**SUPPLEMENTARY NOTE 1**

**ASSEMBLY METRICS**

3    We sequenced and assembled 15 wheat genomes. **Supplementary Table 1** provides passport
4    information for these lines. Seed stocks of the assembled lines are available at the UK Germplasm
5    Resources Unit (https://www.seedstor.ac.uk/).

6

7    **RQA Assemblies**

8

9    Paired-end, mate-pair, 10X Genomics Chromium linked reads, and Hi-C data were generated and
10   assembled using the DenovoMagic3.0 pipeline as described previously for other cereals (See Methods).
11   The assemblies were 14.2-14.8 Gb in length and consisted of 487,203-1,058,304 contigs and 99,465-
12   701,145 scaffolds (**Supplementary Table 2**). POPSEQ ([ftp://ftp.ipk-gatersleben.de/barley-popseq](ftp://ftp.ipk-gatersleben.de/barley-popseq)) and
13   Hi-C was then used to break chimeric scaffolds (corrected assembly), and then order and orient the
14   scaffolds into chromosome pseudomolecules that were ~14 Gb in length, as was done previously for
15   Chinese Spring RefSeq v1.0 (**Supplementary Table 2**). The assemblies contained ~98% of the BUSCO
16   gene complement (**Supplementary Table 2**) and were highly collinear to the assembly of Chinese Spring
17   (**Extended Data Fig. 1**), thereby validating their quality and completeness.

18

19   **Scaffolded Assemblies**

20

21   Paired-end and mate-pair libraries were constructed and sequenced at the Earlham Institute by the
22   Genomics Pipelines Group. A total of 2 μg of DNA was sheared targeting 1 kb fragments on a Covaris- S2,
23   size selected on a Sage Science Blue Pippin 1.5% cassette to remove DNA molecules <600 bp, and
24   amplification-free, paired-end libraries were constructed using the Kapa Biosciences Hyper Prep Kit.
25   Mate-pair libraries with insert sizes >7 kb were constructed from 9 μg of DNA based on the Illumina
26   Nextera mate-pair kits. Sequencing was performed on an Illumina HiSeq 2500 instrument, which
27   generated 250 bp reads at a read depth of ~45 for the pair-end and ~25 for the mate-paired sequencing
28   libraries.

29

30   A read depth between 44 and 51 was generated per line. Contig construction was performed using the
31   w2rap-contigger using k=200. Two mate-pair libraries were produced for each line except Weebill 1,
32   where five libraries were used. Mate-pairs were processed, filtered, and used to scaffold contigs as
33   described in the W2RAP pipeline (https://github.com/bioinfologics/w2rap). Scaffolds >500 bp were
34   removed from the final assemblies (**Supplementary Table 3**). The k-mer Analysis Toolkit
35   ([https://github.com/TGAC/KAT](https://github.com/TGAC/KAT)) was used to validate the scaffolds by generating a k-mer histogram
36   from the matrix of k-mers shared between the paired-end reads and the scaffolds.

37

38

| | **SUPPLEMENTARY NOTE 2** |
|---|---|

39  **SUPPLEMENTARY NOTE 2**
40  **OXFORD NANOPORE SEQUENCING OF CDC LANDMARK**
41
42  The order and orientation of the scaffolds in each of the pseudomolecules for the CDC Landmark RQA
43  were determined using Hi-C, as was performed previously for the published genomes of Chinese Spring,
44  Svevo (durum wheat), and Zavitan (wild emmer wheat). This method was used for pseudomolecule
45  assembly (**Supplementary Note 1**) and for calling large scale (>1 Mb) inversions (**Supplementary Table
46  13; Extended Data Fig. 8**) due to the large distances (several Mb) that can be captured by Hi-C contact
47  maps and high base-level accuracy of Illumina sequencing. To validate this method, we used long reads
48  generated from the Oxford Nanopore Technology (ONT) sequencing platform to identify reads spanning
49  the gaps between the neighbouring scaffolds, thereby confirming the scaffold placement by Hi-C. We
50  performed long read sequencing of CDC Landmark; high molecular weight DNA was isolated from wheat
51  seedlings and size selected as was performed for the 10x Genomics Chromium sequencing. Libraries for
52  whole-genome ONT sequencing were prepared following the standard ligation protocol (LSK109) and
53  sequenced on a GridION instrument using standard parameters and R9 flow cells. The sequencing
54  generated ~72 million reads with a read N50 of ~15 kb, which is equivalent to a ~30 fold sequencing
55  read-depth of CDC Landmark. The sequence reads were aligned to the RQA of CDC Landmark and
56  Chinese Spring using Minimap2 v.2.1.0. Next, we extracted the read IDs for only those ONT long reads
57  that mapped within a 5 kb region at both ends of each scaffold.  We then compared read IDs between
58  scaffold ends to identify those reads that mapped to both neighbouring scaffolds placed in the
59  pseudomolecules, as well as to the ends of scaffolds from other locations in the RQA. The results
60  demonstrate that adjacent scaffolds were 94 times more likely to share read IDs, thus corroborating the
61  scaffold placement by Hi-C. The overlap in read IDs between scaffold junctions was used to construct a
62  scaffold-scaffold contact map across the ordered and oriented scaffolds in the CDC Landmark
63  pseudomolecules (**Extended Data Fig. 2a**).
64
65  We also manually inspected the breakpoints of structural variants that we identified from aligning the
66  Hi-C data from CDC Landmark to the Chinese Spring RQA (**Supplementary Table 13**). For example,
67  chromosome alignments and Hi-C predicted inversions on chromosomes 2A, 3B, and 3D when
68  comparing CDC Landmark to Chinese Spring (**Extended Data Fig. 2b**). Manual inspection of the ONT long
69  reads mapped to the assembly of Chinese Spring identified reads that whose alignment were
70  interrupted at the breakpoints of this inversion event, which mapped to the (+) and (-) strand of the
71  assembly on either side of the breakpoint, thereby supporting the inversion events between assemblies
72  (**Extended Data Fig. 2b**).
73
74  Furthermore, we validated the gene structure of *Sm1* using the ONT read mappings, given that CDC
75  Landmark is a carrier of that gene (**Fig. 3**). Inspection of the region containing the candidate gene for
76  *Sm1* on chromosome 2B of CDC Landmark also confirmed a uniform distribution of read depth, as well
77  as many long reads that were able to completely span the candidate gene (**Extended Data Fig. 10b**).
78  However, alignment of CDC Landmark long reads to corresponding region in the Chinese Spring
79  assembly that only contains a partial gene sequence revealed irregular patterns of read depth, as well as
80  an abundance of point mutations, thereby supporting an alternate haplotype (**Extended Data Fig. 10b**).
81
82  In summary, we used the established RQA assembly pipeline that was used previously used for Chinese
83  Spring, Svevo, and Zavitan, to generate new RQA for wheat lines of interest for global breeding
84  programs. Using ONT long reads, we demonstrate that this assembly approach is able to properly order
85  and orient scaffolds in pseudomolecules and identify structural variation between assemblies. We also
86  demonstrate that the RQA assembly pipeline was able to properly assemble the gene candidate for *Sm1*.

**SUPPLEMENTARY NOTE 3**

**VARIATION IN GENES**

90   **Gene Projections**

91   We used the previously published RefSeq v1.0 high-confidence gene models for Chinese Spring to assess
92   the gene content in each assembly (See Online Methods). Projections were filtered using several
93   iterations to ensure the best projected match was used. The first iteration required a contiguous ORF
94   and informant coverage ≥ 99% on an orthologous chromosome, the second removed the chromosome
95   orthology requirement, the third iteration relaxed the informant coverage to equal or greater than 90%,
96   and the fourth step allowed disrupted/non-contiguous ORFs restricting informant coverage to ≥ 95%.
97   For the fifth iteration, we projected all remaining matches with a contiguous ORF comprising a start and
98   stop codon, an informant coverage ≥ 5% and a genome-wide mapping frequency <50 copies. In a last
99   step, we projected all matches from informants missed in previous iterations if their alignment coverage
100  exceeded 80%. To provide uniformity in the data, Informants were also projected onto the Chinese
101  Spring assembly. This involved first masking existing gene coding regions and iteratively adding genes
102  with the first iteration including high coverage (≥ 99%) matches, then adding matches with a contiguous
103  ORF and a mapping frequency ≤ 10, and then genes with up to 50 matches per genome.

104

105  We identified tandem gene clusters and OGs (orthologous groups; See Online Methods) to quantify
106  unique genes and those displaying PAV and CNV patterns. Next, we surveyed CDS similarities between
107  all pairwise combinations of cultivars using BLATn v3.5. Next, we generated an undirected graph with
108  gene accessions as nodes and edges representing matches between two genes and identified OGs and
109  tandem clusters as subgraphs/connected components restricted to particular edge attributes. OGs were
110  restricted to RBH relations, while tandem clusters required an e-value ≤ $1^{-30}$, an alignment coverage of at
111  least 50% for both, the query and database accession, and a maximal genomic distance of nine
112  unrelated genes between both accessions. PAV genes per line were determined from the orthologous
113  groups limited to one-to-one relations and comprising 2-10 cultivars. An absence variant was counted
114  for each group and for each line if a group missed an orthologous gene of the respective line. Hence, the
115  reported PAV count (**Supplementary Table 5**) is the sum of the number of genes per line that are
116  missing in the union of all orthologous groups. In contrast, computation of CNV genes per line has been
117  restricted to orthotandem clusters (i.e. the orthologous groups above complemented by tandem
118  duplications) that comprised at least one gene of the respective line. Thereby, CNV counts ignore any
119  cluster which misses the respective line and rather describe the expansion of copies. Expansion is
120  defined as the number of gene copies per line which are redundant to the minimum copy number that
121  each of the lines contribute to the cluster. We limit our analysis to the one-to-one groups and their
122  extension to orthotandem clusters to ensure high accuracy. It should be emphasized that the group
123  'other' comprises genes that are either singletons in a line or participate in complex relationships where
124  at least one line contributes more than one copy to the orthologous group. Albeit in the latter cases
125  there are certainly copy number variations, ambiguous orthologous assignments will highly increase the
126  errors in PAV counts. Overall, it should be stressed that the PAV counts are likely inflated by the
127  orthologous groups with very low line counts (2-7) which trigger in each case for a large number of lines
128  an absence count. From our experience, many of these groups represent genes with questionable or
129  spurious functionality.

130

131  **Genetic Diversity of A, B and D Homeologs Suggests Increased Targets of Selection**

132  To estimate the pattern of genome-wide polymorphisms of wheat cultivars, we analysed the coding
133  sequences by using gene alignments of the *de novo* genome assemblies. To derive codon-based

134    alignments, we aligned amino acid sequences of one OG (excluding gene models with disrupted ORFs;
135    see above) using MUSCLE v3.8 with default parameters and then back translated aligned codons to
136    nucleotide alignments. We also only used alignments where each homeolog was found in all RQA used
137    in the analysis. Average pairwise nucleotide diversity ($\pi$) and $\theta_W$ was estimated and their units are
138    reported as per base pair. By taking the mean values of $\pi$ for each subgenome, we found that the
139    genetic diversity of the A and B subgenomes (0.0012 and 0.0021, respectively) are higher than that of
140    the D subgenome (0.0004) (**Supplementary Table 6**). The values were highly similar in both datasets
141    with and without PI190962 (spelt wheat), although lower Tajima's *D* in the dataset with PI190962
142    indicates an excess of rare variants in PI190962. Because PI190962 may have a complex history involving
143    hybridization with wild species, we focused only on the analysis of just the bread wheat lines. The peak
144    $\pi$ distribution for the three subgenomes was similar to previous studies, which indicates that the RQA
145    were able to capture similar patterns of diversity as larger diversity panels representing global breeding
146    programs.

147

148    Synonymous (or silent site) nucleotide diversity ($\pi\_sil$) was approximately double the total site $\pi$ for
149    each subgenome ($\pi\_sil$ = 0.0024, 0.0046, and 0.0009 for A, B, and D subgenomes, respectively
150    (**Supplementary Table 7**). The level of polymorphism of the A and B subgenomes is similar to the natural
151    allopolyploid *Arabidopsis kamchatica* (0.0014-0.0015 in total sites, 0.0044-0.0049 in synonymous sites),
152    suggesting that bread wheat retained considerable global variation comparable to wild species, despite
153    domestication and polyploidization. The patterns are consistent with a recent report showing higher
154    variation in the A and B subgenomes than in the D subgenome. The mean and median of Tajima's *D* of
155    the D subgenome is lower than the A and B subgenomes (**Extended Data Fig. 3 a,b**), indicating an excess
156    of rare variants in the D subgenome.

157

158    Next, we tested whether homeologous pairs experienced similar or different evolutionary trajectories in
159    a genome-wide manner. In wheat, it is known that the level of polymorphism along each chromosome is
160    positively correlated with the distance from centromere, which should result in a positive correlation of
161    $\pi\_sil$ between homeologs. Yet, the correlations among homeologs were low (0.11-0.29) (**Supplementary
162    Table 8**), suggesting different evolutionary trajectories of homeologs. More importantly, the correlation
163    of the neutrality statistic, Tajima's *D*, was very low between pairs of homeologs (*r* = 0.02-0.06) among the
164    three subgenomes, which again supports that homeologous copies experienced different selective
165    pressure. These results are in line with our current understanding that selective sweeps rarely occur in
166    homeologous regions in bread wheat.

167
168

**SUPPLEMENTARY NOTE 4**
170 **SEQUENCE DIVERSITY OF *RF* LOCI IN WHEAT**

171

172 Despite the promise of higher yield and better stress resistance of hybrid varieties compared to
173 conventional lines, hybrid breeding in wheat remains underexploited. An efficient pollination control-
174 system that would be easy to apply on a commercial scale is missing. The application of cytoplasmic
175 male sterility (CMS) and restorer-of-fertility (*Rf*) genes have been successful in other plant species such
176 as rice and maize, but is difficult to use in wheat due to poor effectiveness of the known wheat *Rf* genes.
177 In most flowering plants, the majority of *Rf* genes belong to the pentatricopeptide repeat (PPR) family
178 encoding mitochondrial sequence-specific RNA binding proteins. The PPR family can be split into two
179 classes based on their motif architecture, and this distinction correlates with function: P-class PPR
180 proteins are implicated in a wide range of RNA processing activities whilst PLS-class PPR proteins are
181 almost exclusively implicated in RNA editing. *Rf* genes are found within a subclade (Rf-like, or RFL genes)
182 of the P-class.

183

184 **Discovery of RFL-mTERF Clade and its Expansion in Wheat**

185 Members of another family of sequence-specific organellar nucleic acid binding proteins, the
186 mitochondrial transcription termination factor (mTERF) family, may also act in fertility restoration in
187 plants. The wheat genome contains ~400 mTERF sequences of which more than 300 are found in
188 clusters overlapping with RFL gene clusters (**Supplementary Table 9**). For comparison, in *Arabidopsis*,
189 there are 35 mTERF genes, of which 25 are distributed evenly across the genome and are implicated in
190 RNA-associated processes in chloroplasts and mitochondria. Based on sequence similarity, these 25
191 genes are putatively orthologous to the ~28-30 scattered mTERF sequences found on each of the three
192 subgenomes in wheat (**Supplementary Table 9**). A single cluster of ~10 mTERF genes of unknown
193 function from *Arabidopsis* may correspond to the huge clusters of mTERF genes found in wheat and
194 other cereals including rye and barley. The discovery of clustered mTERF sequences that share genome
195 locations and the same patterns of evolution as *RFL*-type PPR genes is a strong indication that they play
196 a major and hitherto unappreciated role in fertility restoration in cereals.

197

198

199

**SUPPLEMENTARY NOTE 5**

201               *AEGILOPS VENTRICOSA* 2N$^V$S SEGMENT FROM *VPM-1*

202

203     Based on pedigree and marker analyses, it was previously known that the *Aegilops ventricosa* 2N$^v$S
204     segment from the *Vpm-1* introgression is present in several wheat cultivars including Jagger. We
205     generated a RQA of Jagger (**Supplementary Note 1**) and delineated the 2N$^v$S segment to be ~33 Mb
206     based on patterns of RLC Angela elements (**Extended Data Fig. 5**) and chromosome alignments to
207     Chinese Spring chromosome 2A (**Extended Data Fig. 6a**), which was further corroborated using genomic
208     *in situ* hybridization (GISH) technology with NN genome probes (results not shown). The pattern in *RLC
209     Angela* elements and chromosome alignments also revealed that in addition to Jagger, CDC Stanley, SY
210     Mattis, and Mace are also carriers of the same introgression. We also observed a region within the spelt
211     genome that was dissimilar to both chromosome 2A and the 2N$^v$S introgression, suggesting that
212     alternative haplotypes may exist in this region (**Extended Data Fig. 6a, arrow**). An alternative haplotype
213     is also supported by the analysis of unique RLC Angela elements in that region in the spelt genome
214     assembly (**Extended Data Fig. 5**).

215

216     We collected various tissue types from Jagger after being grown until different growth stages and
217     conducted RNA sequencing. Sequencing data was deposited to NCBI SRA (**Supplementary Table 17**). We
218     annotated the Jagger 2N$^v$S segment using a combination of *ab initio* predictions and RNA-seq evidence.
219     First, we mapped the ~3 billion RNA-seq reads against the Jagger genome assembly using STAR v2.6.0b.
220     Intron and exon structures were predicted based on RNA-seq alignments, which were combined with *ab
221     initio* gene prediction by AUGUSTUS v3.2.3. Second, we *de novo* assembled RNA-seq transcripts using
222     Trinity and mapped these back to the genome using GMAP (v 2017-06-20). EvidenceModeler v1.1.1 was
223     used to combine the *ab initio* predictions and mapped transcripts, resulting in a set of candidate gene
224     models. We further differentiated the candidates derived from EvidenceModeler into classes of bona-
225     fide genes, non-coding transcripts, pseudogenes and TEs. We then applied a confidence classification
226     protocol similar to the one applied to Chinese Spring RefSeqv1.0, based on coverage and hits in the
227     PTREP, UniPoa, and UniMag databases. The result was a set of high-confidence (HC) genes
228     (**Supplementary Table 18**). Finally, we assigned a functional annotation and human readable description
229     to the HC genes using AHRD v1.6 (https://github.com/groupschoof/AHRD) (**Supplementary Table 18**).
230     Analysis revealed several groups of related genes, including genes encoding disease resistance proteins
231     (i.e. NB-ARC and NLR genes), cytochrome P450s, transporters, chalcone synthases, glycosyltransferase,
232     sulfotransferase, and proteases. We identified orthologous genes between Jagger 2N$^v$S and Chinese
233     Spring chromosomes 2A, 2B, 2D based on RBH. Orthologous genes and genomic distributions were
234     visualized using Circos (**Extended Data Fig. 6b**).

235

236     We used genotyping-by-sequencing (GBS) to predict the presence of 2N$^v$S in three datasets: (1) Kansas
237     State University winter wheat, (2) USDA Regional Performance Nursery, and (3) International Maize and
238     Wheat Improvement Center (CIMMYT) spring wheat (**Supplementary Table 19**). The prediction was
239     based on relative count of wheat and alien specific GBS tag numbers. The frequency of 2N$^v$S carriers
240     increased in the three datasets (**Extended Data Fig. 6c**), reaching ~80% for CIMMYT and Kansas State
241     University breeding lines in recent years. These results suggest that 2N$^v$S carrier varieties were
242     collectively grown in tens of millions of hectares throughout the world. We also studied the relationship
243     between 2N$^v$S presence and wheat grain yield, our results suggest that the 2N$^v$S segment is providing a
244     yield benefit in majority of the years studied (**Extended Data Fig. 6d; Supplementary Table 20**). We
245     analysed the data by fitting a mixed linear model with the presence of 2N$^v$S as having fixed effects and

246    sites having random effects (lines were tested in ~20 locations each year). The percentages of yield
247    benefit across the years appear to be larger compared to previously reported, likely due to different
248    environmental factors such as watering and disease pressures. The yield benefits appear to be stable
249    across years, under different breeding stages and across different performance tests in regional,
250    national (U.S.A), and global scales. The release of RQA for $2N^vS$ carriers with both spring and winter
251    growth habits, from multiple breeding programs and continents, provides new resources that can be
252    used to characterize this introgression. Additionally, our data suggests that this translocation is
253    increasing in frequency and is having an impact on wheat productivity.

254

255

256

**SUPPLEMENTARY NOTE 6**
**CYTOLOGICAL KARYOTYPING**

260   Mitotic metaphase chromosomes were prepared by the conventional acetocarmine-squash and non-
261   denaturing fluorescence *in situ* hybridization (ND-FISH) of three repetitive sequence probes was
262   performed using three probes: Oligo-pSc119.2-1 (Tamra-5'-CCGTT TTGTG GACTA TTACT CACCG CTTTG
263   GGGTC CCATA GCTAT-3'), Oligo-pTa535 (AlexaFluor488-5'-GACGA GAACT CATCT GTTAC ATGGG CACTT
264   CAATG TTTTT TAAAC TTATT TGAAC TCCA-3'), and Oligo-pTa713 (AlexaFluor647-5'-AGACG AGCAC GTGAC
265   ACCAT TCCCA CCCTG TCTTA GCGTA ACGCG AGTCG-3'). Polymorphisms detected with respect to the
266   karyotype of Chinese Spring were summarized in **Supplementary Table 24**. Overall, 44 polymorphisms
267   were detected in hexaploid wheat. A hierarchal clustering of the accessions based on the detected
268   polymorphisms indicated that the accessions were largely divided into clusters (**Extended Data Fig. 7c**)
269   that were in agreement with our other phylogenetic analyses (**Fig. 1a,b; Extended Data Fig. 3d**) The
270   most striking karyotypic difference to Chinese Spring was the translocation between chromosomes 5B
271   and 7B (**Fig. 2e-g**). The presence of polymorphic FISH signals on chromosome 5BS between SY Mattis
272   and Arina*LrFor* suggested that the translocated chromosomes may have different origins or diversified
273   after the translocation event (**Supplementary Table 24**).

275   Some of the structural rearrangements detected by sequence comparisons (**Extended Data Fig. 1**) and
276   Hi-C (**Supplementary Table 23**) are supported by the karyotyping (**Supplementary Table 24**) and
277   *LTR_Angela* analysis (**Supplementary Tables 13-16**). Namely, a large inversion detected on
278   chromosomes 4B was represented by absence of a block of pTa713 signal on the short arm in Arina*LrFor*,
279   CDC Landmark, CDC Stanley, LongReach Lancer, Robigus, and Paragon. Those without inversions were
280   Julius, Norin 61, Jagger, SY Mattis, Mace, Cadenza, Claire, and Weebill 1. An inversion in the distal region
281   of 3DL in CDC Landmark and CDC Stanley made the distance between two pTa535 signals larger than in
282   other accessions. Inversions at the distal end of 4AL may not be a simple inversion, but associated with
283   loss of FISH signals (loss of pSc119 signal in Norin 61, Julius, Claire, and SY Mattis; loss of pTa713 signal in
284   CDC Stanley). Only Arina*LrFor* was lacking a pTa535 signal at the distal end of chromosome 1AS, which
285   we detected as variable TEs and is consistent with local alignment and *RLC_Angela* analyses
286   (**Supplementary Table 14**). The pTa535 signal at the distal region of 3DL in LongReach Lancer shifted
287   proximal, likely due to insertion of *Th. ponticum* chromatin (**Fig. 2a**). Together, these cytological
288   observations provide experimental evidence to support the observed differences we observed between
289   the RQA.

**SUPPLEMENTARY NOTE 7**
**INHERITANCE OF THE 5B/7B TRANSLOCATION IN NORTH-WESTERN EUROPEAN WHEAT VARIETIES**

292

293   A panel of 538 north-west European wheat varieties was genotyped for the presence of the
294   chromosome 5B/7B translocation using representative marker *RAC875_c5965_1153* from the Illumina
295   90K wheat SNP array. These lines include almost the entire diversity of the UK wheat gene pool from the
296   1920s to the present day, together with their northern European ancestors and relatives. The
297   translocation is widespread, being found in 66% of all lines, ranging from 50% of varieties released
298   before 1970 to 73% of varieties released in the 2000s. Currently, 66% of post-2010 varieties have the
299   translocation. 90% of all lines in the panel that contain the introgression have traceable pedigrees
300   through at least one ancestor. With one or two possible minor exceptions, the ancestry of the
301   translocation in all of these lines can be traced back to French landraces (e.g. "Saumur" types) through
302   early 20$^{th}$ century German ("Heines Kolben") or French ("Vilmorin-23", "Vilmorin-27") ancestral varieties,
303   strongly suggesting that the source of the translocation in modern European wheat germplasm is north-
304   western European landraces. The translocation is absent in subsequent introgressions from other grass
305   species and parents of other wide crosses. We investigated whether the high prevalence of the
306   translocation in this European material might be a result of indirect selection by wheat breeders, using a
307   subset of 135 varieties whose parents are known, were genotyped with the 90K wheat SNP array, and
308   differ for the presence of the translocation. 74/135 of these varieties (55%) inherited the translocation,
309   not significantly different than by chance ($p$=0.15). Breaking this down by decade of release, the most
310   extreme decade is the 2000s, where 22/35 varieties with different parents have the translocation
311   ($p$=0.09). However, we have previously shown across the whole pedigree that for varieties with simple
312   biparental parentage, breeders have strongly selected for the favourable parent. To account for this
313   possible co-factor, the translocation analysis was recalculated, taking into account the overall bias
314   towards one parent in each cross. Again, no significant effect was detected, whether across the whole
315   pedigree ($p$=0.10, 1,000 simulations), or in any particular decade (2000s, $p$=0.16, 1,000 simulations). It is
316   also possible that the presence of the translocation itself may be deleterious and be simultaneously
317   selected for by breeders due to linkage to advantageous haplotypes. Results from an 8-parent UK
318   multiparent advanced generation intercross (MAGIC) population that was genotyped with the same 90K
319   wheat SNP array, which has not been subjected to breeder selection, suggest that the translocation is
320   perfectly neutral (5/8 =62.5% of parents have the translocation, 62.1% of progeny have it). Furthermore,
321   in genome wide association scan (GWAS) studies using the above wheat panel, we have never detected
322   a QTL associated with the haplotype containing *RAC875_c5965_1153*. In summary, the 5B/7B major
323   translocation of almost a whole chromosome arm appears to be selectively neutral, both naturally and
324   with respect to selection by breeders, and likely owes its high frequency in north-western European
325   germplasm to its presence in the landraces commonly used in the earliest European plant breeding
326   programs.

327

328

**SUPPLEMENTARY NOTE 8**

330 **HAPLOTYPE BLOCK ANALYSES AND VISUALIZATION**

331

332 Genotyping and marker assisted selection have quickly become standard practice in most modern wheat
333 breeding programs. Genotyping involves the use of DNA markers that are in linkage with the gene or
334 locus of interests; however there may be several other genes that are also in linkage with the marker
335 that will also get inherited, many of which will impact crop performance. For breeding purposes, it is
336 important to identify regions of the genome that are similar between lines as well as regions that may
337 be in genetic linkage and are therefore inherited together. Towards this, we compared the genomes of
338 The 10+ Wheat Genomes lines to identify haplotypes that can be used in applied breeding programs.

339

340 Haplotype blocks were calculated using a combination of whole chromosome level pairwise alignments
341 using MUMmer v4.0 and gene-based pairwise alignments using BLAST v2.8. For the MUMmer
342 alignments, chromosome level pairwise alignments between all RQA were performed for each
343 chromosome. For alignments between RQA and scaffold-level assemblies, scaffolds were filtered to
344 contain at least one RefSeqv1.0 gene model projection from the corresponding chromosome. Pairwise
345 alignments between scaffold-level assemblies were performed to the RQA, but not to each other.
346 Haplotype blocks from the MUMmer alignments were called using the script
347 assign_mummer_blocks_whole_genome.r. All scripts used for haplotype construction were deposited to
348 github (https://github.com/Uauy-Lab/pangenome-haplotypes). For the zoomed in haplotypes
349 surrounding the *Sm1* gene (**Fig. 3b**), blocks were defined as above but using a reduced bin size of 250 kb.
350 To complement the MUMmer alignments and allow for direct pairwise comparisons between scaffold-
351 level assemblies, pairwise BLAST alignments of projected genes +/- 2000 bp were conducted for all
352 genome assemblies. Alignments were filtered to remove any alignments containing Ns in the aligned
353 sequence. For each pairwise comparison, gene-based alignments were ordered based on the Chinese
354 Spring RefSeqv1.0 physical position. Haplotype blocks were then called using a sliding window approach
355 using the script assign_BLAST_blocks_whole_genome.r. Haplotype blocks called using MUMmer v4.0
356 and BLAST were combined using the script combine_mummer_and_BLAST.r. To account for slight
357 differences in the absolute positions of haplotype blocks in **Fig. 3b**, chromosomes were scaled according
358 to the largest chromosome 2B across the RQA, and the coordinates of the haplotype blocks were
359 averaged across the assemblies. The positions of the *Sm1* gene and associated markers was determined
360 using BLAST alignments of the gDNA and marker sequence, respectively, against all assemblies. A
361 database and interactive visualization of the haplotypes has been made available to facilitate gene
362 discover and breeding efforts (http://www.crop-haplotypes.com/).

363

364 To complement the haplotype database, we constructed additional genome visualization tools for
365 examining larger structural variation. Pairwise gene comparisons by BLASTn were combined into larger
366 blocks using MCScanX v2.0 and the annotated positions of the projected gene annotations. The data was
367 then imported into AccuSyn (https://accusyn.usask.ca/) and SynVisio (https://synvisio.github.io/#/)
368 visualization tools, with menu options to select genomes for pairwise comparison
369 (https://kiranbandi.github.io/10wheatgenomes/). Pretzel (https://github.com/plantinformatics/pretzel)
370 was also used to visualize and compare the RQA and the projected gene annotations
371 (http://10wheatgenomes.plantinformatics.io/). These tools provide access to linear, multi-dimensional,
372 and circular visualizations comparing the RQA, as well as options to upload additional data tracks by
373 research scientists and breeders using these genomes.