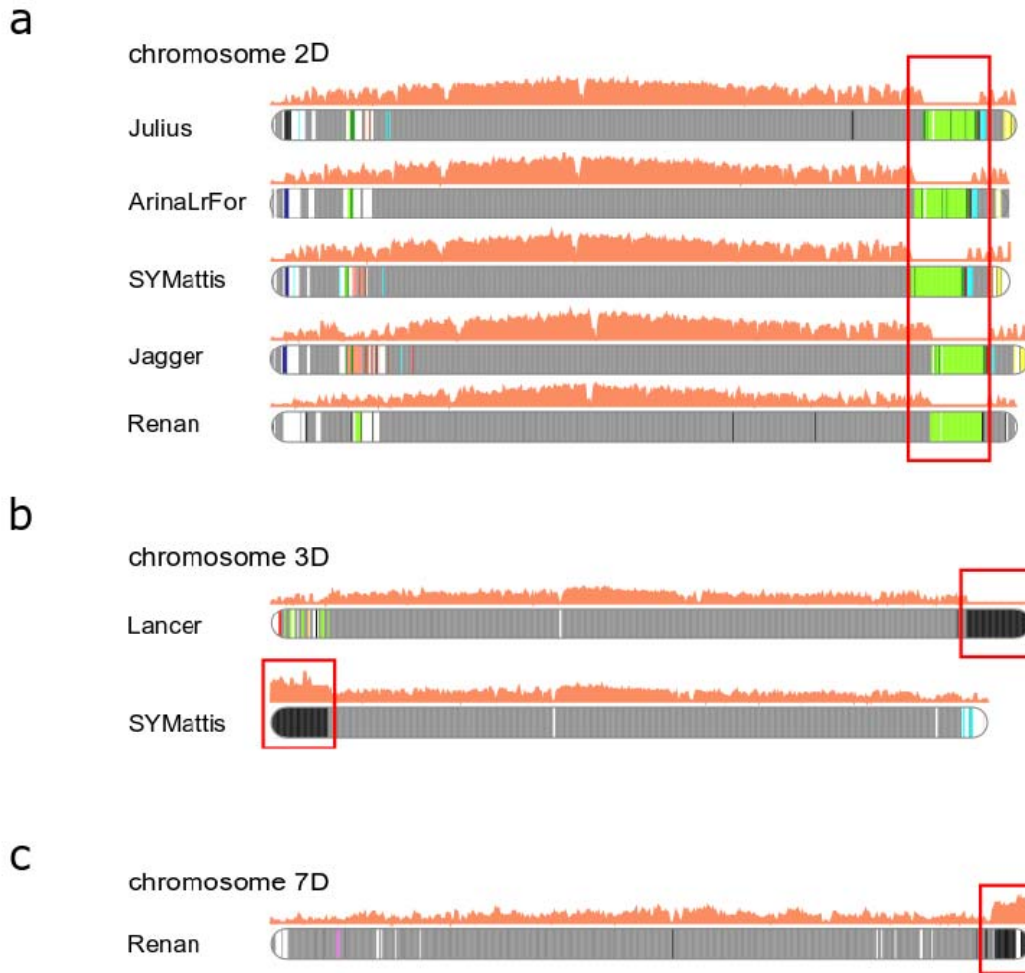752    or specific to a given cultivar.

753

754 **Figure 5.** Haplotypic blocks in wheat chromosomes. Colors represent common regions in
755 wheat cultivars. The gray and black regions represent haplotypes respectively shared by at
756 least 10 cultivars or specific to a given cultivar. The orange curve, when present, represents
757 coverage with *Ae. ventricosa* short reads. The red boxes frame the introgressions. **a.** Known
758 introgressions in chromosomes 3D and 2B in Lancer. Regions in black represent genomic
759 regions that are specific to Lancer and are respectively *T. ponticum* and *T. timopheevii*
760 introgressions as described previously[8]. **b.** *Ae. ventricosa* introgression on chromosome
761 3D in Stanley, Mace, SY Mattis and Jagger. This known introgression is also present in
762 Renan. The dark blue block represents the region shared across the five cultivars. **c.**
763 Validation of the introgression in Renan (chromosome 2A from 1 to 34.2Mb) using Bionano
764 maps. **d.** Comparison of the contig composition of the first megabases from the introgression

765 point                                                                                          in Jagger
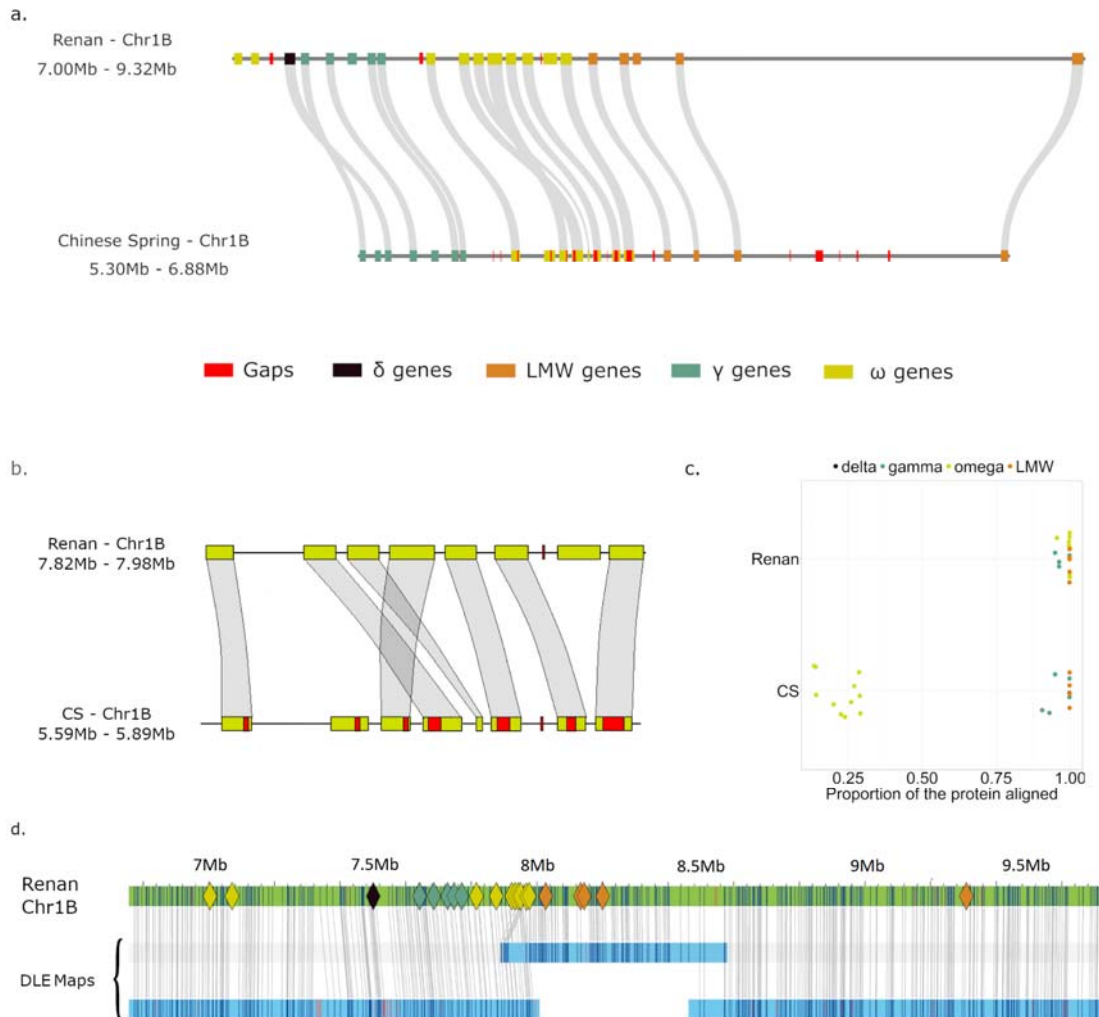766 and                                                                                            Renan
                                                                                                cultivars.

768 **Figure 6.** Haplotypic blocks in wheat chromosomes. Colors represent common regions in
769 wheat cultivars. The gray and black regions represent haplotypes respectively shared by at
770 least 10 cultivars or specific to a given cultivar. The orange curve represents coverage with
771 *Ae. ventricosa* short reads. The red boxes frame the introgressions. **a.** Candidate
772 introgression (green block) on chromosomes 2D in Julius, ArinaLrFor, SY Mattis, Jagger and
773 Renan. **b.** Candidate introgressions (black blocks) on chromosome 3D in Lancer and SY
774 Mattis. **c.** *Ae. ventricosa* introgression (black block) on chromosome 7D in Renan.
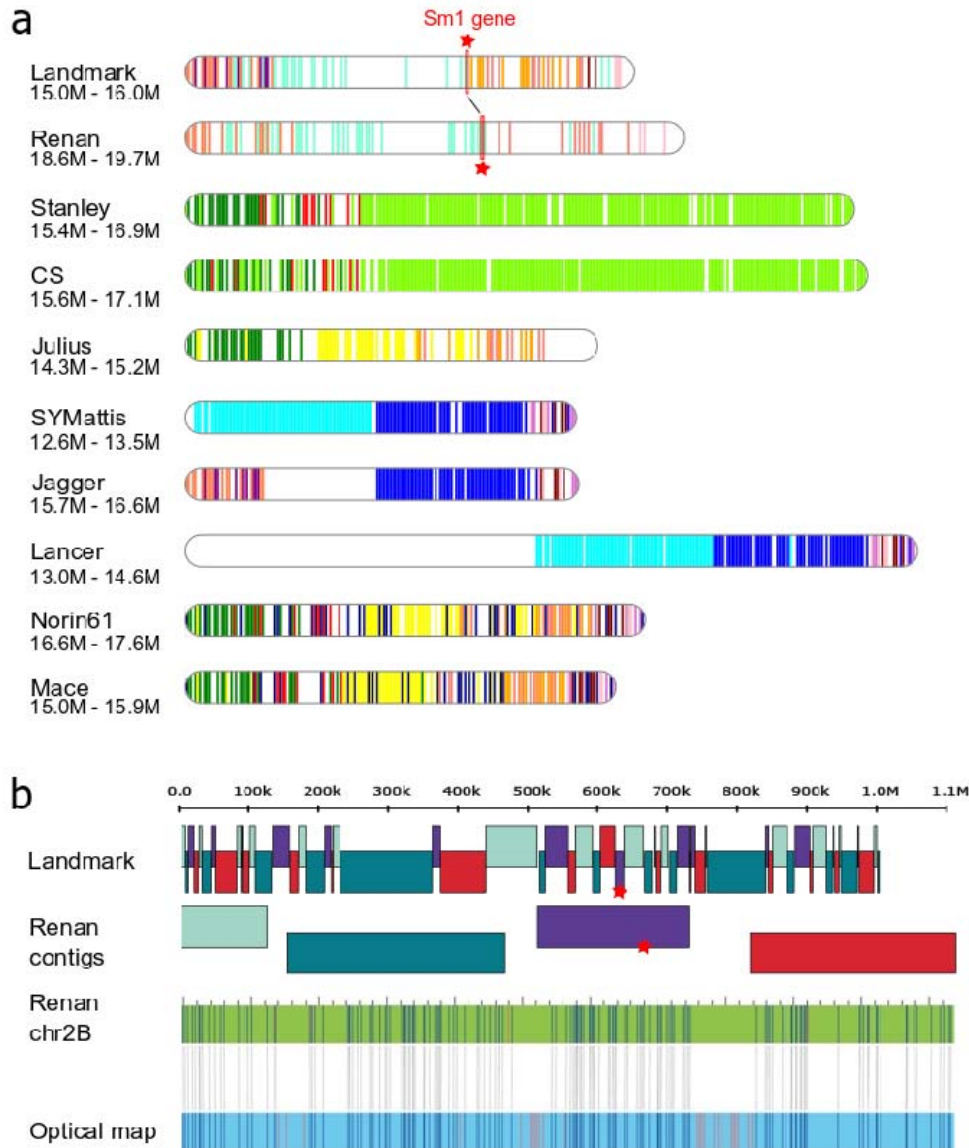


775

**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 two assemblies of Renan and CS. **b.** Zoomed view on the omega gliadin gene cluster **c.** Proportion of the length of the proteins that were aligned in the genomic region of Renan and CS. **d.** Alignment view of Bionano maps on the Renan cluster, colored diamond shapes represent genes belonging to the omega 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.

**Figure 8.** Comparison of the *Sm1* loci. **a.** Representation of haplotype blocks (5kb bins) of the region surrounding the *Sm1* gene on chromosome 2B. Colors represent common regions in wheat cultivars. The genomic region of Landmark (15Mb to 16Mb) was aligned against other cultivars to localize the *Sm1* loci. The *Sm1* gene in Landmark and Renan, the two *Sm1* carrier cultivars, is represented by a red star. **b.** Comparison of the contig composition in the *Sm1* region of 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..

# References

1. Dubcovsky J, Dvorak J. Genome Plasticity a Key Factor in the Success of Polyploid Wheat Under Domestication. *Science*. American Association for the Advancement of Science; 2007; doi: 10.1126/science.1143986.

2. Marcussen T, Sandve SR, Heier L, Spannagl M, Pfeifer M, International Wheat Genome Sequencing Consortium, et al.. Ancient hybridizations among the ancestral genomes of bread wheat. *Science*. 2014; doi: 10.1126/science.1250092.

3. Guan J, Garcia DF, Zhou Y, Appels R, Li A, Mao L. The Battle to Sequence the Bread Wheat Genome: A Tale of the Three Kingdoms. *Genomics Proteomics Bioinformatics*. 2020; doi: 10.1016/j.gpb.2019.09.005.

4. Chapman JA, Mascher M, Buluç A, Barry K, Georganas E, Session A, et al.. A whole-genome shotgun approach for assembling and anchoring the hexaploid bread wheat genome. *Genome Biol*. 2015; doi: 10.1186/s13059-015-0582-8.

5. Zimin AV, Puiu D, Hall R, Kingan S, Clavijo BJ, Salzberg SL. The first near-complete assembly of the hexaploid bread wheat genome, Triticum aestivum. *GigaScience*. 2017; doi: 10.1093/gigascience/gix097.

6. Clavijo BJ, Venturini L, Schudoma C, Accinelli GG, Kaithakottil G, Wright J, 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; doi: 10.1101/gr.217117.116.

7. Consortium (IWGSC) TIWGS, Appels R, Eversole K, Stein N, Feuillet C, Keller B, et al.. Shifting the limits in wheat research and breeding using a fully annotated reference genome. *Science*. American Association for the Advancement of Science; 2018; doi: 10.1126/science.aar7191.

8. Walkowiak S, Gao L, Monat C, Haberer G, Kassa MT, Brinton J, et al.. Multiple wheat genomes reveal global variation in modern breeding. *Nature*. 2020; doi: 10.1038/s41586-020-2961-x.

9. Miga KH, Koren S, Rhie A, Vollger MR, Gershman A, Bzikadze A, et al.. Telomere-to-telomere assembly of a complete human X chromosome. *Nature*. 2020; doi: 10.1038/s41586-020-2547-7.

10. Belser C, Istace B, Denis E, Dubarry M, Baurens F-C, Falentin C, et al.. Chromosome-scale assemblies of plant genomes using nanopore long reads and optical maps. *Nat Plants*. Nature Publishing Group; 2018; doi: 10.1038/s41477-018-0289-4.

11. Rousseau-Gueutin M, Belser C, Da Silva C, Richard G, Istace B, Cruaud C, et al.. Long-read assembly of the Brassica napus reference genome Darmor-bzh. *GigaScience*. 2020; doi: 10.1093/gigascience/giaa137.

12. Li G, Wang L, Yang J, He H, Jin H, Li X, et al.. A high-quality genome assembly highlights rye genomic characteristics and agronomically important genes. *Nat Genet*. 2021; doi: 10.1038/s41588-021-00808-z.

13. Liu J, Seetharam AS, Chougule K, Ou S, Swentowsky KW, Gent JI, et al.. Gapless assembly of maize chromosomes using long-read technologies. *Genome Biol*. 2020; doi: 10.1186/s13059-020-02029-9.

14. Tørresen OK, Star B, Mier P, Andrade-Navarro MA, Bateman A, Jarnot P, et al.. Tandem repeats lead to sequence assembly errors and impose multi-level challenges for genome and protein databases. *Nucleic Acids Res*. 2019; doi: 10.1093/nar/gkz841.

15. Li C, Xiang X, Huang Y, Zhou Y, An D, Dong J, et al.. Long-read sequencing reveals genomic structural variations that underlie creation of quality protein maize. *Nat Commun*. Nature Publishing Group; 2020; doi: 10.1038/s41467-019-14023-2.

16. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods*. Nature Publishing Group; 2020; doi: 10.1038/s41592-019-0669-3.

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

18. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using

849 repeat graphs. *Nat Biotechnol*. Nature Publishing Group; 2019; doi: 10.1038/s41587-019-
850 0072-8.
851 19. Vaser R, Sović I, Nagarajan N, Šikić M. Fast and accurate de novo genome assembly
852 from long uncorrected reads. *Genome Res*. 2017; doi: 10.1101/gr.214270.116.
853 20. Aury J-M, Istace B. Hapo-G, haplotype-aware polishing of genome assemblies with
854 accurate reads. *NAR Genomics Bioinforma*. 2021; doi: 10.1093/nargab/lqab034.
855 21. Istace B, Belser C, Aury J-M. BiSCoT: improving large eukaryotic genome assemblies
856 with optical maps. *PeerJ*. PeerJ Inc.; 2020; doi: 10.7717/peerj.10150.
857 22. Zhu T, Wang L, Rimbert H, Rodriguez JC, Deal KR, Oliveira RD, et al.. Optical maps
858 refine the bread wheat Triticum aestivum cv. Chinese Spring genome assembly. *Plant J*.
859 2021; doi: 10.1111/tpj.15289.
860 23. Rimbert H, Darrier B, Navarro J, Kitt J, Choulet F, Leveugle M, et al.. High throughput
861 SNP discovery and genotyping in hexaploid wheat. *PLOS ONE*. Public Library of Science;
862 2018; doi: 10.1371/journal.pone.0186329.
863 24. Istace B, Belser C, Falentin C, Labadie K, Boideau F, Deniot G, et al.. Sequencing and
864 Chromosome-Scale Assembly of Plant Genomes, Brassica rapa as a Use Case. *Biology*.
865 Multidisciplinary Digital Publishing Institute; 2021; doi: 10.3390/biology10080732.
866 25. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.
867 *J Mol Biol*. 1990; doi: 10.1016/S0022-2836(05)80360-2.
868 26. De Oliveira R, Rimbert H, Balfourier F, Kitt J, Dynomant E, Vrána J, et al.. Structural
869 Variations Affecting Genes and Transposable Elements of Chromosome 3B in Wheats.
870 *Front Genet*. Frontiers; 2020; doi: 10.3389/fgene.2020.00891.
871 27. Wick RR, Judd LM, Holt KE. Performance of neural network basecalling tools for Oxford
872 Nanopore sequencing. *Genome Biol*. 2019; doi: 10.1186/s13059-019-1727-y.
873 28. Daron J, Glover N, Pingault L, Theil S, Jamilloux V, Paux E, et al.. Organization and
874 evolution of transposable elements along the bread wheat chromosome 3B. *Genome Biol*.
875 2014; doi: 10.1186/s13059-014-0546-4.
876 29. Wicker T, Gundlach H, Spannagl M, Uauy C, Borrill P, Ramírez-González RH, et al..
877 Impact of transposable elements on genome structure and evolution in bread wheat.
878 *Genome Biol*. 2018; doi: 10.1186/s13059-018-1479-0.
879 30. Leroy P, Guilhot N, Sakai H, Bernard A, Choulet F, Theil S, et al.. TriAnnot: A Versatile
880 and High Performance Pipeline for the Automated Annotation of Plant Genomes. *Front Plant
881 Sci*. Frontiers; 2012; doi: 10.3389/fpls.2012.00005.
882 31. Brinton J, Ramirez-Gonzalez RH, Simmonds J, Wingen L, Orford S, Griffiths S, et al.. A
883 haplotype-led approach to increase the precision of wheat breeding. *Commun Biol*. 2020;
884 doi: 10.1038/s42003-020-01413-2.
885 32. Hao M, Zhang L, Ning S, Huang L, Yuan Z, Wu B, et al.. The Resurgence of
886 Introgression Breeding, as Exemplified in Wheat Improvement. *Front Plant Sci*. 2020; doi:
887 10.3389/fpls.2020.00252.
888 33. Kondrashov FA. Gene duplication as a mechanism of genomic adaptation to a changing
889 environment. *Proc R Soc B Biol Sci*. Royal Society; 2012; doi: 10.1098/rspb.2012.1108.
890 34. Panchy N, Lehti-Shiu M, Shiu S-H. Evolution of Gene Duplication in Plants. *Plant
891 Physiol*. 2016; doi: 10.1104/pp.16.00523.
892 35. Huo N, Zhang S, Zhu T, Dong L, Wang Y, Mohr T, et al.. Gene Duplication and Evolution
893 Dynamics in the Homeologous Regions Harboring Multiple Prolamin and Resistance Gene
894 Families in Hexaploid Wheat. *Front Plant Sci*. Frontiers; 2018; doi: 10.3389/fpls.2018.00673.
895 36. Xu J-H, Messing J. Organization of the prolamin gene family provides insight into the
896 evolution of the maize genome and gene duplications in grass species. *Proc Natl Acad Sci U
897 S A*. 2008; doi: 10.1073/pnas.0807026105.
898 37. Lang D, Zhang S, Ren P, Liang F, Sun Z, Meng G, et al.. Comparison of the two up-to-
899 date sequencing technologies for genome assembly: HiFi reads of Pacific Biosciences
900 Sequel II system and ultralong reads of Oxford Nanopore. *GigaScience*. 2020; doi:
901 10.1093/gigascience/giaa123.
902 38. Hon T, Mars K, Young G, Tsai Y-C, Karalius JW, Landolin JM, et al.. Highly accurate
903 long-read HiFi sequencing data for five complex genomes. *Sci Data*. 2020; doi:

904     10.1038/s41597-020-00743-4.

905     39. Belser C, Baurens F-C, Noel B, Martin G, Cruaud C, Istace B, et al.. Telomere-to-
906     telomere gapless chromosomes of banana using nanopore sequencing. *Commun Biol*.
907     2021; doi: 10.1038/s42003-021-02559-3.

908     40. Lv X, Chen Z, Lu Y, Yang Y. An End-to-end Oxford Nanopore Basecaller Using
909     Convolution-augmented Transformer. 2020 Nov.

910     41. Huang N, Nie F, Ni P, Luo F, Wang J. An attention-based neural network basecaller for
911     Oxford Nanopore sequencing data. *2019 IEEE Int Conf Bioinforma Biomed BIBM*.

912     42. Alberti A, Poulain J, Engelen S, Labadie K, Romac S, Ferrera I, et al.. Viral to metazoan
913     marine plankton nucleotide sequences from the Tara Oceans expedition. *Sci Data*. Nature
914     Publishing Group; 2017; doi: 10.1038/sdata.2017.93.

915     43. Rhie A, Walenz BP, Koren S, Phillippy AM. Merqury: reference-free quality,
916     completeness, and phasing assessment for genome assemblies. *Genome Biol*. 2020; doi:
917     10.1186/s13059-020-02134-9.

918     44. Durand NC, Shamim MS, Machol I, Rao SSP, Huntley MH, Lander ES, et al.. Juicer
919     Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Syst*.
920     2016; doi: 10.1016/j.cels.2016.07.002.

921     45. Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, et al.. De novo
922     assembly of the Aedes aegypti genome using Hi-C yields chromosome-length scaffolds.
923     *Science*. American Association for the Advancement of Science; 2017; doi:
924     10.1126/science.aal3327.

925     46. Pedersen BS, Quinlan AR. Mosdepth: quick coverage calculation for genomes and
926     exomes. *Bioinforma Oxf Engl*. 2018; doi: 10.1093/bioinformatics/btx699.

927     47. Abyzov A, Urban AE, Snyder M, Gerstein M. CNVnator: an approach to discover,
928     genotype, and characterize typical and atypical CNVs from family and population genome
929     sequencing. *Genome Res*. 2011; doi: 10.1101/gr.114876.110.

930     48. Kent WJ. BLAT—The BLAST-Like Alignment Tool. *Genome Res*. 2002; doi:
931     10.1101/gr.229202.

932     49. Wu TD, Watanabe CK. GMAP: a genomic mapping and alignment program for mRNA
933     and EST sequences. *Bioinformatics*. 2005; doi: 10.1093/bioinformatics/bti310.

934     50. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
935     features. *Bioinformatics*. 2010; doi: 10.1093/bioinformatics/btq033.

936     51. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment and
937     genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol*. 2019; doi: 10.1038/s41587-
938     019-0201-4.

939     52. Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL. StringTie
940     enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*.
941     2015; doi: 10.1038/nbt.3122.

942     53. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al..
943     Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and
944     isoform switching during cell differentiation. *Nat Biotechnol*. 2010; doi: 10.1038/nbt.1621.

945     54. Holley G, Melsted P. Bifrost: highly parallel construction and indexing of colored and
946     compacted de Bruijn graphs. *Genome Biol*. 2020; doi: 10.1186/s13059-020-02135-8.

947     55. Hao Z, Lv D, Ge Y, Shi J, Weijers D, Yu G, et al.. RIdeogram: drawing SVG graphics to
948     visualize and map genome-wide data on the idiograms. *PeerJ Comput Sci*. PeerJ Inc.; 2020;
949     doi: 10.7717/peerj-cs.251.

950     56. Tang H, Bowers JE, Wang X, Ming R, Alam M, Paterson AH. Synteny and Collinearity in
951     Plant Genomes. *Science*. American Association for the Advancement of Science; 2008; doi:
952     10.1126/science.1153917.

953     57. Zulkower V, Rosser S. DNA Features Viewer, a sequence annotations formatting and
954     plotting library for Python. *bioRxiv*. Cold Spring Harbor Laboratory; 2020; doi:
955     10.1101/2020.01.09.900589.

956