

# Reference genome assemblies reveal the origin and evolution of allohexaploid oat

**Yuanying Peng** (✉ [yy.peng@hotmail.com](mailto:yy.peng@hotmail.com))

Sichuan Agricultural University <https://orcid.org/0000-0002-3304-3301>

**Honghai Yan**

Sichuan Agricultural University

**Laichun Guo**

Baicheng Academy of Agricultural Sciences

**Cao Deng**

DNA Stories Bioinformatics Center <https://orcid.org/0000-0001-9823-9121>

**Lipeng Kang**

Institute of Genetics and Developmental Biology <https://orcid.org/0000-0003-0055-9149>

**Chunlong Wang**

Baicheng Academy of Agricultural Sciences

**Pingping Zhou**

Sichuan Agricultural University

**Kaiquan Yu**

Sichuan Agricultural University

**Xiaolong Dong**

Sichuan Agricultural University

**Jun Zhao**

Sichuan Agricultural University

**Yun Peng**

Sichuan Agricultural University

**Xiaomeng Liu**

Sichuan Agricultural University

**Di Deng**

Sichuan Agricultural University

**Yinghong Xu**

Sichuan Agricultural University

**Ying Li**

Sichuan Agricultural University

**Qiantao Jiang**

Sichuan Agricultural University

**Yan Li**

Rice Research Institution <https://orcid.org/0000-0002-6443-9245>

**Liming Wei**

Baicheng Academy of Agricultural Sciences

**Jirui Wang**

Sichuan Agricultural University

**Jian Ma**

Sichuan Agricultural University

**Ming Hao**

Sichuan Agricultural University

**Wei Li**

Sichuan Agricultural University

**Houyang Kang**

Triticeae Research Institute, Sichuan Agricultural University

**Youliang Zheng**

Sichuan Agricultural University

**Yuming Wei**

Sichuan Agricultural University

**Fei Lu**

Institute of Genetics and Developmental Biology <https://orcid.org/0000-0002-3596-1712>

**Changzhong Ren**

Baicheng Academy of Agricultural Sciences

---

**Biological Sciences - Article**

**Keywords:** common oat, avena sativa, plant evolution, genome sequencing

**Posted Date:** June 29th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-664692/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Nature Genetics on July 18th, 2022. See the published version at <https://doi.org/10.1038/s41588-022-01127-7>.

1 **Reference genome assemblies reveal the origin and evolution**  
2 **of allohexaploid oat**

3 Yuanying Peng<sup>1,2,3,4,10\*</sup>, Honghai Yan<sup>1,3,10</sup>, Laichun Guo<sup>2,4,10</sup>, Cao Deng<sup>5,6,10</sup>, Lipeng Kang<sup>7,8,10</sup>,  
4 Chunlong Wang<sup>2,4</sup>, Pingping Zhou<sup>1,3</sup>, Kaiquan Yu<sup>3</sup>, Xiaolong Dong<sup>3</sup>, Jun Zhao<sup>3</sup>, Yun Peng<sup>3</sup>, Xiaomeng  
5 Liu<sup>3</sup>, Di Deng<sup>3</sup>, Yinghong Xu<sup>3</sup>, Ying Li<sup>3</sup>, Qiantao Jiang<sup>1,3</sup>, Yan Li<sup>1</sup>, Liming Wei<sup>2,4</sup>, Jirui Wang<sup>1,3</sup>, Jian  
6 Ma<sup>1,3</sup>, Ming Hao<sup>1,3</sup>, Wei Li<sup>1,3</sup>, Houyang Kang<sup>1,3</sup>, Youliang Zheng<sup>1,3</sup>, Yuming Wei<sup>1,3\*</sup>, Fei Lu<sup>7,8,9\*</sup>,  
7 Changzhong Ren<sup>2,4,11\*</sup>

8

9 <sup>1</sup> State Key Laboratory of Crop Gene Exploration and Utilization in Southwest China, Sichuan  
10 Agricultural University, Chengdu, China.

11 <sup>2</sup> National Oat Improvement Center, Baicheng Academy of Agricultural Sciences, Baicheng, China.

12 <sup>3</sup> Triticeae Research Institute, Sichuan Agricultural University, Chengdu, China.

13 <sup>4</sup> China Oat and Buckwheat Research Center, Baicheng, China.

14 <sup>5</sup> The Key Laboratory of Animal Disease and Human Health of Sichuan Province, College of  
15 Veterinary Medicine, Sichuan Agricultural University, Chengdu, China.

16 <sup>6</sup> Departments of Bioinformatics, DNA Stories Bioinformatics Center, Chengdu, China

17 <sup>7</sup> State Key Laboratory of Plant Cell and Chromosome Engineering, Institute of Genetics and  
18 Developmental Biology, Innovative Academy of Seed Design, Chinese Academy of Sciences, Beijing,  
19 China.

20 <sup>8</sup> University of Chinese Academy of Sciences, Beijing, China.

21 <sup>9</sup> CAS-JIC Centre of Excellence for Plant and Microbial Science (CEPAMS), Institute of Genetics and  
22 Developmental Biology, Chinese Academy of Sciences, Beijing, China.

23 <sup>10</sup> These authors contributed equally to this work.

24 <sup>11</sup> Lead contact.

25 \*Correspondence: Changzhong Ren (renchangzhong@163.com), Fei Lu (flu@genetics.ac.cn),

26 Yuming Wei (ymwei@sicau.edu.cn), Yuanying Peng (yy.peng@hotmail.com)

27 **Common oat (*Avena sativa*) is one of the most important cereal crops serving as a**  
28 **valuable source of forage and human food. While reference genomes of many**  
29 **important crops have been generated, such work in oat has lagged behind,**  
30 **primarily owing to its large, repeat-rich, polyploid genome. By using Oxford**  
31 **Nanopore ultralong sequencing and Hi-C technologies, we have generated the**  
32 **first reference-quality genome assembly of hulless common oat with a contig N50**  
33 **of 93 Mb. We also assembled the genomes of diploid and tetraploid *Avena***  
34 **ancestors, which enabled us to identify oat subgenome, large-scale structural**  
35 **rearrangements, and preferential gene loss in the C subgenome after**  
36 **hexaploidization. Phylogenomic analyses of cereal crops indicated that the oat**  
37 **lineage descended before wheat, offering oat as a unique window into the early**  
38 **evolution of polyploid plants. The origin and evolution of hexaploid oat is**  
39 **deduced from whole-genome sequencing, plastid genome and transcriptomes**  
40 **assemblies of numerous *Avena* species. The high-quality reference genomes of**  
41 ***Avena* species with different ploidies and the studies of their polyploidization**  
42 **history will facilitate the full use of crop gene resources and provide a reference**  
43 **for the molecular mechanisms underlying the polyploidization of higher plants,**  
44 **helping us to overcome food security challenges.**

45 Common oat (*Avena sativa* L.,  $2n = 6x = 42$ ) is the sixth most important cereal crop  
46 cultivated worldwide <sup>1</sup> and has long been prized by consumers largely because it is  
47 one of the richest sources of protein, fat and VB1 among all crops <sup>2</sup>. In addition, it is  
48 the most widely grown cool-season annual forage species, representing a major source  
49 of forage for livestock globally <sup>3</sup>. Polyploid plants often have significant advantages  
50 concerning biomass production, vigorousness, and robust adaptation to environmental  
51 changes and contribute to the emergence of important agronomic traits in food crops  
52 <sup>4-8</sup>. Therefore, crop polyploidization may play an important role in next-generation  
53 crop improvement aimed at overcoming food security challenges <sup>8</sup>. Many  
54 commercially important crops have been sequenced and assembled, which has  
55 improved the understanding of crop evolutionary history and the development of  
56 efficient approaches for the selection of important traits <sup>9</sup>. Oat research has lagged  
57 behind in this regard, primarily due to the large hexaploid oat genome <sup>10</sup> containing  
58 highly repetitive DNA sequences, in which two of the three subgenomes are too  
59 homologous to distinguish from each other <sup>11</sup>. Relatively little is currently known  
60 about the position and distribution of genes on each of the oat chromosomes and their  
61 evolution during the polyploidization events that gave rise to hexaploid species, which  
62 limits the full and effective utilization of oat germplasms.

63 With the recent advances made by Oxford Nanopore Technologies (ONT), the ONT  
64 system now offers ultralong sequence reads, delivering high contiguity with low  
65 assembly errors caused by long repetitive regions <sup>12</sup>. This technology has facilitated  
66 complete telomere-to-telomere (T2T) genome assembly in various species, including  
67 *Homo sapiens* by resolving long, complex repetitive regions <sup>13,14</sup>. It is very suitable  
68 for the assembly of the large, complex polyploid oat genome, with a high content of  
69 repetitive sequences and high subgenomic homology. Here, we employed the ONT

70 system to sequence the genome of the Sanfensan oat variety (*Avena sativa* ssp. *nuda*  
71 cv. Sanfensan, abbreviated SFS; see Methods), an ancient, important hexaploid oat  
72 landrace that originated from the diversity centre of hulless oat, together with  
73 assembling the genomes of additional diploid and tetraploid *Avena* species to address  
74 the phylogenomic relationship of cereal crops and gain insights into the evolutionary  
75 processes that established the dominant subgenome in allopolyploid oat species with  
76 large, complex genomes.

77

### 78 ***Assembly and annotation of the oat genome***

79 Using the PromethION platform and the ultralong read sequencing strategy<sup>13,15</sup>, we  
80 assembled the SFS genome into 329 contigs based on 1028 Gb cleaned ultralong  
81 reads (max length: 1.38 Mb), representing ~100× theoretical coverage with a contig  
82 N50 of 93.26 Mb and a maximum length of 405.55 Mb (Supplementary Tables 1-3).  
83 The longest 119 contigs contained 90% of the sequences, with the 34 longest contigs  
84 covering half the genome. The total assembly size was 10.76 Gb, including 10.44 Gb  
85 anchored to 21 pseudochromosomes using ~1.37 billion valid unique interaction pairs  
86 generated from 1296 Gb of cleaned Hi-C data (Supplementary Table 1, Extended Data  
87 Fig. 1a). The quality and contiguity of the genome assembly were assessed through  
88 alignments with the hexaploid consensus map and BUSCO pipeline. Most markers  
89 from the same linkage group were anchored to the corresponding pseudochromosome  
90 (Extended Data Fig. 2a), and 97.75% of the 1375 BUSCO genes were identified in the  
91 SFS genome assemblies, indicating a high-quality genome assembly (Extended Data  
92 Fig. 2b). A key feature of this assembly was its long-range organization, with contigs  
93 longer than 20.9 Mb representing 90% of the genome and an average of 15 contigs  
94 representing one chromosome (Table 1). The longest contig spanned 405 Mb, which

95 is larger than the size of rice genome <sup>16</sup>. This assembly represents a substantial  
96 improvement over the hexaploid oat genome sequenced by Pepsi Co. in 2020 (>1900  
97 contigs with a contig N50 of 30.27 Mb,  
98 [https://wheat.pw.usda.gov/jb/?data=/ggds/oat-ot3098-pepsico.](https://wheat.pw.usda.gov/jb/?data=/ggds/oat-ot3098-pepsico)) and recent assemblies  
99 reported for polyploid species with large, complex genomes <sup>17,18</sup>.

100 Protein-coding genes were annotated using a combination of ab initio prediction  
101 and transcript evidence gathered from RNA sequences from multiple tissues based on  
102 PacBio isoform sequencing (Iso-seq) approaches and reference protein sequences  
103 from the other six closely related grass species (see Methods). In total, 120,769  
104 protein-coding genes were identified (Table 1, Supplementary Table 4), 88.41% of  
105 which were assigned to a predicted function (Supplementary Table 5). Among these  
106 genes, the majority were retained in duplicate, largely owing to the allohexaploid  
107 nature of the genome. The C subgenome contained fewer gene loci (33181) than the A  
108 (40,085) and D (41,633) subgenomes. We also identified 59,916 noncoding RNAs,  
109 including 5,386 rRNAs, 50,536 small RNAs, 3,712 rRNAs and 283 cis-regulatory  
110 elements (Supplementary Table 6). Moreover, a total of 9.5 Gb (88.64%) of the SFS  
111 assembly was annotated as repetitive sequences (Table 1, Supplementary Table 7).  
112 Transposable element (TE)-related sequences constituted 86.2% of the total genome  
113 assembly. This TE content was higher than those previously reported for barley  
114 (80.80%) <sup>19</sup> and bread wheat (84.70) <sup>9</sup>. Long terminal repeat (LTR) retrotransposons,  
115 including Gypsy elements, unclassified retrotransposon elements and Copia elements,  
116 were the most abundant TEs and constituted 71.66% of the assembled SFS genome.

117 To distinguish the subgenomes accurately and clarify the polyploidization history  
118 of hexaploid oat, we also sequenced, assembled and annotated the most likely  
119 ancestral diploid species, *A. longiglumis*, and tetraploid species, *A. insularis* <sup>20</sup>,

120 resulting in >60-fold genome coverage for *A. longiglumis* (218.67 Gb) and *A.*  
121 *insularis* (374.77 Gb). The assembled *A. longiglumis* genome was 3.74 Gb in size,  
122 and 99.20% (3.71 Gb) of its sequences were anchored on the 7 chromosomes of *A.*  
123 *atlantica*<sup>21</sup>. The tetraploid *A. insularis* genome was 7.52 Gb in size, 95.08% (7.15 Gb)  
124 of which was arranged into 14 chromosomes by Hi-C analysis (Extended Data Fig.  
125 1b). We annotated 43,352 protein-coding genes in the *A. longiglumis* genome and  
126 87,154 in the *A. insularis* genome. BUSCO analysis also revealed that a majority  
127 (97.53% and 98.11%, respectively) of conserved genes were identified in these two  
128 genome assemblies. Accordingly, most of the genome sequences of the *A. longiglumis*  
129 (87.82%) and *A. insularis* (87%) were composed of repetitive elements.

130 Based on the genomic synteny and similarity among the hexaploid, diploid, and  
131 tetraploid genomes, we successfully divided the 21 pseudochromosomes of hexaploid  
132 oat into A, C, and D subgenomes (Extended Data Fig. 3a-d). Am1 (a repeat selectively  
133 hybridized to the C subgenome) and As120a (a repeat selectively hybridized to the A  
134 subgenome) satellite repeats were overrepresented in the C and A subgenomes,  
135 respectively (Fig. 1). The subgenome assignments were further validated by mapping  
136 the molecular markers from a recently published high-density linkage map of  
137 hexaploid oat (Extended Data Fig. 2a) and by mapping the whole-genome sequencing  
138 reads from a range of *Avena* species (Extended Data Fig. 3e-f). The nomenclature  
139 system for bread wheat was adapted for naming the chromosomes of SFS (Extended  
140 Data Fig. 4).

141

#### 142 ***Evolutionary position of oat among cereal crops***

143 We next used our oat genome assemblies to assess the phylogenomic relationships  
144 among common oat and other cereal crops (Supplementary Table 8) by using 644

145 conserved core genes (Fig. 2a, Supplementary Fig. 1). Molecular dating estimated that  
146 the divergence of monocot cereal crops occurred ~45 million years ago (mya) and was  
147 concentrated in Poaceae. Pooideae and Oryzoideae split ~40 mya, whereas  
148 Panicoideae and Chloridoideae diverged ~33 mya. In Pooideae, the tribes Aveneae  
149 and Lolieae were indicated to be more closely related than to Triticeae. The oat  
150 lineage descended (~16 mya) earlier than that of wheat (~7 mya), indicating a longer  
151 time of evolution and offering a unique window into the early evolution of polyploid  
152 cereal crops.

153 A high-confidence orthologous gene set for the 43 species was obtained, 2202 gene  
154 families shared by all the genomes/subgenomes were identified, and 4387, 419, 459,  
155 24, 152, 450, and 40 gene families were indicated to be specific to Polygonaceae,  
156 Chenopodiaceae, Chloridoideae, Panicoideae, Aveninae, and Triticinae, respectively  
157 (Fig. 2b, Supplementary Table 9). Most of these *Avena*-specific gene families were  
158 assigned to molecular binding or metabolic process functions (Supplementary Table  
159 10). In addition, 10 of them may play roles in fat biosynthesis, including 479 genes  
160 involved in fatty acid biosynthesis. These genes might contribute to the much higher  
161 concentration of oil in oat grains relative to other cereal crops (Supplementary Table  
162 10).

163 The ancestral grass karyotype (AGK) was inferred by comparing modern species  
164 via a completely automated method reconstructing the order of ancestral genes within  
165 contiguous ancestral regions and is structured into 7 protochromosomes containing  
166 7,010 protogenes<sup>22</sup>. We therefore used the 7 AGK protochromosomes as an ancestral  
167 reference to investigate the chromosomal evolution of rice, wheat and oat (Fig. 2c). A  
168 total of 731 syntenic blocks (>5 genes) encompassing 19112 orthologous gene pairs  
169 were identified between oat and the 7 AGK protochromosomes (Supplementary Table

170 11). Essentially, the third homologous group of oat was derived from a single ancient  
171 chromosome, AGK1, (which also steadily gave rise to the complete chromosomes of  
172 rice and wheat), with the exception of a segment of AGK2 that was translocated to  
173 chromosome 3C in oat. Group six was derived from the insertion of AGK5 into  
174 AGK4. The remaining five of the seven oat homologous groups were each derived  
175 from at least three ancestral chromosomes via complex translocations (Fig. 2c). When  
176 the oat assembly was compared with the three subgenomes of common wheat, the  
177 ancestral chromosomes were found to be roughly similar, but a large number of  
178 chromosomal rearrangements occurred, leading to many structural variations and  
179 lower collinearity.

180

### 181 ***Polyploidization history and reticulate evolution***

182 Previous phylogenetic studies aimed at identifying the progenitors of these species  
183 have often analysed only a limited set of molecular markers, yielding inconsistent  
184 results. To more accurately identify the relatives of each subgenome donor and the  
185 extant ancestral hexaploid species, we sequenced and de novo assembled  
186 transcriptomes, conducted whole-genome sequencing and assembled the complete  
187 plastid genomes of other *Avena* species. These *Avena* species represented different  
188 genomic subtypes (As, Al, Ac, Ad, Cv, and Cp) and different ploidies.

189 We sequenced 14 different *Avena* species (generating 6767.64 Gb of sequence data,  
190 Supplementary Table 1) and mapped the obtained data onto the reference SFS genome.  
191 The genomic variants of these 14 species as well as the assembled *A. longiglumis* and  
192 *A. insularis* genomes relative to the SFS sequences were called based on the mapped  
193 reads. Approximately 39~129 million SNPs were identified in these species  
194 (Supplementary Table 12). The ratios of SNP numbers in the A, C and D subgenomes

195 of each species showed that the C subgenome is differed greatly from the A, B and D  
196 subgenomes (Fig. 3a). These SNP-based phylogenetic analyses indicated that the  
197 species clustered according to the genome composition except for the AB genome  
198 tetraploids, which were positioned closer to different A genome diploids (Fig. 3b).  
199 Three maximum likelihood trees based on the A-, C- and D-type SNPs further implied  
200 that *A. longiglumis* was the most closely related A genome donor for hexaploids,  
201 while no extant diploids showed a particularly close relationship with the C and D  
202 subgenomes (Extended Data Fig. 5). The identity distribution of different sequenced  
203 reads clearly showed that the A1 genome of the diploid species *A. longiglumis*, and the  
204 C and D subgenomes of *A. insularis* have the closest relationships with the A, C and  
205 D subgenomes of hexaploids, respectively (Fig. 3c). These results further confirmed  
206 that the diploid *A. longiglumis* and the tetraploid *A. insularis* were the most closely  
207 related extant species to hexaploid oat.

208 We sequenced and de novo assembled 11 transcriptomes of different diploid species  
209 (Supplementary Table 13) to further clarify the relationships between extant diploid  
210 species and cultivated hexaploid oat. The phylogenetic analyses based on 2863  
211 nuclear genes using barley as the outgroup again provided strong support for *A.*  
212 *longiglumis* as the extant diploid with the closest relationship to the A subgenome of  
213 hexaploids and for no diploid being a D genome donor. However, the A subgenome of  
214 hexaploid *A. sativa* differentiated from the ancestral genome earlier than the A  
215 genome of diploid *A. longiglumis* (Fig. 2a, Fig. 3d). This finding suggests that the A  
216 genome progenitor of hexaploids was more likely the A diploid ancestor than an  
217 extant A genome diploid, which may be the reason that previous studies<sup>23</sup> have  
218 considered different A genome diploids to be the donors of hexaploid genomes.

219 To determine the maternal origin of hexaploid oat, we assembled approximately

220 complete plastid genomes of 136 kb for each *Avena* species. Phylogenetic analysis of  
221 18 complete *Avena* plastid genomes together with the previous 26 *Avena* chloroplast  
222 genomes (Supplementary Tables 14-15)<sup>24,25</sup> showed that the C genome was  
223 undoubtedly the male parent in polyploid formation and that the D genome, rather  
224 than the A genome, was the maternal donor of hexaploid oat. The evolutionary order  
225 of the different A genome subtypes was Ac-Ad-Al-As (Fig 3e, Supplementary Fig. 2).  
226 The relatively low collinearity between C genome diploid species and the C  
227 subgenome of *Avena* polyploids was consistent with the nuclear-cytoplasmic  
228 interaction hypothesis suggesting that the paternally inherited genome of an  
229 allopolyploid is usually more prone to genetic changes than the maternally derived  
230 genome<sup>26</sup>.

231 Our phylogenetic analyses revealed a common ancestor of A and D genome  
232 diploids, indicating that the D genome ancestor is showed a closer relationship to the  
233 A genome and may be extinct. The C and A/D lineages diverged from the ancestor of  
234 the genus *Avena* ~17 mya, followed by the A genome subtypes, and the D genome  
235 split from a common ancestor ~9 mya (Fig. 3d, 3f). Cultivated oat initially originated  
236 from an allotetraploidy event between a paternal C genome and a maternal D genome  
237 diploid (~2 mya, Supplementary Fig. 4), followed by hybridization between a paternal  
238 A genome diploid progenitor (closely related to *A. longiglumis*) and maternal CD  
239 genome tetraploid (closely related to *A. insularis*) to form the extant ACD genome  
240 hexaploids (~0.5 mya, Supplementary Table 16, Supplementary Fig. 3). Our findings  
241 clarify the evolutionary history of oats based on various lines of evidence (Fig. 3f)  
242 and have broad implications for the understanding of genome function and cultivar  
243 improvement in oat.

244

245 ***Subgenome structure, content, and response of subgenomes to polyploidization***

246 Synteny analyses between the subgenomes of hexaploid oat revealed a much lower  
247 level of conservation in oat (62.70% of genes organized in collinear blocks,  
248 75727/120769) than in bread wheat (72.57%, 77292/106508). The A and D  
249 subgenomes of oat showed a higher synteny than the C and A or D subgenomes, as  
250 indicated by the numbers of collinear genes preserved between the three subgenomes  
251 (Extended Data Fig. 6a). Such low levels of conservation between the C and D  
252 subgenomes have also been observed in *A. insularis* (Extended Data Fig. 6b),  
253 supporting the hypothesis that both the A and D subgenomes were derived from a  
254 common diploid lineage that diverged from the C subgenome. The low levels of  
255 conservation between the hexaploid subgenomes, particularly between C and A or D,  
256 may reflect the preexisting differences in the subgenomes before polyploidization but  
257 could also be attributed to chromosomal rearrangements after polyploidization. To  
258 explore broad-scale structural evolution after polyploidization, we performed a  
259 whole-genome synteny comparison, which revealed low collinearity between the  
260 sequenced diploid C and tetraploid C subgenomes, especially for chromosome7C (Fig.  
261 4a), suggesting that large chromosomal rearrangements might have accumulated  
262 during evolution. We also found evidence of large chromosomal rearrangements  
263 during hexaploidization (Fig. 4a, Extended Data Fig. 7a, b). For instance, one C to A  
264 introgression and four C to D introgressions (Fig. 4c, Extended Data Fig. 7d) occurred.  
265 These large translocations were confirmed by FISH assays (Fig. 4d, Extended Data  
266 Fig. 7e). Moreover, the structural rearrangements in oat appeared biased among the  
267 three subgenomes (Fig. 4b): most of them occurred between subgenomes A and D,  
268 whereas fewer occurred between subgenomes A and C and between subgenomes D  
269 and C. For instance, half of the tetraploid 6D chromosome was transferred to

270 hexaploid 6A, 7A and 2D chromosomes, resulting in a much smaller hexaploid 6D  
271 chromosome. These results may support the hypothesis that the existence of two  
272 well-conserved homologous genomes would facilitate inter-subgenome recombination  
273 and rearrangement after polyploidization, as observed in a previous study <sup>27</sup>.

274 The C subgenome (3.94 Gb) is 20% larger than the A (3.28 Gb) and D (3.22 Gb)  
275 subgenomes in hexaploid (Table 1). TE abundance accounted for almost all of the size  
276 differences between the subgenomes (Fig. 1, Extended Fig. 8a). We further found that  
277 C genomes contained twice as many full-length LTRs as A and D genomes (Extended  
278 Data Fig. 8b), and that the average length of LTRs was shorter in the C subgenome  
279 than in the A and D subgenomes (Extended Data Fig. 8c). The amplification of  
280 full-length LTRs in C genomes occurred over a relatively long period, peaking  
281 approximately 1.8 mya, whereas LTRs were inserted into A and D genomes within the  
282 last 1 mya (Extended Data Fig. 8d), which indicates that the C genome has a longer  
283 evolutionary history than the A and D genomes. The C subgenome contained a  
284 smaller number of gene loci (33181; 29%) than the A (40,085, 35%) and D (41,633,  
285 36%) subgenomes in hexaploid (Table 1), whereas the tetraploid (89,995) showed  
286 almost the same gene number as the A (43,477) and C (46,925) diploids together  
287 (Table 1), suggesting that a remarkable gene loss occurred in the C subgenome after  
288 hexaploidization (33,181 in hexaploid vs 43,243 in tetraploid). These results were  
289 supported by the gene loss ratio in the subgenomes of *A. insularis* and SFS (Extended  
290 Data Fig. 9a, b). We further found that more gene families underwent decreases in  
291 hexaploid C subgenomes than in A and D subgenomes (Extended Data Fig. 9c), and  
292 that gene losses mainly affected genes with expanded families, rather than singleton  
293 genes (Extended Data Fig. 9d, e). Accordingly, the hexaploid C subgenome contained  
294 more pseudogenes (29297) than the A (26141) and D (27262) subgenomes (Table 1,

295 Extended Data Fig. 9f). Together, these results suggest preferential gene loss or  
296 pseudogenization in the C subgenome after hexaploidization.

297 Subgenome dominance is a striking whole-genome feature common to polyploids,  
298 presumably associated with biased fractionation<sup>28</sup>. Genes located in the dominant  
299 subgenome always present a lower nonsynonymous ratio and higher expression levels  
300 than orthologues in other subgenomes<sup>29,30</sup>. To test whether subgenome dominance has  
301 occurred in hexaploid oat, we first compared nonsynonymous ( $Ka$ ) with synonymous  
302 ( $Ks$ ) substitution rates between (sub)genomes using 2,767 one-to-one orthologous  
303 gene sets among the polyploid oat and the progenitors, and the results suggested that  
304 an accelerated evolution rate exists in polyploids relative to their diploid progenitors  
305 (Supplementary Fig. 5). In SFS, the C subgenome presented a significantly higher  
306  $Ka/Ks$  value (0.37) than the A (0.17) and D (0.25) subgenomes (Fig. 4e), suggesting  
307 that the C subgenome is subject to more relaxed selective constraints than the other  
308 two subgenomes. We further examined the expression of 13,774 single-copy genes in  
309 the three subgenomes, referred to as triads, in diverse organs and under six different  
310 abiotic stresses. We found that most genes were expressed in a subgenome-specific  
311 manner (Fig. 4f). Pairwise comparisons of the homologous gene triads revealed that a  
312 similar number of gene pairs were preferentially transcribed in the A and D  
313 subgenomes; however, the C subgenome had a significantly lower number of  
314 preferentially transcribed gene pairs than the A ( $P=1.3E-3$ ) and D ( $P=6.1E-3$ )  
315 subgenomes (Fig. 4g, Supplementary Tables 17, 18). TE abundance plays an  
316 important role in gene expression bias because TE insertion can have deleterious  
317 effects on genes leading to lower expression<sup>31</sup>. In oat, we found that the C  
318 subgenome contained more TEs and showed higher overall TE densities near genes  
319 than the A and D subgenomes; these characteristics are negatively correlated with

320 gene expression (Extended Data Fig. 10), supporting the hypothesis that subgenome  
321 expression dominance is influenced by TE-density differences between subgenomes  
322 as observed in other allopolyploids <sup>7,32</sup>. All of these results revealed a biased gene  
323 fractionation in the hexaploid C subgenome, which might be closely related to the  
324 high global TE amount and TE density near genes, and the relaxation of selection  
325 pressure of the C subgenome.

326 In conclusion, the three reference-quality genomes presented here constitute  
327 important community resources for cereal genomics and provide new insights into the  
328 evolutionary position of oat among cereal crops. This is the most comprehensive  
329 molecular phylogenetic analysis of the genus *Avena* conducted to date, because it  
330 included samples representing all extant *Avena* genomes and the greatest number (and  
331 number of different types) of molecular markers evaluated thus far. The model for the  
332 chronological formation of polyploid oat has been clarified, and the investigation of  
333 oat subgenomes evolution during polyploidization events offers a unique window into  
334 the early evolution of polyploid crops, which will accelerate improvements of oat and  
335 help to make oat a useful model for studying polyploid genome evolution.

336

337 1 FAOSTATS. Food and agriculture organization of the united nations--statistics  
338 division. <http://www.fao.org> (2021).

339 2 Fu, J. *et al.* Concise review: Coarse cereals exert multiple beneficial effects on  
340 human health. *Food chem* **325**, 126761, doi:10.1016/j.foodchem.2020.126761  
341 (2020).

342 3 Fraser, J. & McCartney, D. *Fodder oats in North America*. in *Fodder oats: A*  
343 *World Overview* (eds Suttie, J. M & Reynolds, S. G.) 19-35 (FAO, 2004).

344 4 Comai, L. The advantages and disadvantages of being polyploid. *Nat Rev*

- 345            *Genet* **6**, 836-846, doi:10.1038/nrg1711 (2005).
- 346    5        Fang, Z. & Morrell, P. L. Domestication: Polyploidy boosts domestication.  
347            *Nat Plants* **2**, 16116, doi:10.1038/nplants.2016.116 (2016).
- 348    6        Van de Peer, Y., Mizrachi, E. & Marchal, K. The evolutionary significance of  
349            polyploidy. *Nat Rev Genet* **18**, 411-424, doi:10.1038/nrg.2017.26 (2017).
- 350    7        Edger, P. P. *et al.* Origin and evolution of the octoploid strawberry genome.  
351            *Nat Genet* **51**, 541-547, doi:10.1038/s41588-019-0356-4 (2019).
- 352    8        Yu, H. *et al.* A route to de novo domestication of wild allotetraploid rice. *Cell*  
353            **184**, 1156-1170. e14, doi:10.1016/j.cell.2021.01.013 (2021).
- 354    9        Appels, R. *et al.* Shifting the limits in wheat research and breeding using a  
355            fully annotated reference genome. *Science* **361**, eaar7191,  
356            doi:10.1126/science.aar7191 (2018).
- 357    10       Yan, H. *et al.* Genome size variation in the genus *Avena*. *Genome* **59**,  
358            209-220, doi:10.1139/gen-2015-0132 (2016).
- 359    11       Jellen, E., Gill, B. & TS, C. Genomic in situ hybridization differentiates  
360            between A/D- and C-genome chromatin and detects intergenomic  
361            translocations in polyploid oat species (genus *Avena*). *Genome* **37**, 613-618,  
362            doi:10.1139/g94-087 (1994).
- 363    12       Lang, D. *et al.* Comparison of the two up-to-date sequencing technologies for  
364            genome assembly: HiFi reads of Pacific Biosciences Sequel II system and  
365            ultralong reads of Oxford Nanopore. *GigaScience* **9**,  
366            doi:10.1093/gigascience/giaa123 (2020).
- 367    13       Miga, K. H. *et al.* Telomere-to-telomere assembly of a complete human X  
368            chromosome. *Nature* **585**, 79-84, doi:10.1038/s41586-020-2547-7 (2020).
- 369    14       Logsdon, G. A. *et al.* The structure, function and evolution of a complete

370 human chromosome 8. *Nature* **593**, 101-107,  
371 doi:10.1038/s41586-021-03420-7 (2021).

372 15 Jain, M. *et al.* Nanopore sequencing and assembly of a human genome with  
373 ultra-long reads. *Nat Biotechnol* **36**, 338-345, doi:10.1038/nbt.4060 (2018).

374 16 Sasaki, T. & International Rice Genome Sequencing, P. The map-based  
375 sequence of the rice genome. *Nature* **436**, 793-800, doi:10.1038/nature03895  
376 (2005).

377 17 Lovell, J. T. *et al.* Genomic mechanisms of climate adaptation in polyploid  
378 bioenergy switchgrass. *Nature* **590**, 438-444,  
379 doi:10.1038/s41586-020-03127-1 (2021).

380 18 Huang, G. *et al.* Genome sequence of *Gossypium herbaceum* and genome  
381 updates of *Gossypium arboreum* and *Gossypium hirsutum* provide insights  
382 into cotton A-genome evolution. *Nat Genet* **52**, 516-524,  
383 doi:10.1038/s41588-020-0607-4 (2020).

384 19 Mascher, M. *et al.* A chromosome conformation capture ordered sequence of  
385 the barley genome. *Nature* **544**, 427-433, doi: 10.1038/nature22043 (2017).

386 20 Yan, H. *et al.* High-density marker profiling confirms ancestral genomes of  
387 *Avena* species and identifies D-genome chromosomes of hexaploid oat. *Theor*  
388 *Appl Genet* **129**, 2133-2149, doi:10.1007/s00122-016-2762-7 (2016).

389 21 Maughan, P. J. *et al.* Genomic insights from the first chromosome-scale  
390 assemblies of oat (*Avena* spp.) diploid species. *BMC Biol* **17**, 92,  
391 doi:10.1186/s12915-019-0712-y (2019).

392 22 Murat, F., Armero, A., Pont, C., Klopp, C. & Salse, J. Reconstructing the  
393 genome of the most recent common ancestor of flowering plants. *Nat Genet*  
394 **49**, 490-496, doi:10.1038/ng.3813 (2017).

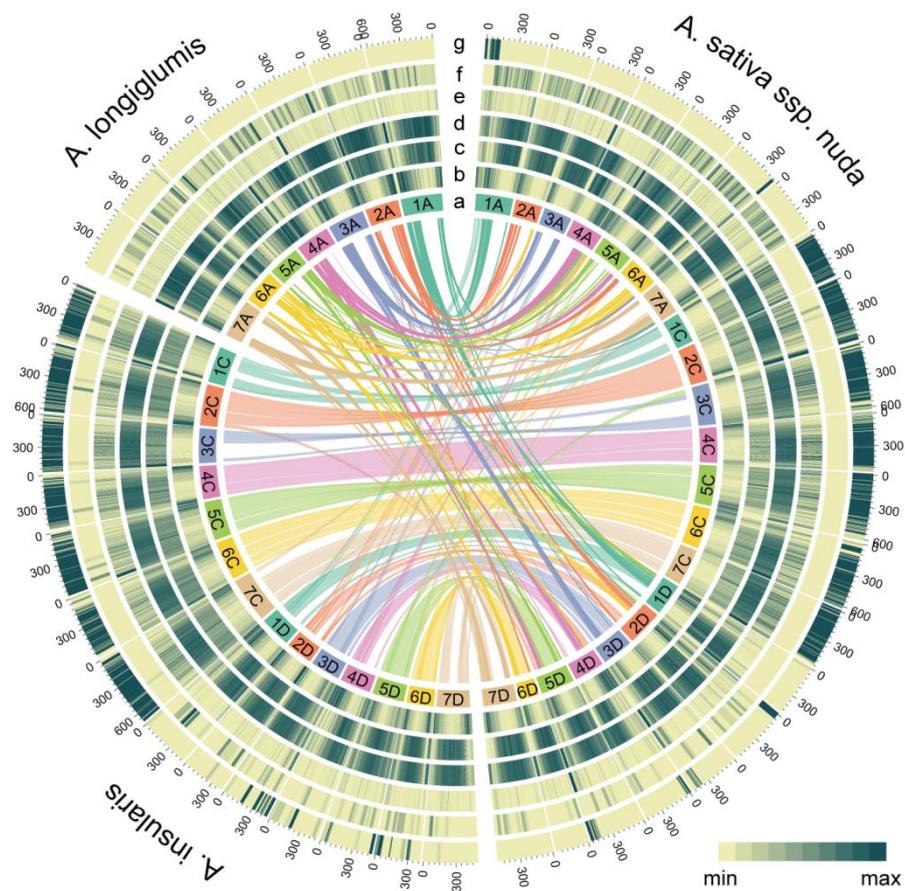
- 395 23 Peng, Y. *et al.* Phylogenetic inferences in *Avena* based on analysis of *FL*  
396 *intron2* sequences. *Theor Appl Genet* **121**, 985-1000,  
397 doi:10.1007/s00122-010-1367-9 (2010).
- 398 24 Fu, Y. B. Oat evolution revealed in the maternal lineages of 25 *Avena* species.  
399 *Sci Rep* **8**, 4252, doi:10.1038/s41598-018-22478-4 (2018).
- 400 25 Liu, Q. *et al.* Comparative chloroplast genome analyses of *Avena*: insights into  
401 evolutionary dynamics and phylogeny. *BMC Plant Biol* **20**, 406,  
402 doi:10.1186/s12870-020-02621-y (2020).
- 403 26 Gill, B. S. & Friebe, B. Nucleocytoplasmic interaction hypothesis of genome  
404 evolution and speciation in polyploid plants revisited: polyploid  
405 species-specific chromosomal polymorphisms in wheat. in *Polyploid and*  
406 *Hybrid Genomics* (eds Chen, Z. J & Birchler, J. A) 213-221,  
407 doi:10.1002/9781118552872.ch13 (Wiley, 2013).
- 408 27 Zhuang, W. *et al.* The genome of cultivated peanut provides insight into  
409 legume karyotypes, polyploid evolution and crop domestication. *Nat Genet* **51**,  
410 865-876, doi:10.1038/s41588-019-0402-2 (2019).
- 411 28 Alger, E. I. & Edger, P. P. One subgenome to rule them all: underlying  
412 mechanisms of subgenome dominance. *Curr Opin Plant Biol* **54**, 108-113,  
413 doi:10.1016/j.pbi.2020.03.004 (2020).
- 414 29 Parkin, I. A. P. *et al.* Transcriptome and methylome profiling reveals relics of  
415 genome dominance in the mesopolyploid *Brassica oleracea*. *Genome Biol* **15**,  
416 R77, doi:10.1186/gb-2014-15-6-r77 (2014).
- 417 30 Liu, S. *et al.* The *Brassica oleracea* genome reveals the asymmetrical  
418 evolution of polyploid genomes. *Nat Commun* **5**, 3930,  
419 doi:10.1038/ncomms4930 (2014).

- 420 31 Hollister, J. D. & Gaut, B. S. Epigenetic silencing of transposable elements: a  
421 trade-off between reduced transposition and deleterious effects on neighboring  
422 gene expression. *Genome Res* **19**, 1419-1428, doi:10.1101/gr.091678.109  
423 (2009).
- 424 32 VanBuren, R. *et al.* Exceptional subgenome stability and functional divergence  
425 in the allotetraploid Ethiopian cereal teff. *Nat Commun* **11**, 884,  
426 doi:10.1038/s41467-020-14724-z (2020).

**Table 1 | Assembly statistics for diploid, tetraploid and hexaploid *Avena* species.**

	<i>A. longiglumis</i> CN 58138 (2n = 2x = 14)	<i>A. insularis</i> CN 108634 (2n = 4x = 28)	<i>A. sativa</i> ssp. <i>nuda</i> cv. Sanfensan (2n = 6x = 42)
Illumina (Gb)	204.68	451.89	649.68
ONT (Gb)	268.74	481.39	-
ONT ultralong (Gb)	-	-	1260.30
HiC (Gb)	-	816.93	1312.83
IsoSeq (Gb)	25.74	49.94	81.14
Total assembly size (bp)	3,736,548,545	7,519,018,440	10,757,433,345
Longest contigs (bp)	29,014,927	36,557,065	405,550,188
Number of contigs	956	1,924	326
N50 contig length (bp)	7,297,603	7,836,599	93,262,735
L50 contig count	160	297	34
N90 contig length (bp)	2,123,884	2,126,637	20,933,943
L90 contig count	523	1,003	119
Number of contigs per (sub)genome	A: 943   Total: 943	C: 1,019 D: 848 T: 1,867	A: 84 C: 155 D: 77 T: 31
Sequences assigned per (sub)genome (bp)	A: 3,708,832,268   T: 3,708,832,268	C: 4,020,068,809 D: 3,131,824,267 T: 7,151,893,076	A: 3,279,788,166 C: 3,942,742,583 D: 3,215,854,646 T: 10,438,385,395
Genome completeness (BUSCO)	97.53	98.11	97.75
Repeats (%)	87.09%	87.11%	87.12%
Protein-coding genes	A: 43,477   T: 43,477	C: 46,925 D: 43,477 T: 89,995	A: 40,085 C: 33,181 D: 41,633 T: 120,769
Pseudogenes	A: 14,058   T: 14,058	C: 20,847 D: 17,766 T: 38,613	A: 23,013 C: 26,510 D: 24,009 T: 73,532

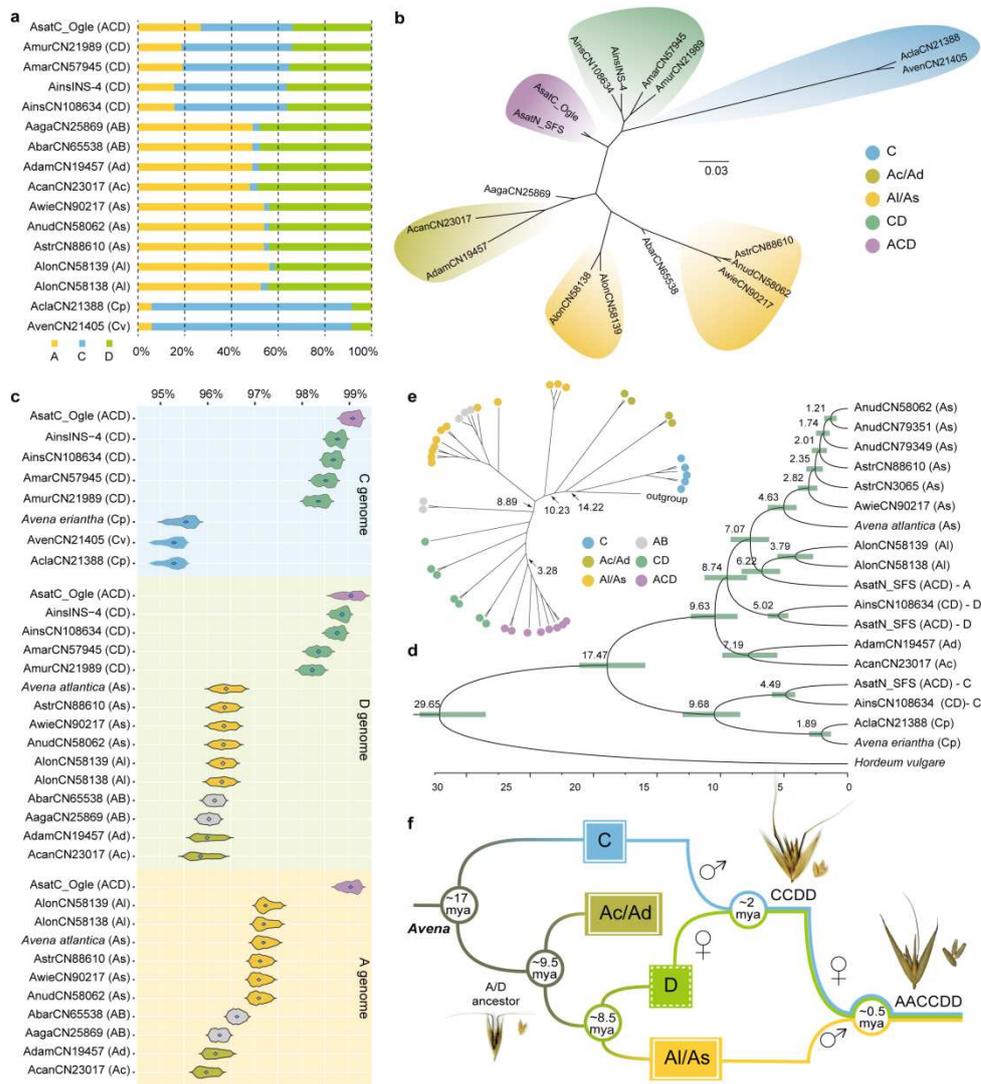
Note: A: A subgenome, C: C subgenome, D: D subgenome, T: Total.



429

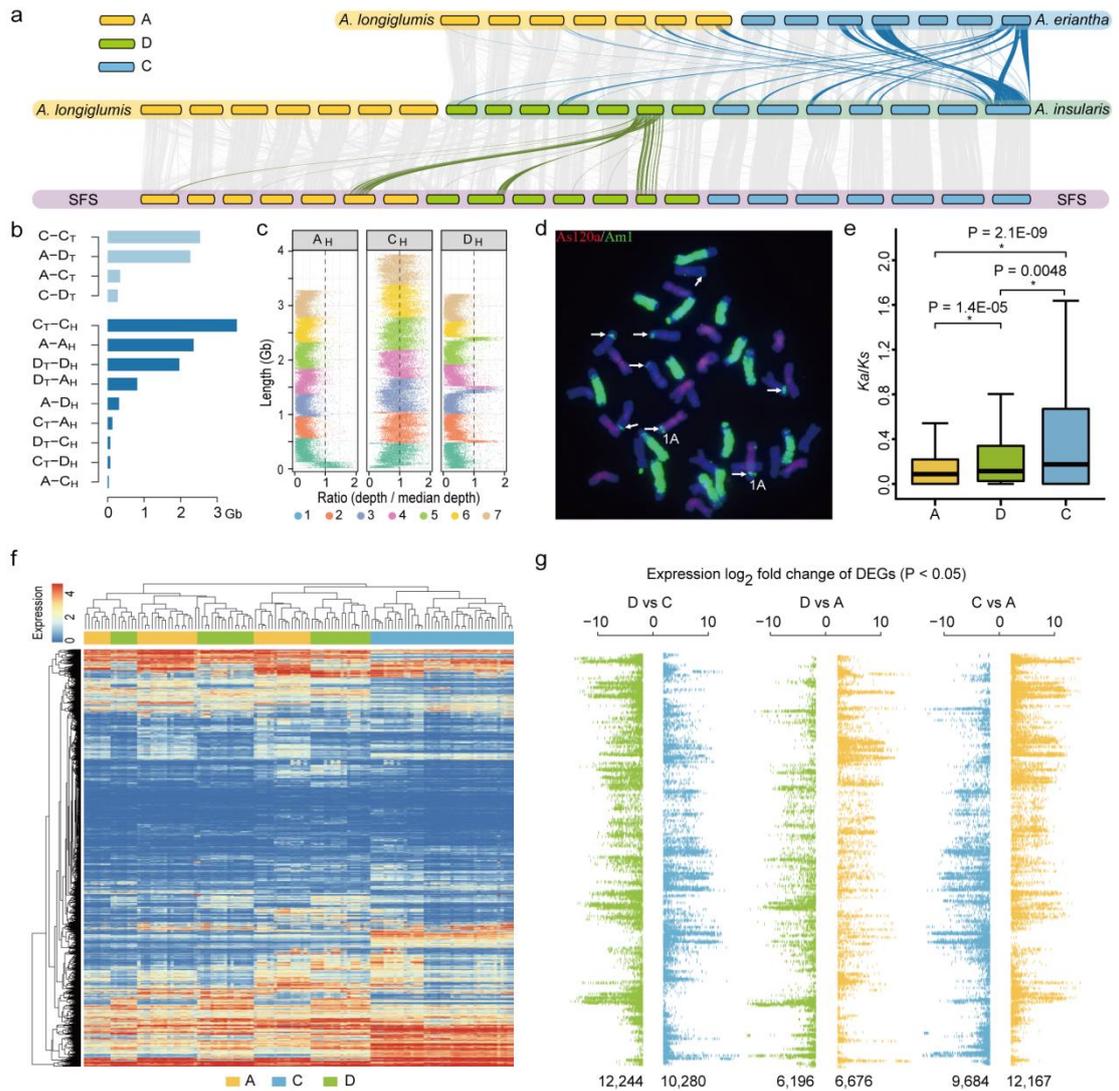
430 **Fig. 1 | Genomic features of oat.** Circos display of important features of the  
 431 assembled diploid, tetraploid and hexaploid *Avena* species genomes. **a**, Chromosome  
 432 names and sizes. **b**, Gene density. **c**, Long terminal repeat (LTR) retrotransposon  
 433 density. **d**, Tandem repeat (TR) density. **e**, *k*-mer frequencies. **f**, The A  
 434 genome-specific repeat As120a. **g**, The C genome-specific repeat Am1. Links  
 435 between syntenic genes are shown in different colours.





446

447 **Fig. 3 | Phylogeny and polyploidization of *Avena* species.** **a**, Comparable numbers  
 448 of A-, C- and D-type SNPs for each taxon. **b**, Maximum likelihood tree generated  
 449 from the SNPs. **c**, Identity distribution of different sequenced reads on different  
 450 subgenomes of hexaploids. **d**, Phylogenetic analyses based on 2,863 nuclear genes  
 451 using barley as an outgroup. **e**, Phylogenetic studies of complete plastid *Avena*  
 452 genomes. **f**, Model of the phylogenetic history of oat (*Avena sativa*; AACCCDD)



453

454 **Fig. 4 | Genomic structures of hexaploid oat and its relatives.** **a**, Different colours  
 455 of lines depict synteny between subgenomes of SFS and putative tetraploid and  
 456 diploid ancestors. Thick lines connecting chromosomes 7C and 6D depict observed  
 457 large-scale chromosome rearrangements. **b**, Length of individual subgenome DNA  
 458 sequences inherited from ancestral genomes. Subscript letter represents subgenomes  
 459 of hexaploid (H) or tetraploid (T). **c**, Reads of the C genome diploid mapped to the  
 460 SFS reference reveal at least four large C to D genomic exchanges and one C to A  
 461 genomic exchange after polyploidization. **d**, FISH using the C genome-specific repeat  
 462 as a probe confirms the C to D and C to A translocations. Red and green represent  
 463 FISH signals detected with an A-specific repeat (As120a) or a C-specific repeat (Am1)  
 464 as probes. White arrows indicate C to D or C to A subgenome translocations. **e**,

465 Comparison of *Ka/Ks* value distributions between three subgenomes of SFS. The  
466 central line for each box plot indicates the median. The top and bottom edges of the  
467 box indicate the 25th and 75th percentiles, and the whiskers extend 1.5 times the  
468 interquartile range beyond the edges of the box. The asterisks represent significant  
469 differences (\* Wilcoxon rank-sum test,  $P < 0.05$ ). **f**, Two-dimensional hierarchical  
470 cluster analysis of expression among single-copy homologous oat genes compared  
471 with organ-specific gene expression. **g**, Analysis of log<sub>2</sub>-fold changes in pairwise  
472 gene expression between homologous genes.  
473

## 474 **Methods**

### 475 **Plant materials**

476 The hexaploid species SFS ( $2n=6x=42$ , AACDD), the diploid species *A. longiglumis*  
477 (accession CN 58139,  $2n=2x=14$ , AIAI) and the tetraploid species *A. insularis*  
478 (accession CN 108634,  $2n=4x=28$ , CCDD) were selected for genome sequencing and  
479 assembly. SFS is a traditional hulless oat landrace that has a long cultivation history in  
480 Shanxi, China, which is thought to be the region of origin of hulless oat varieties. It  
481 has also been widely used as a parental line in hulless oat breeding programs, while *A.*  
482 *insularis* and *A. longiglumis* have been considered the most likely tetraploid and  
483 diploid ancestors of hexaploid oat<sup>20</sup>.

### 484 **Genome sequencing and assembly**

485 Illumina HiSeq X-Ten or MGISEQ2000 platform was used to generate short  
486 paired-end reads for SFS, *A. insularis* and *A. longiglumis*. For Illumina sequencing,  
487 paired-end libraries were constructed by using the TruSeq Nano DNA HT Sample  
488 preparation kit (Illumina, USA) according to the manufacturer's instructions. For  
489 MGI sequencing, method described in Supplementary Note was used for paired-end  
490 libraries construction. A total of 649.7 Gb, 451.9 Gb and 204.7 Gb of raw reads were  
491 generated for SFS, *A. insularis* and *A. longiglumis*, respectively (Supplementary Table  
492 1). Oxford Nanopore Technology (ONT) platforms were used to generate long reads  
493 for these three taxa. Considering the large and complex genome of hexaploid oat, the  
494 ONT ultralong reads were generated for the SFS, whereas regular ONT long reads  
495 were generated for the genomes of *A. longiglumis* and *A. insularis*. For ONT ultralong  
496 reads library construction, high-molecular-weight genomic DNA molecules were  
497 extracted from young leaves using the SDS method without the purification step to  
498 maintain the length of the DNA. Then, 8-10  $\mu\text{g}$  of high-quality DNA was

499 size-selected (>50 kb) and processed using the Ligation Sequencing 1D kit  
500 (SQK-LSK109, Oxford Nanopore Technologies, UK) according to the manufacturer's  
501 instructions. For regular Nanopore library construction, 3-4 µg of  
502 high-molecular-weight DNA was size selected (>20 kb) and ligated using the  
503 SQK-LSK109 kit with the recommended protocol. All the prepared libraries were  
504 sequenced on the PromethION platform at the Genome Center of Grandomics  
505 (Wuhan, China). A total of 1260.3 Gb, 481.4 Gb and 268.7 Gb of raw ONT  
506 (ultra-)long reads for SFS, *A. insularis* and *A. longiglumis* were produced from 71, 7  
507 and 8 libraries, covering approximately 100-, 60- and 60- fold of their genomes,  
508 respectively (Supplementary Table 2).

509 The raw ONT long reads were subjected to self-correction using the NextCorrect  
510 module implemented in NextDenovo (<https://github.com/Nextomics/NextDenovo>)  
511 (v2.0-beta.1). The corrected reads were then assembled into contigs by using  
512 NextDenovo (Supplementary Table 3). To improve the accuracy of the assembled  
513 contigs, a two-step polishing strategy was applied: the corrected Nanopore reads were  
514 first used for initial polishing with Racon<sup>33</sup> (three rounds), and the highly accurate  
515 short reads were then used to further correct the assemblies with NextPolish (four  
516 rounds).

517 We employed Hi-C technology to obtain chromosome-level genome assemblies of  
518 SFS and *A. insularis*. Hi-C libraries were created from tender leaves of SFS and *A.*  
519 *insularis*. In brief, tender leaves were fixed with formaldehyde, and the cross-linked  
520 DNA was then isolated, purified, and digested with *DpnII*. Sticky ends were filled and  
521 marked with biotin-14-dATP; the resulting blunt-end fragments were ligated to form  
522 chimaeric junctions, physically sheared, and enriched for fragments with a size of  
523 300-600 bp. Chimaeric fragments representing the original cross-linked long distance

524 physical interactions were then processed to obtain paired-end sequencing libraries,  
525 which were sequenced on the Illumina platform. A total of 1312.8 Gb and 816.9 Gb of  
526 raw Hi-C data were generated for SFS and *A. insularis*. The polished contigs of SFS  
527 and *A. insularis* were further clustered, ordered, anchored to pseudochromosomes  
528 using Hi-C data with the LACHESIS program<sup>34</sup>, whereas RaGOO<sup>35</sup> with the default  
529 parameters was employed to anchor the contigs of *A. longiglumis* to seven  
530 pseudochromosomes with the previously published *As* genome of diploid *A. atlantica*  
531 as the reference. Marker sequences from the oat consensus genetic linkage map that  
532 was derived from ten populations<sup>36</sup> were aligned to the final SFS assemblies to  
533 evaluate the consistency between the Hi-C map and the genetic map (Extended Data  
534 Fig. 2). The completeness of the assembled genomes was evaluated with the BUSCO  
535 (v3.1.0)<sup>37</sup> pipeline.

### 536 **Subgenome assignment**

537 A reference-guided strategy based on subgenome homology was used to distinguish  
538 the subgenomes of *A. insularis* and SFS. The *A. longiglumis* reference genome was  
539 split into 100 bp markers and mapped onto the hexaploid genome assemblies. Unique  
540 mapped markers were retained (Extended Data Fig. 3a). A syntenic block was defined  
541 based on the presence of at least five synteny markers. The chromosomes with the  
542 highest homology to *A. longiglumis* were assigned to the A subgenomes, and the  
543 chromosomes with the second-highest homology were assigned to the D subgenomes  
544 because of the high homology between the *Avena* A and D genomes according to  
545 previous studies; the remaining chromosomes with the lowest homology were  
546 accordingly assigned to the C subgenomes (Extended Data Fig. 3c). The  
547 chromosomes of *A. insularis* were aligned with those of SFS and assigned to C and D  
548 subgenomes (Extended Data Fig. 3b, d). The subgenome assignments were further

549 validated by two independent approaches. First, trimmed short reads from *A.*  
550 *longiglumis* and *A. insularis* were mapped to the reference SFS genome, and the depth  
551 coverage of paired reads from these two species was quantified (Extended Data Fig.  
552 3e-f). Second, the abundance and distribution of two types of DNA repeats, As120a  
553 and Am1, in the *A. insularis* reference and SFS genomes were investigated; these  
554 repeat types have been reported to be overrepresented in the *Avena* A and C genomes,  
555 respectively (Fig. 1).

### 556 **Repeat and protein-coding gene annotation**

557 Each of the whole-genomes was searched for repetitive sequences including tandem  
558 repeats (TRs) and transposable elements (TEs). TRs were identified by using GMATA  
559 <sup>38</sup> and Tandem Repeats Finder <sup>39</sup>. A species-specific de novo repeat library was  
560 constructed using MITE-Hunter <sup>40</sup>, LTR\_FINDER <sup>41</sup> and RepeatModeler  
561 (<https://github.com/Dfam-consortium/RepeatModeler>). Then RepeatMasker <sup>42</sup> was  
562 adapted to search for TEs in the reference genome against Repbase <sup>43</sup> and the  
563 species-specific *de novo* repeat library.

564 Protein coding genes were predicted using an evidence-based annotation workflow  
565 by integrating evidence from transcriptomic data, homologue searches and ab initio  
566 prediction. Transcriptomic data were generated by performing PacBio Iso-seq using  
567 total RNA isolated from mixed organs. Raw reads were processed with IsoSeq3,  
568 LIMA, and REFINE to identify full-length, nonchimeric circular consensus sequences  
569 (CCSs). The resulting high-quality CCSs were mapped onto the reference genome for  
570 de-redundancy. The nonredundant isoforms were then used to determine the locations  
571 of potential intron-exon boundaries using GeneMarkST <sup>44</sup>. Protein sequences from *A.*  
572 *atlantica*, *A. eriantha*, *Brachypodium distachyon*, *Hordeum vulgare*, *Oryza sativa*, and  
573 *Triticum aestivum* were used as protein evidence. *Ab initio* gene prediction was

574 performed using GeneMark-ET <sup>45</sup> and AUGUSTUS <sup>46</sup> with two rounds of iterative  
575 training. All gene predictions were integrated using the recommended settings of  
576 EVIDENCEModeler (EVM) <sup>47</sup> after removing transposable element-related genes,  
577 pseudogenes and noncoding genes using TransposonPSI <sup>48</sup> with the default settings.

578 Noncoding RNAs (ncRNAs), including microRNAs, small nuclear RNAs, rRNAs  
579 and regulatory elements, were identified using the Infernal <sup>49</sup> program to search  
580 against the Rfam database <sup>50</sup>. The rRNAs, tRNAs and miRNAs were identified using  
581 RNAmmer <sup>51</sup>, tRNAscan-SE <sup>52</sup>, and miRanda v3.0, respectively. Functional  
582 assignments of the predicted protein-coding genes were obtained with BLAST by  
583 aligning the coding regions to sequences in public protein databases, including the  
584 NCBI nonredundant (NR) protein, Kyoto Encyclopedia of Genes and Genomes  
585 (KEGG), Eukaryotic Orthologous Groups of proteins (KOG), Gene Ontology (GO)  
586 and SwissProt databases. The putative domains and GO terms of the predicted genes  
587 were identified using the InterProScan <sup>53</sup> program with the default settings.

### 588 **Phylogenomic analyses of cereal crops**

589 Foughty-three plant species that containing the main cereal crops with high-quality  
590 reference genomes were downloaded to infer the phylogenetic relationships of cereal  
591 crops. For this purpose, single-copy gene families were identified using reciprocal  
592 best hit (RBH)-based methods. In brief, each proteome of 42 species (51 subgenomes,  
593 Supplementary Table 8) was subjected to BLAST searches against *Amborella*  
594 *trichopoda* sequences according to an E-value  $\leq 1e^{-5}$ . The RBHs in each pair were  
595 obtained and the gene families present in all the subgenomes were retained.  
596 Conserved CDS alignments were extracted by using Gblocks <sup>54</sup> and were  
597 concatenated to generate a supermatrix. 4DTV sites were extracted from this  
598 supermatrix and subject to analysis with RAXML (v.8.2.0) <sup>55</sup> to generate the

599 maximum likelihood tree with the GTR+I+ $\Gamma$  model. Divergence times were estimated  
600 using the MCMCTree program in the PAML4.7 package <sup>56</sup>.

601 To explore the evolution of gene families in these cereal crops, orthologous groups  
602 were constructed with OrthoFinder2 (v2.2.7) <sup>57</sup> using default settings based on the  
603 all-vs-all BLASTP results of the 52 subgenomes. The genes that were exclusively  
604 found in each subfamily (>60% species presented) were identified. Significantly  
605 overrepresented GO terms in each group were identified using the R package  
606 “topGO”.

607 To investigate the chromosomal evolution of the *Avena* genomes, representative  
608 species from five subfamilies in Poaceae with chromosome-level genome assemblies  
609 were selected. The synteny between these extant genomes and the reconstructed  
610 AGKs <sup>22</sup> was identified using MCScanX <sup>58</sup> with the default settings, and the identified  
611 syntenic blocks were then used to deduce the homologous relationship between AGK  
612 marker genes and the protein sequences of *Avena* and the related cereal crop species.  
613 The collinearity between species was identified and plotted using MCScanX (python  
614 version).

### 615 **The evolution and allopolyploidization history of oat**

616 For whole-genome sequencing, 14 other *Avena* accessions were chosen, including six  
617 A genome diploids: *A. canariensis* (CN 23017, Ac), *A. damascena* (CN 19457, Ad), *A.*  
618 *longiglumis* (CN58138, Al), *A. strigosa* (CN 88610, As), *A. wiestii* (CN 90217, As)  
619 and *A. nuda* (CN 58062, As); two C genome diploids: *A. clauda* (CN 21388, Cp) and  
620 *A. ventricosa* (CN 21405, Cv); two AB genome tetraploids: *A. barbata* (CN 65538)  
621 and *A. agadiriana* (CN 25869); three CD genome tetraploids: *A. insularis* (INS-4), *A.*  
622 *maroccana* (CN 57945), *A. murphyi* (CN 21989); and one ACD genome hexaploid:  
623 *A. sativa* (Ogle). These accessions represent all genome types found among extant

624 *Avena* species. All sequencing was performed with an Illumina HiSeq X-Ten  
625 instrument, using 400 bp paired-end libraries.

626 Raw reads of the 14 newly sequenced as well as *A. longiglumis* and *A. insularis*  
627 were trimmed using Trimmomatic v 0.40<sup>59</sup> and mapped onto the reference SFS  
628 genome using BWA<sup>60</sup>. The genomic variants were called based on the mapped short  
629 paired-end reads using the call function of BCFtools<sup>61</sup>, and the obtained variants were  
630 further filtered. The qualified SNPs were used as the RAxML input for constructing  
631 the phylogenetic relationships using maximum likelihood and 100 bootstrap replicates.  
632 The identity between SFS and each sequenced accession was calculated by mapping  
633 ~1X clean reads of each accession to the reference SFS genome.

634 All sequenced diploid accessions were further subjected to transcriptome  
635 sequencing. For this purpose, RNA was isolated from each accession independently  
636 from a combination of sample from seven tissues/conditions, including seedlings (two  
637 weeks old), flag leaves at the booting (Zodoks 45) and heading (Zodoks 58) stages,  
638 panicles at the booting (Zodoks 45), heading (Zodoks 50, 58) and grain dough  
639 (Zodoks 83) stages. Then, the seven types of RNA samples were mixed in equal  
640 amounts and sequenced in paired-end, 150-bp reads on an MGI system. The raw reads  
641 were cleaned with Trimmomatic. De novo assembly was performed using Trinity with  
642 the default settings. Coding sequences (CDSs) were predicted using TransDecoder<sup>62</sup>.  
643 Each proteome from these diploid species and five *Avena* reference genomes was  
644 subjected to BLAST searches against *Hordeum vulgare* according to an E-value  $\leq$   
645 1e-5. The RBHs in each pair were obtained and the gene families present in all the  
646 evaluated species were retained. Phylogenetic tree inference and divergence time  
647 estimation were conducted using the same methods described above.

648 The chloroplast genomes of all sequenced taxa were assembled using the clean

649 short reads with NOVOPlasty<sup>63</sup>. Another 26 *Avena* chloroplast genomes previously  
650 published by Fu *et al.*<sup>24</sup> and Liu *et al.*<sup>25</sup> were also downloaded (Supplementary Table  
651 15). Multiple sequence alignments were performed using MUSCLE<sup>64</sup>, and the  
652 identified informative sites were used for phylogenetic tree construction using  
653 RAxML with 100 bootstrap replicates under the GTR+I+ $\Gamma$  evolutionary model, where  
654 the chloroplast genome sequence from *Triticum aestivum* was used as the outgroup.  
655 Divergence times were estimated according to independent rates and JC69 models  
656 using the MCMCTree program in the PAML4.7 package.

### 657 **Synteny and comparative genomics**

658 The subgenome synteny between SFS and *A. insularis* was analysed by plotting the  
659 positions of homologous pairs in the subgenome pairs within the context of 21 and 14  
660 chromosomes using the Circos<sup>65</sup> software (Extended Data Fig. 6a, b). The  
661 chromosome painting of SFS chromosomes with *A. insularis* and *A. longiglumis*, as  
662 well those of *A. insularis* chromosomes with *A. longiglumis* and *A. eriantha* were  
663 performed by aligning genomic chunks (100 bp markers) to their potential ancestors  
664 using BWA with the default settings. The uniquely mapped markers were retained. We  
665 then processed the markers on each chromosome by requiring at least five consecutive  
666 markers supporting homology to the same chromosome. We consolidated each group  
667 of five consecutive potential markers as one confirmed block. These confirmed blocks  
668 with less than 2 Mb were further consolidated as superblocks. The total lengths of the  
669 subgenomes of *A. insularis* and SFS inherited from the ancestral genomes were  
670 calculated by summing the lengths of superblocks.

671 The interchromosomal exchanges between *A. insularis* and SFS after  
672 polyploidization was also analysed by individually mapping reads from *A.*  
673 *longiglumis* and *A. eriantha* to the *A. insularis* reference genome and reads from *A.*

674 *longiglumis*, *A. eriantha* and *A. insularis* to the SFS reference genome. The  
675 single-base depth coverage of the properly paired reads obtained from the *A.*  
676 *longiglumis*, *A. eriantha* and *A. insularis* mapping was calculated using the Mosdepth  
677 program. The median depth within a sliding window (window size: 1 Mb, step size:  
678 0.5 Mb) was calculated and plotted along with the chromosomes of the reference  
679 genome (Fig. 4c, Extended Data Fig. 7a-d).

#### 680 **FISH analysis**

681 Major interchromosomal exchanges between the C and D subgenomes of *A. insularis*  
682 and SFS were detected by using FISH technology with the C genome-specific repeat  
683 Am1 as the probe. Metaphase chromosome preparation<sup>66</sup> and FISH analysis<sup>67</sup> were  
684 performed as previously described (Fig. 4d, Extended Data Fig. 7e).

#### 685 ***Ka/Ks* analysis**

686 One-to-one orthologous gene sets among the genome assemblies for *Hordeum vulgare*,  
687 the A and C diploid progenitors, *A. longiglumis* and *A. eriantha*, and the subgenomes  
688 of *A. insularis* and SFS were identified using OrthoFinder2<sup>57</sup>. A total of 2,767  
689 orthologous gene sets were obtained and then used for nonsynonymous (*Ka*) and  
690 synonymous (*Ks*) rate calculations. For this purpose, the homoeologous gene pair list  
691 was used as input and the protein sequences from each gene pair were aligned using  
692 MUSCLE<sup>64</sup>. PAL2NAL<sup>68</sup> was used to convert the peptide alignment to a nucleotide  
693 alignment and *Ka* and *Ks* values were computed between gene pairs by using Codeml  
694 from PAML4.7 in free-ratio mode. All estimates with *Ks*<0.01 were excluded from  
695 the analysis. The significance of the differences in *Ka/Ks* values between genomes  
696 (subgenomes) was estimated using the Wilcoxon rank-sum test for nonnormal  
697 distributions in R.

698 **Full-length long terminal repeat (LTR) retrotransposon analysis**

699 Full-length LTRs were identified using LTR\_FINDER and classified using  
700 hierarchical methods (Supplementary Note) (Extended Data Fig. 8b). To estimate the  
701 insertion times of full-length LTRs, the 5'- and 3'-LTR sequences of the full-length  
702 LTRs were aligned and used to calculate divergence using distmat<sup>69</sup>. The insertion  
703 times were estimated with the formula  $T=K/2r$ , where  $r$  represents the neutral  
704 mutation rate, which is  $1.3 \times 10^{-8}$  mutations per site per year<sup>70</sup> (Extended Data Fig.  
705 8d).

706 **Gene family size comparisons between the SFS subgenomes and the A and C**  
707 **diploid genomes**

708 To evaluate the relationships between gene family sizes in the hexaploid subgenomes  
709 and the diploid genomes of *A. longiglumis* and *A. eriantha*, gene family were  
710 identified through OrthoFinder2<sup>57</sup>. The gene family sizes were compared between  
711 each diploid genome and each of the SFS subgenomes. The scatter dots and  
712 regression lines were plotted (Extended Data Figure 9a, b).

713 **Gene loss and retention**

714 Orthologues between *A. eriantha* and the C subgenome of *A. insularis* were identified  
715 using RBH-based methods. A sliding window approach with a window size of 100  
716 genes and a step size of 10 genes by using the *A. eriantha* genome as the reference  
717 was employed to calculate the percentage of retained genes in the C and D  
718 subgenomes individually and in *A. insularis* (Extended Data Fig. 9e). The gene  
719 retention rates of the subgenomes of SFS were also estimated and plotted using the  
720 same methods (Extended Data Fig. 9f).

721 **Gene expression analysis**

722 For gene expression analyses, RNA from SFS samples from seven tissues/conditions

723 (described above) of SFS was individually isolated and sequenced with three  
724 biological replicates using an MGISEQ2000 sequencer. In addition, two-week-old  
725 seedlings of SFS were exposed to six different abiotic stresses: heat, cold, drought,  
726 waterlogging, alkaline and salt. RNA was isolated from seedlings after stress  
727 treatment for one week and sequenced with three biological replicates.

728 Gene expression levels in each sample were quantified using the HiSAT2-HTSeq  
729 pipeline. Differentially expressed genes (DEGs) between stress and control conditions  
730 were detected with the edgeR software package <sup>71</sup> (FDR < 0.05 and |log<sub>2</sub>-fold change  
731 (FC)| > 0.5). To analyse differences in the expression patterns of homoeologous genes,  
732 we undertook an initial analysis of expression variation among strict single-copy  
733 homoeologous triplets. To this end, we used MCSanX <sup>58</sup> to detect syntenic blocks, and  
734 13,744 orthologous triads with a single gene copy per subgenome were identified.  
735 Triplet expression vectors were created by concatenating the observed gene  
736 expression values for the A, C, and D homoeologs. Triplets that expressed at least one  
737 homolog across the sampled tissues were summarized in a triplet expression matrix.  
738 The expression values of the triplet expression matrix were log<sub>10</sub> transformed  
739 (log<sub>10</sub>(TPM+1)), and the matrix was subjected to two-dimensional hierarchical  
740 clustering using “hclust” implemented in R with the “average” correlation distance  
741 and clustering. Heatmap visualization was performed using the heatmap.2 command  
742 from the R package gplots.

743 33 Vaser, R., Sovic, I., Nagarajan, N. & Sikic, M. Fast and accurate de novo  
744 genome assembly from long uncorrected reads. *Genome Res* **27**, 737-746,  
745 doi:10.1101/gr.214270.116 (2017).

746 34 Burton, J. N. *et al.* Chromosome-scale scaffolding of de novo genome  
747 assemblies based on chromatin interactions. *Nat Biotechnol* **31**, 1119-1125,

748 doi:10.1038/nbt.2727 (2013).

749 35 Alonge, M. et al. RaGOO: fast and accurate reference-guided scaffolding of  
750 draft genomes. *Genome Biol* 20, 224. doi:10.1186/s13059-019-1829-6 (2019).

751 36 Bekele, W. A., Wight, C. P., Chao, S., Howarth, C. J. & Tinker, N. A.  
752 Haplotype-based genotyping-by-sequencing in oat genome research. *Plant*  
753 *Biotechnol J* 16, 1452-1463, doi: 10.1111/pbi.12888 (2018).

754 37 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov,  
755 E. M. BUSCO: assessing genome assembly and annotation completeness with  
756 single-copy orthologs. *Bioinformatics* 31, 3210-3212,  
757 doi:10.1093/bioinformatics/btv351 (2015).

758 38 Wang, X. & Wang, L. GMATA: An integrated software package for  
759 genome-scale SSR mining, marker development and viewing. *Front Plant Sci*  
760 7, doi:10.3389/fpls.2016.01350 (2016).

761 39 Benson, G. Tandem repeats finder: a program to analyze DNA sequences.  
762 *Nucleic Acids Res* 27, 573-580, doi:10.1093/nar/27.2.573 (1999).

763 40 Han, Y. & Wessler, S. R. MITE-Hunter: a program for discovering miniature  
764 inverted-repeat transposable elements from genomic sequences. *Nucleic Acids*  
765 *Res* 38, e199-e199, doi:10.1093/nar/gkq862 (2010).

766 41 Xu, Z. & Wang, H. LTR\_FINDER: an efficient tool for the prediction of  
767 full-length LTR retrotransposons. *Nucleic Acids Res* 35, W265-W268,  
768 doi:10.1093/nar/gkm286 (2007).

769 42 Bedell, J. A., Korf, I. & Gish, W. MaskerAid: a performance enhancement to  
770 RepeatMasker. *Bioinformatics* 16, 1040-1041,  
771 doi:10.1093/bioinformatics/16.11.1040 (2000).

772 43 Jurka, J. et al. Repbase Update, a database of eukaryotic repetitive elements.

773 *Cytogenet Genome Res* **110**, 462-467, doi:10.1159/000084979 (2005).

774 44 Tang, S., Lomsadze, A. & Borodovsky, M. Identification of protein coding  
775 regions in RNA transcripts. *Nucleic Acids Res* **43**, e78,  
776 doi:10.1093/nar/gkv227 (2015).

777 45 Lomsadze, A., Burns, P. D. & Borodovsky, M. Integration of mapped  
778 RNA-Seq reads into automatic training of eukaryotic gene finding algorithm.  
779 *Nucleic Acids Res* **42**, e119-e119, doi:10.1093/nar/gku557 (2014).

780 46 Stanke, M., Diekhans, M., Baertsch, R. & Haussler, D. Using native and  
781 syntenically mapped cDNA alignments to improve de novo gene finding.  
782 *Bioinformatics* **24**, 637-644, doi:10.1093/bioinformatics/btn013 (2008).

783 47 Haas, B. J. *et al.* Automated eukaryotic gene structure annotation using  
784 EVIDENCEModeler and the Program to Assemble Spliced Alignments. *Genome*  
785 *Biol* **9**, R7, doi:10.1186/gb-2008-9-1-r7 (2008).

786 48 Urasaki, N. *et al.* Draft genome sequence of bitter melon (*Momordica*  
787 *charantia*), a vegetable and medicinal plant in tropical and subtropical regions.  
788 *DNA Res* **24**, 51-58, doi:10.1093/dnares/dsw047 (2017).

789 49 Nawrocki, E. P. & Eddy, S. R. Infernal 1.1: 100-fold faster RNA homology  
790 searches. *Bioinformatics* **29**, 2933-2935, doi:10.1093/bioinformatics/btt509  
791 (2013).

792 50 Griffiths-Jones, S. *et al.* Rfam: annotating non-coding RNAs in complete  
793 genomes. *Nucleic Acids Res* **33**, D121-D124, doi:10.1093/nar/gki081 (2005).

794 51 Lagesen, K. *et al.* RNAmmer: consistent and rapid annotation of ribosomal  
795 RNA genes. *Nucleic Acids Res* **35**, 3100-3108, doi:10.1093/nar/gkm160  
796 (2007).

797 52 Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection

798 of transfer RNA genes in genomic sequence. *Nucleic Acids Res* **25**, 955-964,  
799 doi:10.1093/nar/25.5.955 (1997).

800 53 Hunter, S. *et al.* InterPro: the integrative protein signature database. *Nucleic*  
801 *Acids Res* **37**, D211-D215, doi:10.1093/nar/gkn785 (2009).

802 54 Talavera, G. & Castresana, J. Improvement of phylogenies after removing  
803 divergent and ambiguously aligned blocks from protein sequence alignments.  
804 *Syst Biol* **56**, 564-577, doi:10.1080/10635150701472164 (2007).

805 55 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and  
806 post-analysis of large phylogenies. *Bioinformatics* **30**, 1312-1313,  
807 doi:10.1093/bioinformatics/btu033 (2014).

808 56 Arenas, M., Sánchez-Cobos, A. & Bastolla, U. Maximum-likelihood  
809 phylogenetic inference with selection on protein folding stability. *Mol Biol*  
810 *Evol* **32**, 2195-2207, doi:10.1093/molbev/msv085 (2015).

811 57 Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for  
812 comparative genomics. *Genome Biol* **20**, 238, doi:10.1186/s13059-019-1832-y  
813 (2019).

814 58 Wang, Y. *et al.* MCSanX: a toolkit for detection and evolutionary analysis of  
815 gene synteny and collinearity. *Nucleic Acids Res* **40**, e49,  
816 doi:10.1093/nar/gkr1293 (2012).

817 59 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for  
818 Illumina sequence data. *Bioinformatics* **30**, 2114-2120,  
819 doi:10.1093/bioinformatics/btu170 (2014).

820 60 Li, H. & Durbin, R. Fast and accurate long-read alignment with  
821 Burrows-Wheeler transform. *Bioinformatics* **26**, 589-595,  
822 doi:10.1093/bioinformatics/btp698 (2010).

- 823 61 Li, H. *et al.* The sequence alignment/map format and SAMtools.  
824 *Bioinformatics* **25**, 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).
- 825 62 Haas, B. J. *et al.* De novo transcript sequence reconstruction from RNA-seq  
826 using the Trinity platform for reference generation and analysis. *Nat Protoc* **8**,  
827 1494-1512, doi:10.1038/nprot.2013.084 (2013).
- 828 63 Dierckxsens, N., Mardulyn, P. & Smits, G. NOVOPlasty: de novo assembly of  
829 organelle genomes from whole genome data. *Nucleic Acids Res* **45**, e18,  
830 doi:10.1093/nar/gkw955 (2017).
- 831 64 Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and  
832 high throughput. *Nucleic Acids Res* **32**, 1792-1797, doi:10.1093/nar/gkh340  
833 (2004).
- 834 65 Krzywinski, M. *et al.* Circos: an information aesthetic for comparative  
835 genomics. *Genome Res* **19**, 1639-1645, doi:10.1101/gr.092759.109 (2009).
- 836 66 Yan, H. *et al.* New evidence confirming the CD genomic constitutions of the  
837 tetraploid *Avena* species in the section *Pachycarpa* Baum. *PloS One* **16**,  
838 e0240703, doi:10.1371/journal.pone.0240703 (2021).
- 839 67 Fu, S. *et al.* Oligonucleotide probes for ND-FISH analysis to identify rye and  
840 wheat chromosomes. *Sci Rep* **5**, 10552, doi:10.1038/srep10552 (2015).
- 841 68 Suyama, M., Torrents, D. & Bork, P. PAL2NAL: robust conversion of protein  
842 sequence alignments into the corresponding codon alignments. *Nucleic Acids*  
843 *Res* **34**, W609-612, doi:10.1093/nar/gkl315 (2006).
- 844 69 Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European molecular  
845 biology open software suite. *Trends Genet* **16**, 276-277,  
846 doi:10.1016/s0168-9525(00)02024-2 (2000).
- 847 70 Wicker, T. *et al.* Impact of transposable elements on genome structure and

848 evolution in bread wheat. *Genome Biol* **19**, 103,  
849 doi:10.1186/s13059-018-1479-0 (2018).  
850 71 Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor  
851 package for differential expression analysis of digital gene expression data.  
852 *Bioinformatics* **26**, 139-140, doi:10.1093/bioinformatics/btp616 (2009).

853 **Acknowledgements** This research was supported by China Agriculture Research  
854 System of MOF and MARA (CARS07), the National Natural Science Foundation of  
855 China (32072025, 31801430), the Science and Technology Development Program of  
856 Jilin Province (20200402034NC) and Talent Fund Project of Jilin Province. The  
857 authors thank Mawsheng Chern, Department of Plant Pathology and the Genome  
858 Center, University of California, Davis, California, USA, for improving the writing of  
859 this article. We also thank Agriculture & Agri-Food Canada (AAFC) and Dr. Eric N.  
860 Jellen, Brigham Young University, Provo, Utah, USA, for providing the *Avena*  
861 materials.

862

863 **Author contributions** C.R., Yuanying Peng, F.L. and Y.W. conceived the study. C.R.,  
864 Yuanying Peng and H.Y. provided funding. C.R., Yuanying Peng, H.Y., L.G., P.Z.,  
865 C.W., J.Z., Yun Peng, D.D. and L.W. collected and prepared the tissue samples for  
866 sequencing. Yuanying Peng, C.D., H.Y., L.K. and F.L. led the bioinformatics analyses.  
867 P.Z., K.Y., C.D., Y.X. and Yun Peng conducted transcriptome sequencing and analysis.  
868 K.Y. and X.D. constructed database. X.L. and Ying Li conducted the FISH validation  
869 of chromosome translocation. J.M., M.H. and Yan Li collected the pictures of spikes.  
870 C.D., Yuanying Peng and H.Y. developed the figures. Yuanying Peng, H.Y., L.G., and  
871 C.D. drafted the manuscript, Q.J., J.W., L.W., W.L., H.K., Y.W., Y.Z. contributed to  
872 the writing. Yuanying Peng, H.Y., L.G., C.D. and L.K. contributed equally.

873

874 **Competing interests** The authors declare no competing interests.

875

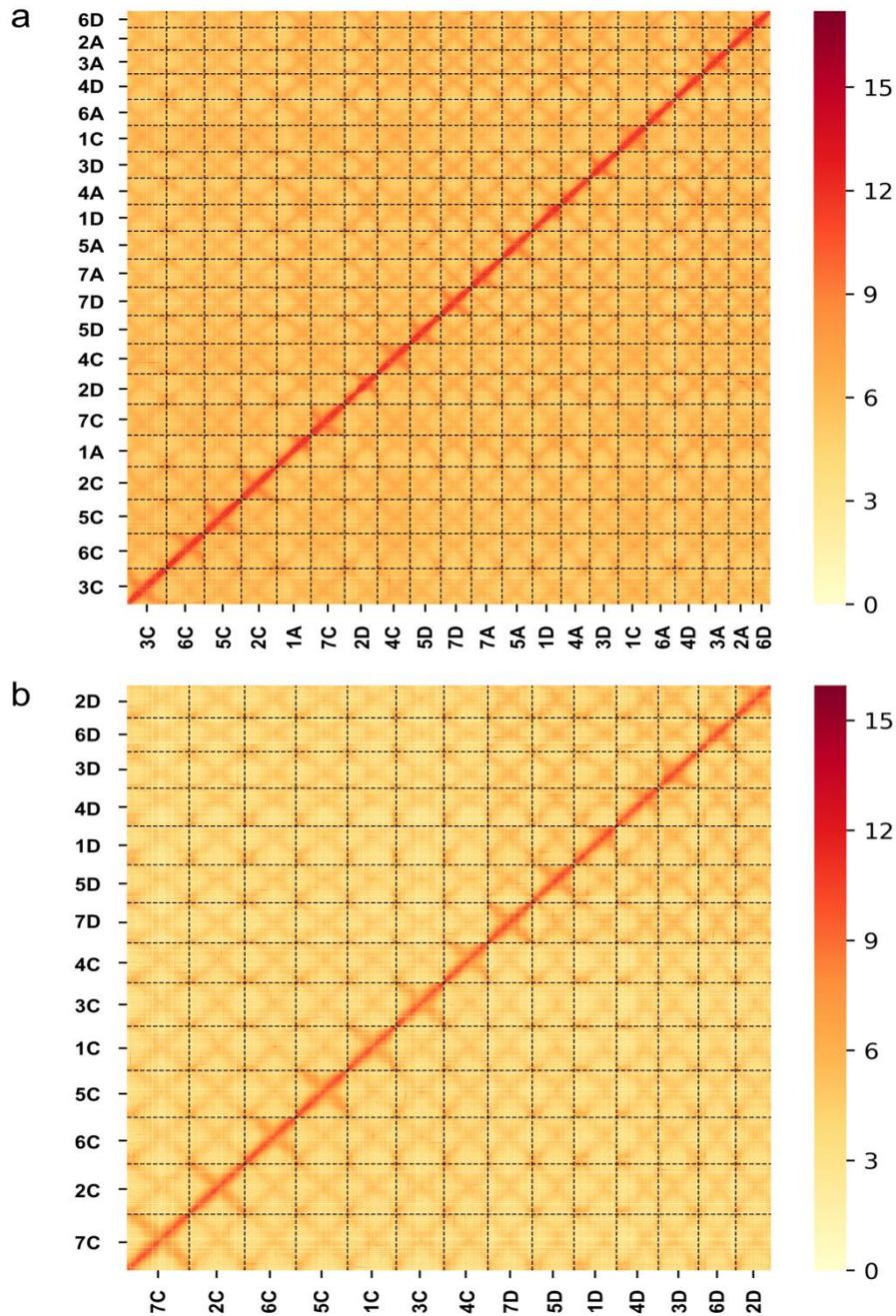
876 **Additional Information**

877 **Supplementary Information** is available for this paper.

878 **Correspondence and requests for materials** should be addressed to R.C.  
879 (renchangzhong@163.com) or Yuanying Peng (yy.peng@hotmail.com).

880 **Data availability**

881 The genome assemblies and sequence data for *A. sativa* ssp. *nuda* cv. Sanfensan, *A.*  
882 *insularis* (CN 108634) and *A. longiglumis* (CN 58139) were deposited at NCBI under  
883 BioProject codes PRJNA 727473, PRJNA731599 and PRJNA716144, respectively.  
884 Sanfensan genome assembly (SAMN19770945), ONT data (SAMN19021785), Hi-C  
885 data (SAMN19340419), NGS data (SAMN19582572), Iso-seq data (SAMN19581880)  
886 and RNA-seq data (SAMN19582573, SAMN19582574); *A.insularis* genome  
887 assembly (SAMN19771048), ONT data (SAMN19291344), Hi-C data  
888 (SAMN19312172), NGS data (SAMN19579880) and Iso-seq data (SAMN19581879);  
889 *A.longiglumis* genome assembly (SAMN19771099), ONT data (SAMN18395928),  
890 NGS data (SAMN19523931) and Iso-seq data (SAMN19581877). All raw data for the  
891 other 14 deep-sequenced accessions including 8 diploids, 5 tetraploids and 1  
892 hexaploid are available under projected numbers can be found in Supplementary  
893 Table 1. All other data are available from the corresponding author upon reasonable  
894 request.



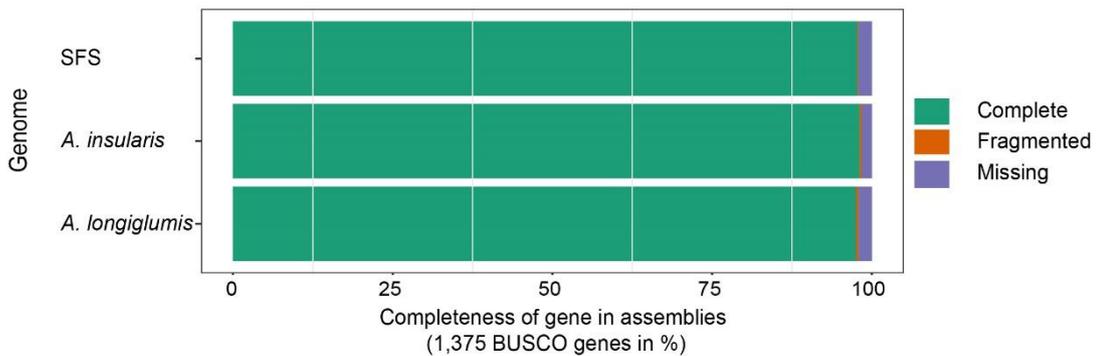
895  
 896 **Extended Data Fig. 1 | Hi-C contact map of *Avena* polyploid species.** Abundant  
 897 intrachromosomal contacts were observed. Chromosomes are sorted by size;  
 898 interchromosomal contacts were also found but with a much lower intensity. **a**, The 21  
 899 chromosomes of hexaploid SFS. **b**, The 14 chromosomes of tetraploid *A. insularis*.

a

		Linkage Group																				
		18	33	23	20	24	5	12	28	13	15	9	3	17	11	1	8	19	21	6	4	2
A. insularis	Total	326	180	209	456	348	317	343	276	454	569	643	745	715	550	694	505	286	552	437	325	559
	1C	248	0	63	1	1	4	3	246	0	11	9	8	5	3	26	2	2	2	3	1	1
	2C	1	15	1	5	0	14	1	0	403	6	3	8	4	17	1	12	5	2	134	1	2
	3C	2	0	9	11	1	1	1	0	2	384	4	1	4	15	1	1	3	1	3	0	1
	4C	2	0	1	31	7	2	2	0	2	2	538	8	4	3	2	0	0	14	0	0	8
	5C	1	0	2	6	32	7	18	4	1	2	11	677	2	4	0	0	1	2	10	0	23
	6C	14	4	0	0	1	24	0	0	1	0	9	6	647	2	1	14	0	3	1	16	7
	7C	1	2	1	12	11	69	16	3	3	1	2	10	5	475	7	5	0	0	8	0	6
	1D	28	9	9	50	64	10	74	8	4	2	3	0	2	8	592	5	17	4	1	4	4
	2D	5	86	5	4	1	28	5	1	12	1	2	1	0	5	4	188	0	4	3	3	5
	3D	0	2	67	25	16	7	13	0	2	13	2	1	1	1	6	4	236	1	3	0	6
	4D	1	4	4	278	7	9	3	0	1	4	23	3	0	2	7	3	6	501	13	2	13
	5D	0	2	38	7	166	9	4	6	2	132	3	6	0	2	8	4	2	1	228	3	6
	6D	20	3	3	8	3	124	10	1	1	3	3	1	18	1	11	249	0	0	2	287	7
7D	1	49	3	7	5	4	188	1	0	1	5	9	9	3	7	8	0	7	1	2	459	

		Linkage Group																				
		18	33	23	20	24	5	12	28	13	15	9	3	17	11	1	8	19	21	6	4	2
SFS	Total	359	261	316	638	457	432	491	301	477	598	699	779	751	575	777	548	302	620	481	341	599
	1A	320	1	1	3	1	0	1	15	2	0	10	2	2	2	126	34	3	1	0	0	6
	2A	0	147	0	1	0	1	1	0	3	1	0	0	0	1	12	0	0	0	0	0	1
	3A	0	0	258	1	2	2	2	11	0	8	1	0	0	0	1	2	4	0	0	0	2
	4A	0	1	0	517	2	3	3	0	1	3	8	0	0	3	14	2	2	54	2	1	2
	5A	0	0	1	3	272	4	1	0	1	0	1	4	0	1	9	0	0	1	51	0	2
	6A	0	0	0	0	345	1	1	1	3	0	0	1	6	21	3	3	0	2	5	45	0
	7A	1	0	2	1	1	2	397	0	0	0	0	5	0	0	18	0	6	0	0	0	30
	1C	2	0	25	0	0	0	0	262	0	4	13	9	0	0	1	1	0	1	2	0	1
	2C	1	4	1	0	0	3	0	0	442	1	3	4	2	0	0	3	0	1	0	0	0
	3C	0	0	7	5	0	0	1	3	1	564	2	2	2	19	0	1	18	2	0	0	2
	4C	0	0	0	11	1	0	1	0	0	4	598	1	0	2	2	0	0	6	0	0	3
	5C	3	0	0	0	11	4	8	3	0	0	2	735	0	0	0	0	0	3	4	0	16
	6C	5	2	0	0	0	9	0	0	1	2	1	4	710	2	0	9	0	0	0	13	3
	7C	0	0	1	2	5	0	6	0	0	0	0	0	1	292	0	1	0	0	3	0	2
	1D	7	2	0	14	24	0	15	2	1	2	1	0	0	3	585	3	12	1	0	0	0
	2D	9	13	1	1	0	6	2	0	10	0	0	3	1	1	4	469	0	0	0	0	2
	3D	0	0	16	5	1	1	6	0	0	3	1	0	1	0	3	1	256	1	0	0	2
	4D	1	2	1	66	1	2	0	0	1	2	24	1	1	0	0	1	0	535	0	1	1
	5D	8	0	0	0	123	4	1	3	2	1	0	2	0	5	0	1	0	0	412	1	1
6D	0	0	0	0	0	13	1	0	1	2	0	0	5	0	0	1	0	0	0	279	3	
7D	1	13	0	2	1	0	42	0	0	0	0	0	2	6	2	4	4	0	1	1	0	

b



900

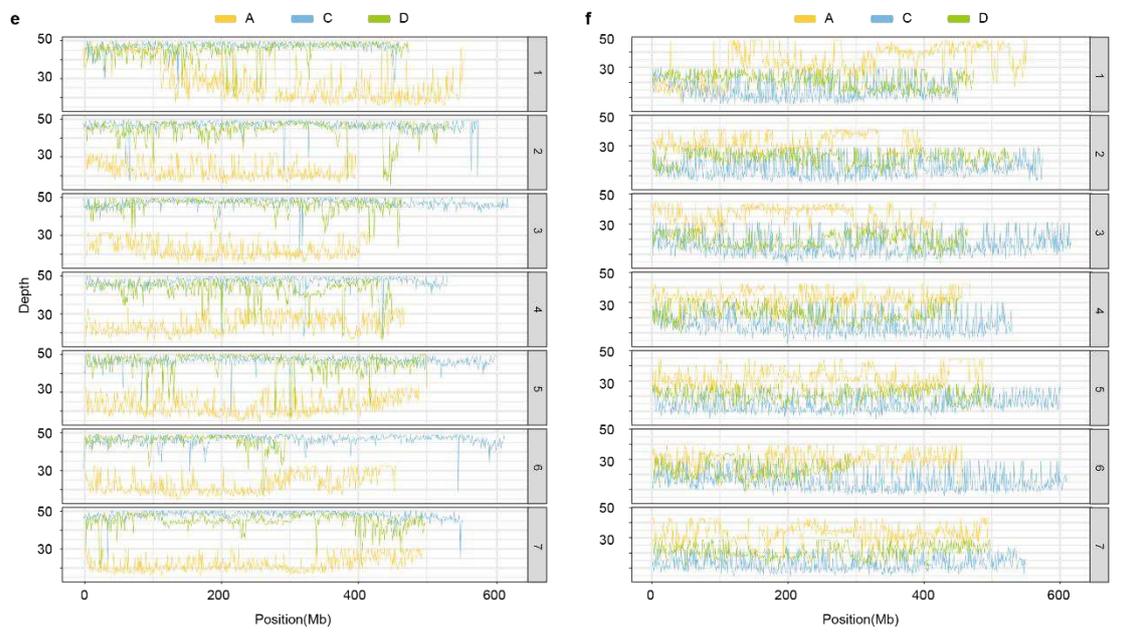
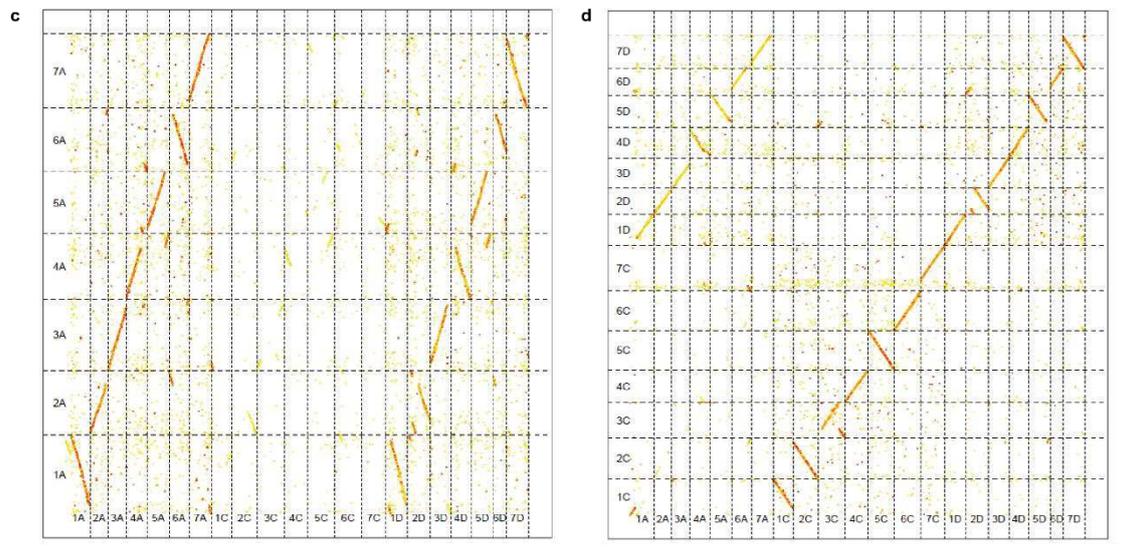
901 **Extended Data Fig. 2 | Assembly validation.** a, The final assemblies were  
 902 evaluated using the consensus hexaploid linkage map. Markers from each of the  
 903 linkage groups of the hexaploid consensus linkage map anchored to the assembled  
 904 SFS (top) and *A. insularis* (bottom) chromosomes. b, Completeness of the three  
 905 assembled genomes as assessed by BUSCO. Bar charts show the percentages of 1,375  
 906 highly conserved plant BUSCO genes that are completely present, fragmented or  
 907 missing in the assembly.

**a**

<i>A. longiglumis</i>							
Chr	1A	2A	3A	4A	5A	6A	7A
1A	238,647	559	5,779	1,637	1,089	2,501	3,007
2A	5,992	139,885	1,001	2,393	1,433	20,747	579
3A	7,159	1,972	219,868	421	2,789	604	2,411
4A	3,010	1,837	11,237	128,677	16,802	23,243	5,654
5A	3,388	1,407	1,307	30,888	162,911	1,961	9,019
6A	526	25,914	18,881	3,529	4,117	11,768	3,735
7A	24,677	991	2,218	3,252	2,378	516	135,213
1C	1,729	92	5,853	140	88	220	184
2C	112	1,318	45	114	86	730	26
3C	188	59	1,483	50	87	242	132
4C	175	48	156	1,475	197	89	127
5C	424	34	88	956	285	106	433
6C	1,417	249	134	193	120	691	2,477
7C	113	32	68	133	383	82	157
1D	84,613	197	831	1,368	15,178	1,391	1,338
2D	8,895	74,898	891	2,693	918	15,723	1,740
3D	4,625	971	56,807	757	630	998	266
4D	1,673	811	4,629	85,668	1,942	13,762	2,482
5D	2,012	318	2,137	18,250	83,570	1,289	937
6D	281	12,199	226	1,137	1,733	54,848	4,203
7D	5,468	590	1,728	2,347	3,687	282	65,588

**b**

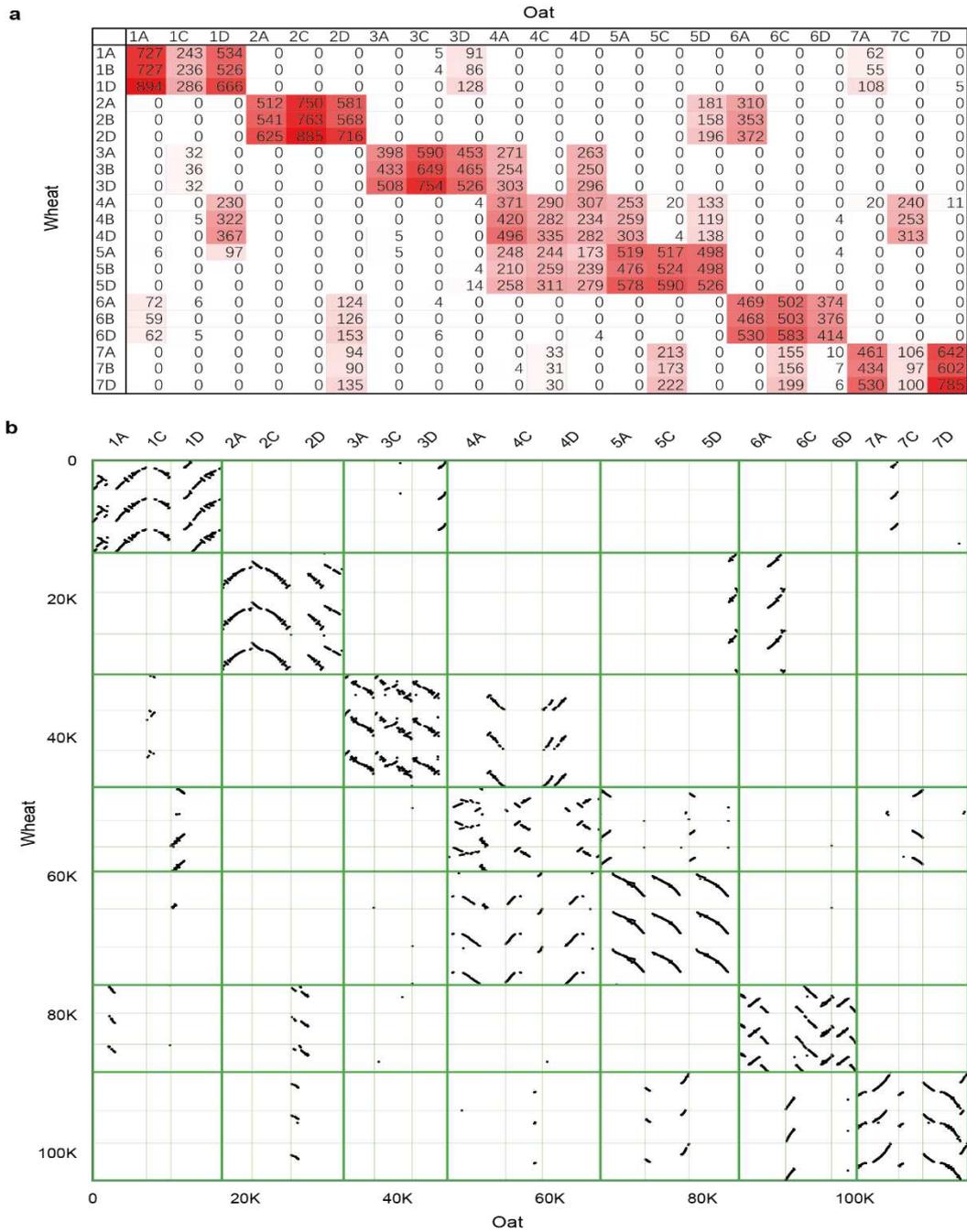
<i>A. insularis</i>														
Chr	1C	2C	3C	4C	5C	6C	7C	1D	2D	3D	4D	5D	6D	7D
1A	88,414	2,507	483	994	111	143	938	37,036	481	100	154	254	4,897	348
2A	6	195	-	16	7	47	113	22	32,095	40	207	112	289	120
3A	4,946	6	51	21	43	-	83	746	244	17,964	193	328	198	145
4A	39	33	92	732	19	27	2,205	3,875	149	1,642	40,595	121	639	143
5A	75	47	95	353	103	22	457	2,063	286	269	770	25,825	194	137
6A	48	400	90	20	30	194	34,146	443	2,150	63	290	354	16,549	168
7A	28	98	11	14	50	-	435	3,605	496	444	232	102	122	26,448
1C	139,252	286	1,182	1,589	696	861	4,648	297	1,027	156	237	153	507	68
2C	1,524	587,218	2,014	610	791	717	2,655	108	413	1,150	160	87	23	68
3C	1,165	338	589,931	2,273	2,743	4,501	2,402	31	217	2,114	531	80,978	84	270
4C	684	1,591	4,624	488,302	2,260	910	794	6	-	20	1,022	80	422	505
5C	2,262	1,403	2,362	507	376,238	1,822	2,290	289	216	648	186	1,195	474	88
6C	2,135	2,518	235	604	2,611	888,768	1,458	26	18	106	104	718	102	189
7C	2,787	2,887	852	2,381	683	721	651,783	-	-	60	34	521	48	12
1D	298	27	36	18	51	6	2,229	457,781	105	332	1,192	2,287	1,120	1,935
2D	183	111	22	63	134	37	1,300	240	365,616	310	775	1,062	84,822	3,214
3D	104	2,614	22	836	22	39	802	4,306	632	468,162	1,135	1,011	20	2,345
4D	117	39	83	496	28	1,309	831	244	256	725	344,828	180	1,760	552
5D	802	47,899	2,114	95	49	195	1,981	1,845	2,196	188	5,894	145,107	2,188	2,185
6D	21	82	226	19	11	167	790	608	1,317	421	1,917	303,994	2,537	2,537
7D	84	83	13	616	87	354	1,125	1,746	800	1,247	2,493	2,132	2,833	460,530



908

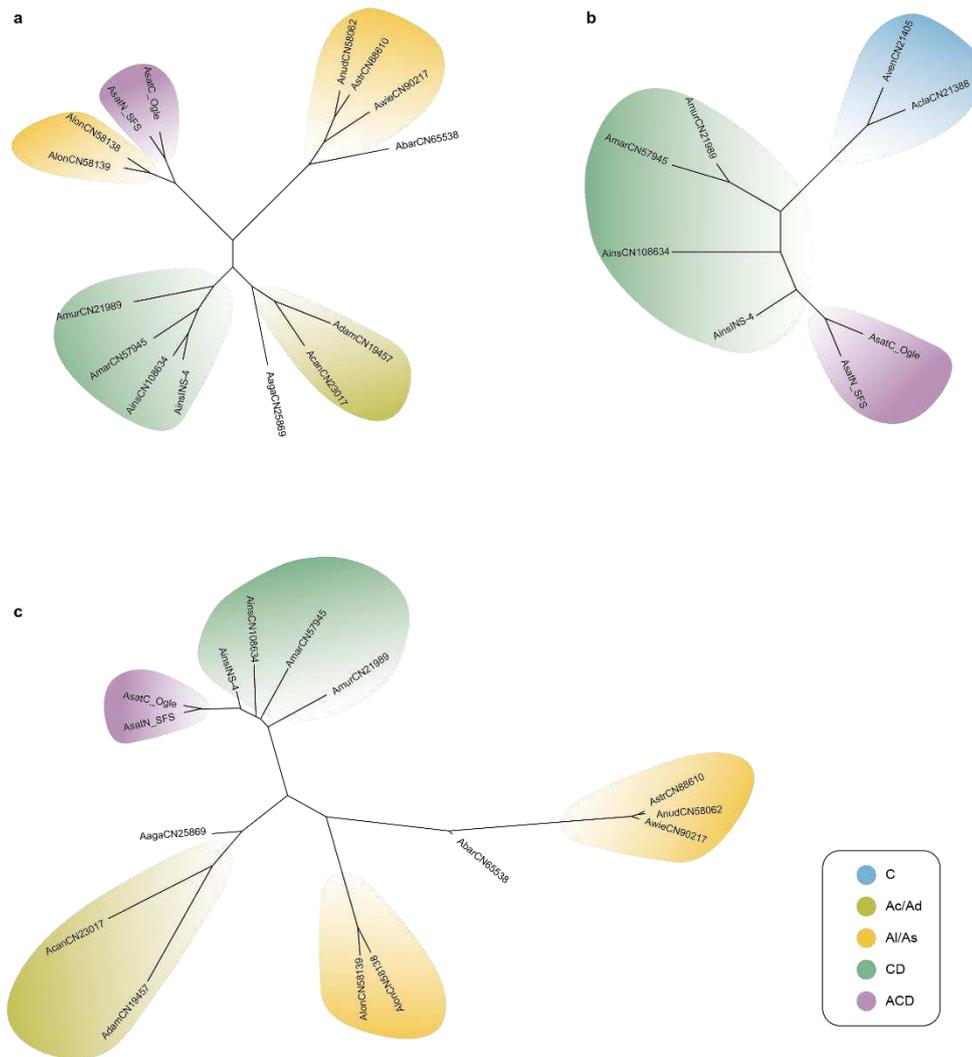
909 **Extended Data Fig. 3 | Subgenome assignment of hexaploid oat.** The genome  
 910 sequences of *A. longiglumis* and *A. insularis* were divided into 100 bp markers and  
 911 then aligned with the reference SFS genome sequence. Unique mapped markers were  
 912 retained. **a**, The greatest numbers of markers shared by the chromosomes of *A.*  
 913 *longiglumis* and SFS are highlighted. **b**, The greatest number of markers shared by the

914 chromosomes of *A. insularis* and SFS are highlighted. **c**, Dot plots display the degrees  
915 of synteny between *A. longiglumis* and SFS. **d**, Dot plots display the degrees of  
916 synteny between *A. insularis* and SFS. Each dot represents a syntenic block of five or  
917 more 100 bp markers. The distance between each pair of adjacent markers is less than  
918 200 kb. **e-f**, Coverage depth obtained along the SFS chromosomes after mapping  
919 Illumina sequence reads from *A. longiglumis* (e) and *A. insularis* (f) to the reference  
920 genome of SFS.



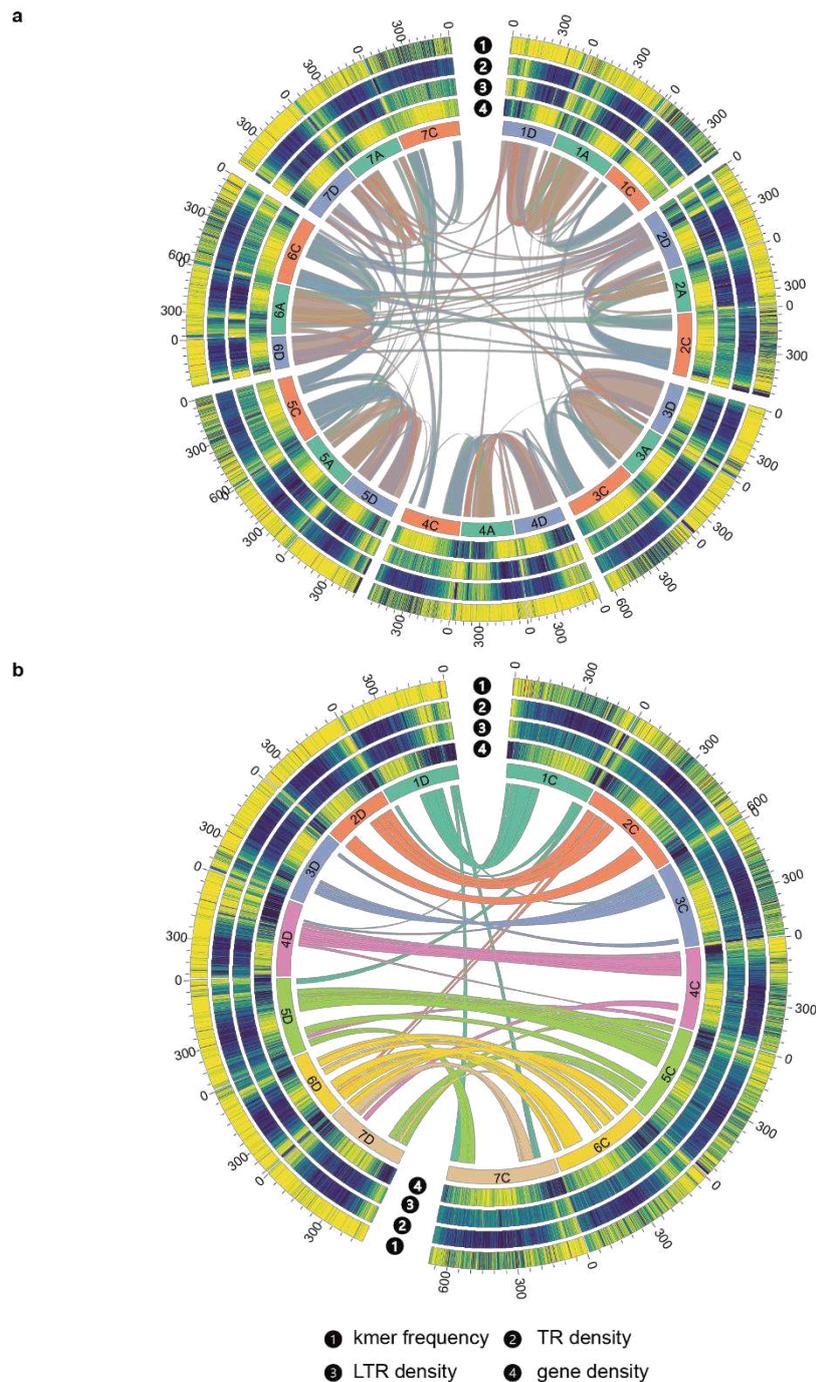
921

922 **Extended Data Fig. 4 | Standard nomenclature for SFS chromosomes.** The  
 923 assignment of SFS chromosomes was performed based on synteny with the  
 924 chromosomes of bread wheat to maintain consistency with future assemblies.  
 925 Chromosome names specify homologous groups (1 to 7) and subgenomes (A/C/D). **a**,  
 926 The matrix shows the number of conserved protein-coding genes between each pair of  
 927 SFS and bread wheat chromosomes. **b**, Dot plots display the degrees of synteny  
 928 between SFS and bread wheat. Each dot represents a syntenic block of five or more  
 929 genes.



930

931 **Extended Data Fig. 5 | Phylogenetic analyses of different *Avena* genomes.** Clean  
 932 short paired reads from all the sequenced *Avena* taxa were aligned to the subgenomes  
 933 of SFS individually, and variants were called and filtered using bcftools. The resulting  
 934 SNPs based on three subgenomes of SFS were used for phylogenetic tree construction.  
 935 **a**, Phylogenetic tree generated from the SNPs of the A subgenome. **b**, Phylogenetic  
 936 tree generated from the SNPs of the C subgenome. **c**, Phylogenetic tree generated  
 937 from the SNPs of the D subgenome.



938

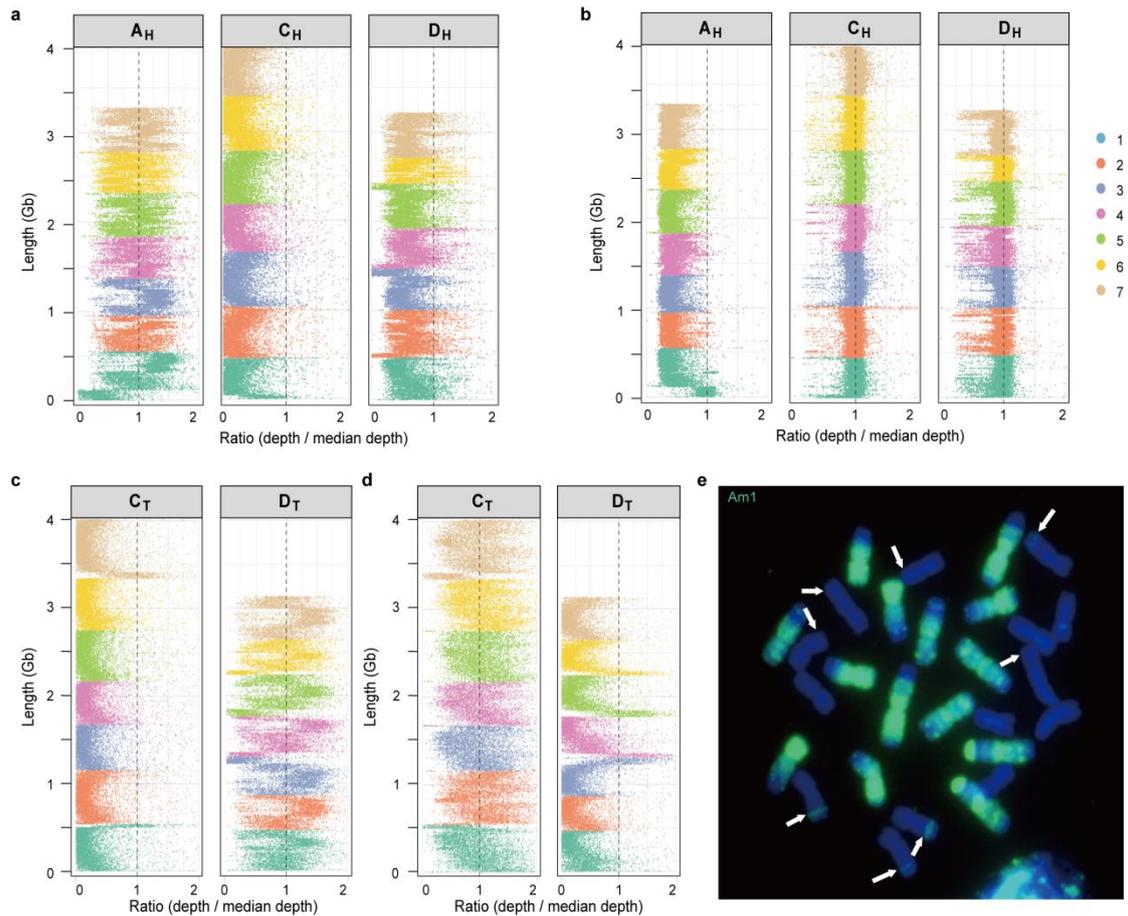
939 **Extended Data Fig. 6 | Structural and conserved synteny landscape of the *A.***

940 *insularis* and SFS subgenomes. Homoeologous gene pairs in syntenic blocks are

941 linked. The rings depict the 31-mer distribution along each chromosome (1), the

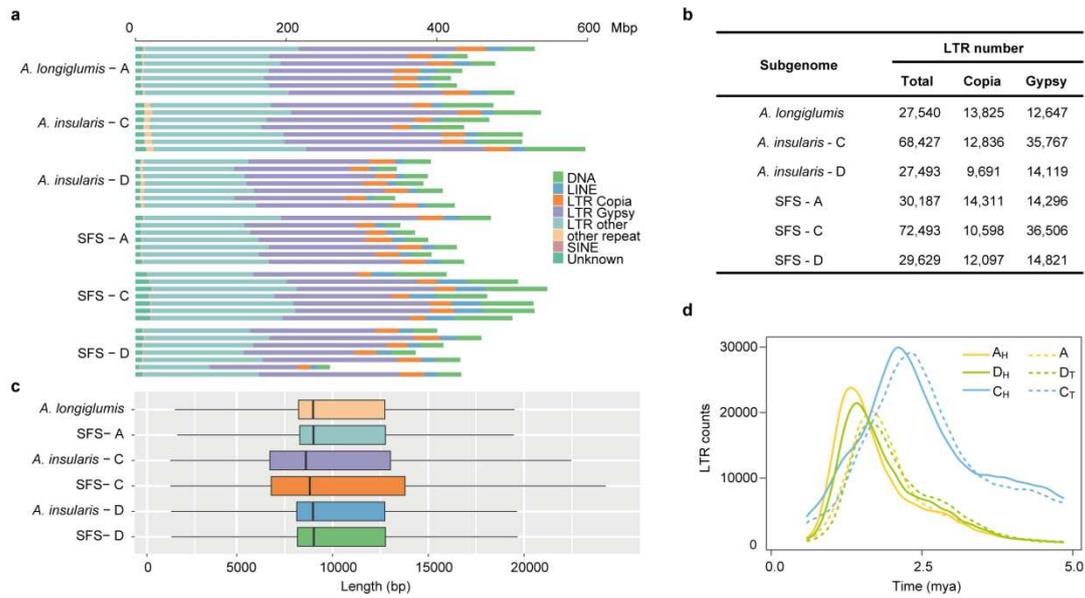
942 density of tandem repeat (TR) (2), the density of long terminal repeats (LTRs) (3), and

943 the density of protein-coding genes (4). **a**, SFS. **b**, *A. insularis*.



944

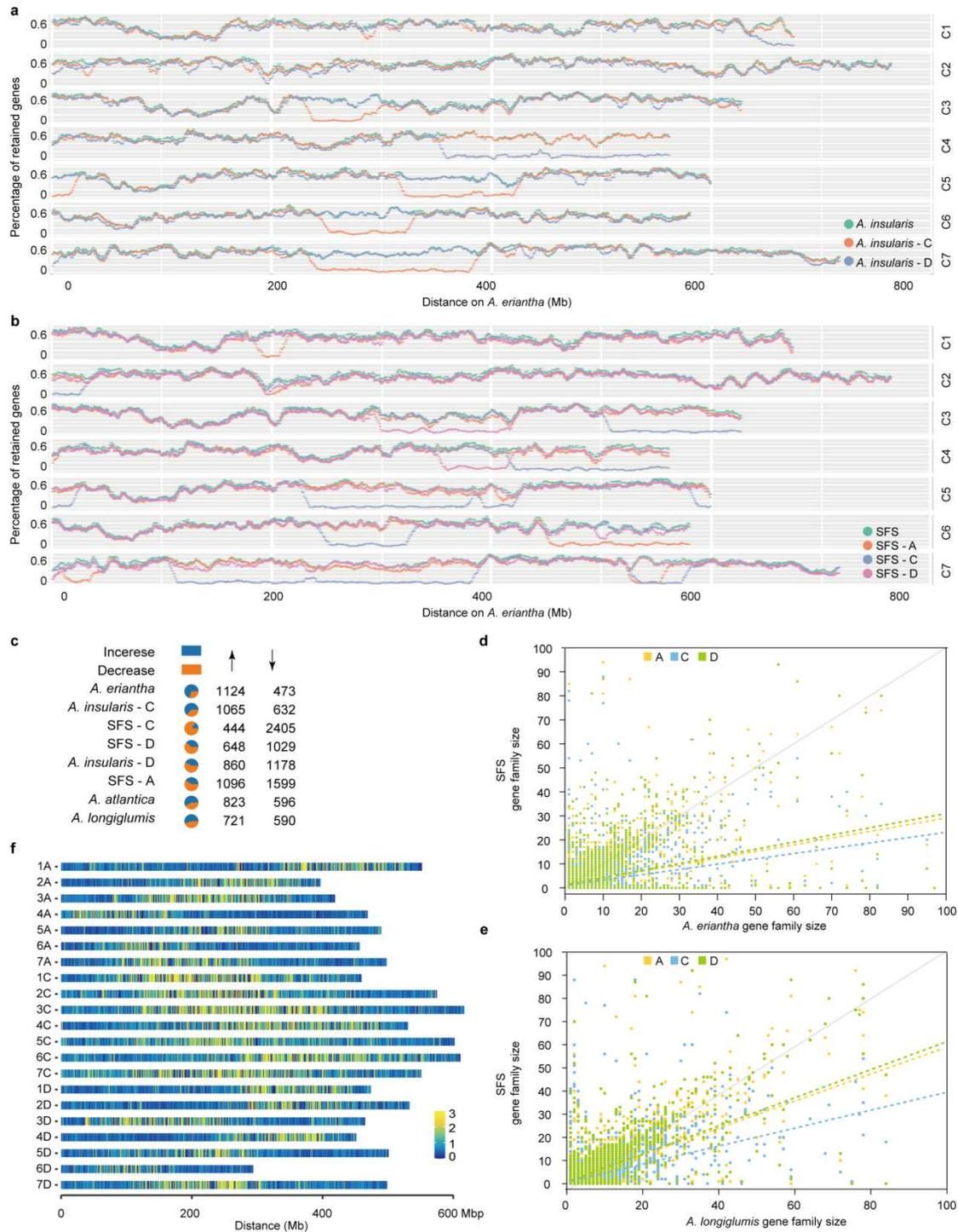
945 **Extended Data Fig. 7 | Intergenomic exchanges between the oat subgenomes. a,**  
 946 Reads of the A genome diploid mapped to the reference SFS genome. **b,** Reads of the  
 947 CD genome tetraploid mapped to the reference SFS genome. **c,** Reads of the A  
 948 genome diploid mapped to the reference *A. insularis* genome. **d,** Reads of the C  
 949 genome diploid mapped to the reference *A. insularis* genome reveal at least four large  
 950 C to D genomic exchanges. **e,** The major C to D translocations (indicated by white  
 951 arrows) in *A. insularis* were confirmed using FISH technology with the C  
 952 genome-specific repeat Am1 (green signals) as the probe.



953

954 **Extended Data Figure 8 | Analysis of TEs in *Avena* genomes.** **a**, Genomic  
 955 constituents of the subgenomes of SFS in comparison with those of the subgenomes  
 956 of *A. insularis*, and the Al (*A. longiglumis*), As (*A. atlantica*), and Cp (*A. eriantha*)  
 957 diploid genomes. **b**, Numbers of full-length long terminal repeats (LTRs) identified  
 958 with LTR\_FINDER in the subgenomes of SFS and *A. insularis*, as well as in the Al,  
 959 As, and Cp genomes. **c**, Average sequence length of LTRs. **d**, Insertion times of  
 960 full-length LTRs in *Avena* genomes.

961



962

963 **Extended Data Fig. 9 | Gene conservation and subgenome fractionation patterns.**

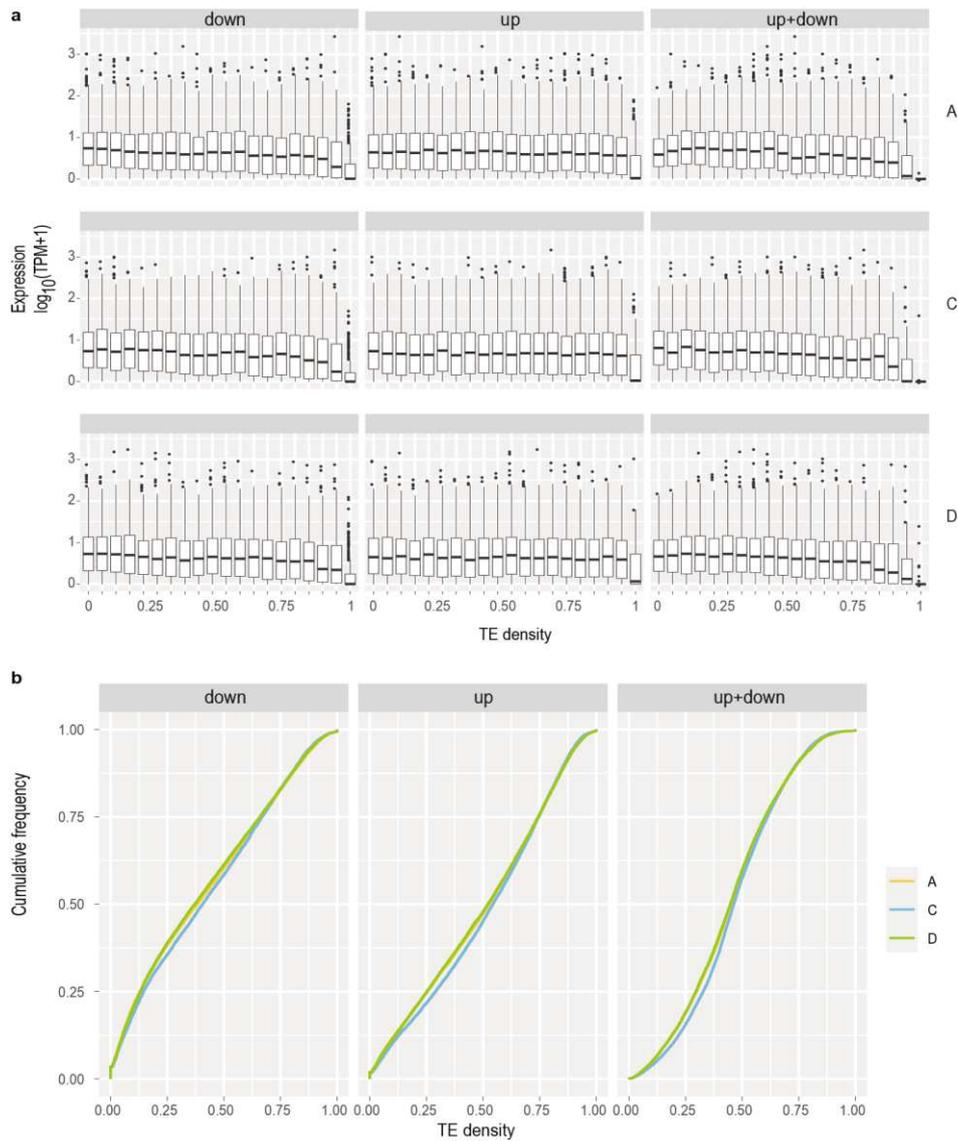
964 **a**, A sliding window approach with a window size of 100 syntenic genes and a step

965 size of 10 syntenic genes was used to show the percentages of retained genes in

966 subgenome C (orange), subgenome D (blue), and both (green) of *A. insularis* using

967 the *A. eriantha* (Cp) genome as a reference. **b**, A sliding window approach with a

968 window size of 100 syntenic genes and a step size of 10 syntenic genes was used to  
969 show the percentage of retained genes in subgenome A (orange), subgenome C (blue),  
970 subgenome D (pink), and all three (green) of SFS using the *A. eriantha* (Cp) genome  
971 as a reference. **c**, Number of expanded and contracted gene families in the  
972 subgenomes of SFS and *A. insularis* and the putative diploid progenitor genomes of *A.*  
973 *longiglumis* and *A. eriantha*. **d**, Relationships between gene family sizes in diploid *A.*  
974 *eriantha* and each subgenome of SFS. **e**, Relationship between gene family sizes in  
975 diploid *A. longiglumis* and each subgenome of SFS. The black line shows a 1:1 gene  
976 copy number relationship for SFS, *A. eriantha*, and *A. longiglumis*, and coloured lines  
977 show the regression fits for the observed gene family size in the subgenomes of SFS  
978 (colours: green, A subgenome; orange, D subgenome; blue, C subgenome). **f**,  
979 Distribution of the identified pseudogenes along the SFS chromosomes.



980

981 **Extended Data Fig. 10 | Impact of TE density on gene expression in SFS. a,** Gene  
 982 expression is negatively correlated with TE density in SFS. Genes were grouped into  
 983 twenty bins (5% increments in TE density). The mean log<sub>10</sub>-transformed gene  
 984 expression of each bin is shown by a boxplot. The central line in each box plot  
 985 indicates the median. The top and bottom edges of the box indicate the first and third  
 986 quartiles and the whiskers extended 1.5 times the interquartile range beyond the edges  
 987 of the box. Black dots represent outliers. **b,** TE density was calculated in 5 kb up- and  
 988 5 kb downstream windows from the gene. TE densities near genes in the C  
 989 subgenome are highest relative to homoeologs in the A and D subgenomes.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryTables12720210627.xlsx](#)
- [SupplementaryInformation20210627.pdf](#)
- [SupplementaryInformation20210627.pdf](#)