

High molecular weight glutenin gene diversity in *Aegilops tauschii* demonstrates unique origin of superior wheat quality

Emily Delorean

Kansas State University

LiangLiang Gao

Kansas State University

Jose Fausto Cervantes Lopez

International Maize and Wheat Improvement Center (CIMMYT)

The Open Wild Wheat Consortium

John Innes Centre

Brande Wulff

John Innes Centre <https://orcid.org/0000-0003-4044-4346>

Maria Itria Ibbá

International Maize and Wheat Improvement Center (CIMMYT)

Jesse Poland (✉ jpoland@ksu.edu)

Kansas State University <https://orcid.org/0000-0002-7856-1399>

Article

Keywords: plant gene diversity, breeding, wheat quality

Posted Date: February 23rd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-230727/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 Communications Biology on November 1st, 2021. See the published version at <https://doi.org/10.1038/s42003-021-02563-7>.

1 **High molecular weight glutenin gene diversity in *Aegilops tauschii***
2 **demonstrates unique origin of superior wheat quality**

3
4 Emily Delorean^{1,4}, LiangLiang Gao^{1,4}, Jose Fausto Cervantes Lopez², Open Wild Wheat Consortium,
5 Brande B. H. Wulff³, Maria Itria Ibba², Jesse Poland^{1,*}

6
7
8
9 ¹ Department of Plant Pathology, Kansas State University, Manhattan, 66506, KS, USA

10 ² Global Wheat Program, International Maize and Wheat Improvement Center (CIMMYT), Apdo Postal 6-
11 641, Mexico, D.F., Mexico

12 ³ The John Innes Centre, Norwich Research Park, Norwich, NR4 7UH, UK

13 ⁴ These authors contributed equally

14 * corresponding author: jpoland@ksu.edu

15 **Abstract**

16

17 Central to the diversity of wheat products was the origin of hexaploid bread wheat, which added the D-
18 genome of *Aegilops tauschii* to tetraploid wheat giving rise to superior dough properties in leavened
19 breads. The polyploidization, however, imposed a genetic bottleneck, with only limited diversity
20 introduced in the wheat D-subgenome. To understand genetic variants for quality, we sequenced 273
21 accessions spanning the known diversity of *Ae. tauschii*. We discovered 45 haplotypes in *Glu-D1*, a
22 major determinant of quality, relative to the two predominant haplotypes in wheat. The wheat allele
23 *2+12* was found in *Ae. tauschii* Lineage 2, the donor of the wheat D-subgenome. Conversely, the
24 superior quality wheat allele *5+10* allele originated in Lineage 3, a recently characterized lineage of *Ae.*
25 *tauschii*, showing a unique origin of this important allele. These two wheat alleles were also quite
26 similar relative to the total observed molecular diversity in *Ae. tauschii* at *Glu-D1*. *Ae. tauschii* is thus a
27 reservoir for unique *Glu-D1* alleles and provides the genomic resource to begin utilizing new alleles for
28 end-use quality improvement in wheat breeding programs.

29

30 Introduction

31

32 Originating in the Fertile Crescent some 10,000 years ago, hexaploid wheat (*Triticum aestivum*) is now
33 grown and consumed around the world ¹. The global consumption of wheat as a staple crop is owed
34 principally to the unique viscoelastic properties of wheat dough that lend it the capacity to make diverse
35 baked products such as leavened bread, tortillas, chapati, pastries, and noodles. The uniqueness of
36 wheat dough can also be described as the strength to resist deformation and elasticity to recover the
37 original shape as well as the viscosity to permanently deform under persistent stress. Elasticity is
38 important for the product to hold shape, while viscosity allows the dough to be worked and formed.
39 The balance of the competing properties determines what baked goods a dough is suitable for, such as a
40 dough with greater strength for leavened pan bread compared to the more extensible dough that is
41 desired for a chapati or tortilla.

42

43 Bread wheat is an allohexaploid with the A-, B- and D-subgenomes contributed by different, but related,
44 species. The closest relative to the wheat A-subgenome is diploid *Triticum urartu*, with other diploid A-
45 genome species including the wild and domesticated Einkorn wheat (*Triticum monococcum*). While the
46 exact ancestor of the B-genome is unknown and presumed extinct, it is believed that *Ae. speltoides* (S-
47 genome) is the closest living relative. These two species were brought together to form a tetraploid
48 wheat species with AABB genome composition, which is known as durum or pasta wheat (*Triticum*
49 *durum*). The D genome from *Aegilops tauschii* was the most recent addition forming the hexaploid
50 genome. This addition of the D-subgenome, to form hexaploid wheat, led to a much broader adaptation
51 and superior bread making quality compared to the tetraploid and diploid ancestors ². However, the
52 original hexaploid species originated from very few *Ae. tauschii* accessions and limited subsequent
53 cross-hybridization likely caused by ploidy barriers with the diploid *Ae. tauschii* ³. This genetic
54 bottleneck resulted in limited genetic diversity in the wheat D-subgenome ⁴.

55

56 The utility of wheat and the variation of wheat products and consumption is driven by the strength and
57 elasticity of the dough which is determined by the structure of the gluten matrix. This matrix is formed
58 from a combination of high-molecular weight (HMW) and low-molecular weight (LMW) glutenin
59 proteins and gliadins ⁵. The backbone of the gluten matrix is developed under mixing by the covalent
60 disulphide bonds between cysteine residues in HMW glutenins ⁶. These glutenins, therefore, are some

61 of the most important genes giving wheat its unique dough properties. They are encoded by a relatively
62 simple locus on the long arm of the group one chromosomes of the Triticeae. Hexaploid wheat,
63 comprised of the A-, B- and D-genomes, thus contains three HMW glutenin loci; *Glu-A1*, *Glu-B1* and *Glu-*
64 *D1*. Each locus harbors two HMW glutenin genes known as the x and y subunit, that are tightly linked
65 but separated by tens to hundreds of kilobase pairs (kb)⁷⁻⁹. Each subunit consists of short, unique N and
66 C terminal domains which flank a central highly repetitive region that accounts for 74-84% of the total
67 protein length¹⁰.

68
69 Allelic differences in all three gluten proteins contribute to the conformation of the gluten matrix and
70 variable end-use quality. The D-subgenome locus, however, is the major driver of bread quality and
71 absence of the D-genome leads to substantially different dough qualities found in tetraploid pasta
72 wheats¹¹. The two common alleles at *Glu-D1* found in bread wheat are *Glu-D1a* (SDS-PAGE allele
73 designation 2+12) and *Glu-D1d* (5+10), with the latter associated with superior breadmaking quality¹²⁻¹⁵.
74 Following the domestication and breeding of wheat, there is limited variation at the *Glu-D1* locus in the
75 D-genome with only these two alleles found in the vast majority of bread wheat throughout the world
76^{16,17}. Of the HMW glutenin alleles on the three sub-genomes, the greatest impact on end-use quality is
77 imparted by the *Glu-D1* locus. Thus, the addition of the wheat D-subgenome and specifically variation at
78 *Glu-D1* has substantial impact on wheat quality globally. This is arguably the single greatest defining
79 feature of bread wheat.

80
81 Reflecting the importance of *Glu-D1* in determining the end-use quality of wheat, focus has been given
82 to understanding the variation present in *Ae. tauschii* for this locus. Much of the work has utilized
83 sodium dodecyl sulfate polyacrylamide gel electrophoresis (SDS-PAGE) protein analysis of *Ae. tauschii*
84 collections¹⁸⁻²¹. From this work, over 37 SDS-PAGE *Glu-D1* alleles have been named in *Ae. tauschii*.
85 However, due to the limited resolution of SDS-PAGE, many of the alleles have indistinguishable SDS-
86 PAGE mobilities from the common *Glu-D1* hexaploid alleles, 2+12 and 5+10, or are difficult to reliably
87 distinguish. By changing the polyacrylamide percentage or acidity in the SDS-PAGE, it was shown that
88 the *Ae. tauschii* 2+12 and 5+10 alleles were slightly different than the common wheat alleles²². These
89 *Ae. tauschii* alleles are therefore given the designations 2t+12t and 5t+10t. In addition to the 2t+12t and
90 5t+10t alleles, a large number of SDS-PAGE alleles have been described, supporting the hypothesis that
91 *Ae. tauschii* could be a vast resource for untapped diversity at *Glu-D1* and that this diversity could be
92 utilized for wheat quality improvement.

93

94 Here we characterized the *Glu-D1* allelic diversity in a panel of 273 sequenced *Ae. tauschii* accessions.
95 The panel spans the known genetic diversity of *Ae. tauschii* and is a powerful resource for association
96 mapping and gene identification²³. From the sequenced *Ae. tauschii* panel, we discovered hundreds of
97 genetic variants which defined dozens of unique haplotypes. This gives the needed molecular
98 information to track these alleles in breeding germplasm, which will in turn enable targeted assessment
99 of the novel *Ae. tauschii* HMW glutenin alleles in hexaploid backgrounds leading to utilization of
100 favorable alleles for wheat quality improvement.

101

102

103 **Results and Discussion**

104

105 **Molecular diversity of *Glu-D1* in *Ae. tauschii***

106 Through the Open Wild Wheat Consortium, we obtained Illumina 150 bp paired-end short reads from
107 234 unique *Ae. tauschii* accessions each sequenced to greater than 7-fold coverage²³. These were
108 aligned to the *Ae. tauschii* AL8/78 reference genome and sequence variants at the annotated *Glu-D1*
109 locus were extracted. We also included three wheat cultivars in this analysis to compare *Ae. tauschii*
110 variants to the common 5+10 (variety 'CDC Stanley') and 2+12 (varieties 'Chinese Spring' and
111 'LongReach Lancer') alleles. From this panel, we identified a total of 310 variants at *Glu-D1*, which were
112 used to generate haplotypes and evaluate molecular diversity at this locus.

113

114 From the *Ae. tauschii* germplasm collection we identified 32 and 33 haplotypes within the coding
115 sequence for the x and y subunits of the *Glu-D1* locus, respectively (Figure 1, Supplemental File S1).
116 When considering the complete *Glu-D1* locus with combination of the x and y subunit, a total of 45
117 haplotypes were identified (Table 1). The various x and y subunit haplotypes were almost exclusively
118 associated with each other, demonstrating the close physical association and limited recombination
119 between the two genes. We included the 2500 bp up- and downstream sequences in our analysis to see
120 if this resulted in further differentiation of alleles as short-read sequences often do not align uniquely to
121 the central, highly repetitive region of the HMW glutenin genes. Including the flanking regions did not
122 result in additional haplotypes. Thus, it appears that the identified variants are sufficient for faithfully
123 differentiating alleles at *Glu-D1*.

124

125 We then calculated genetic distances and determined a gene-level phylogeny at *Glu-D1* for all of the *Ae.*
126 *tauschii* accessions (Figure 1). Haplotypes clustered into three major clades, two of which were
127 associated predominantly with Lineage 2 and one with Lineage 1. A unique group of *Glu-D1* alleles from
128 the newly characterized Lineage 3 accessions were found within a narrow clade with Lineage 2. Among
129 the three major clades, we designated 16 subclades that were clearly distinguished by variants and
130 coincided with a Euclidean distance of 4. Of the 16 subclades, eight were associated exclusively with
131 Lineage 2, five with Lineage 1, and one with Lineage 3. The Lineage 3 accessions all fell within the
132 Lineage 2 major clades, but occupied a unique subclade therein. Thus, the gene-level phylogeny at this
133 locus agrees very closely with the overall previously described population structure of the *Ae. tauschii*
134 lineages^{23,24}. We also observed one clade (9) that had representative accessions from both Lineage 1
135 and 2. This could represent an ancestral haplotype found in both lineages which underwent incomplete
136 lineage sorting, or a case of recent interlineage haplotype exchange. Cases of haplotypes shared across
137 Lineages 1 and 2 were also observed for pest (*Cmc4*) and disease resistance (*Sr46*) genes²³.

138

139 Lineage 2, the recognized ancestral diploid donor of the D-subgenome of hexaploid wheat³, had greater
140 *Glu-D1* molecular haplotype diversity than Lineage 1. Not only were there more subclades associated
141 with Lineage 2, there were also more haplotypes (Supplemental Tables S1 and S2). As expected, the
142 haplotypes of wheat clustered within Lineage 2 subclades (Figure 1). Within Lineage 2, we observed *Ae.*
143 *tauschii* accessions with a matching sequence haplotype to the wheat *2+12* allele consistent with the D-
144 subgenome origin from Lineage 2. Interestingly, the wheat *5+10* allele clustered within the unique
145 Lineage 3 sub-clade. Supporting the inheritance of the *5+10* allele from Lineage 3, Gaurav *et al.* (2021)²³
146 observed genome-wide contribution of Lineage 3 to wheat ancestry. These findings reveal that the
147 Lineage 3 contribution to the wheat D-subgenome included the very valuable *Glu-D1 5+10* allele,
148 arguably one of the most important genes defining the quality of bread wheat.

149

150 Given the large difference in quality between wheat cultivars carrying *2+12* and *5+10* alleles, we
151 hypothesized that these two haplotypes would not be similar at a molecular level. However, we found
152 that *2+12* and *5+10* clustered relatively closely within major-clade III, with much greater overall diversity
153 detected across *Ae. tauschii* particularly when including the Lineage 1 accessions which had very
154 different haplotypes. When comparing the *2+12* and *5+10* haplotypes to those found in Lineage 1, it

155 becomes apparent that *Ae. tauschii* carries alleles that are very unlike anything seen in bread wheat and
156 may offer unique functional characteristics when introgressed into hexaploid backgrounds.

157

158 **Geographic diversity**

159 Given the known geographic structure and distribution of *Ae. tauschii* which is associated with various
160 levels of population structure²⁴, we evaluated the *Glu-D1* diversity relative to the geographic origin of
161 the *Ae. tauschii* accessions. Molecular haplotypes were strongly associated with geographic origin,
162 consistent with the overall genome-wide picture²⁴, and genetic distances between alleles increased
163 with the geographic distance between collection sites of the *Ae. tauschii* accessions (Figure 2). The
164 greatest concentration of haplotype diversity was located along the shores of the Caspian Sea in Iran
165 (Figure 2). Consistent with a hypothesis of admixture between Lineage 1 and Lineage 2 leading to
166 shared gene-level haplotypes across the lineages, the accessions from Lineage 1 and 2 with the same
167 *Glu-D1* haplotype (within subclade 9) were collected very near one another.

168

169 **Molecular haplotypes identify novel *Glu-D1* alleles**

170 We employed SDS-PAGE analysis, the traditional standard for differentiating HMW glutenin loci, to
171 determine if the haplotype molecular sequence diversity would also reflect differences in protein
172 mobility. We evaluated a total of 72 unique accessions with SDS-PAGE and differentiated 9 alleles for
173 the x subunit and 8 alleles at the y subunit from this protein mobility assay. Analysis of the Lineage 1
174 and Lineage 2 variants revealed that molecular haplotypes were consistent with the proteins
175 differentiated by SDS-PAGE (Supplemental Tables 1 and 2). For the majority of the alleles that were
176 differentiated by SDS-PAGE, we were able to unambiguously correlate the observed SDS-PAGE alleles
177 with the molecular variants. Although specific molecular haplotypes were associated with specific SDS-
178 PAGE mobilities, there was little concordance between gene level variation and SDS-PAGE mobility as
179 similar alleles at the molecular level were observed with very different SDS-PAGE mobilities.
180 Alternatively, very different molecular haplotypes were observed with the same SDS-PAGE. This
181 supports our hypothesis that the observed sequence variants are effectively in complete linkage
182 disequilibrium and tagging the size variants from the central repeat region. Similarly, the SDS-PAGE
183 diversity was lower having less differentiating power than the molecular haplotypes. As noted, the
184 same SDS-PAGE mobilities were observed in both Lineage 1 and Lineage 2 haplotypes, but the molecular
185 haplotypes were clearly differentiated (Figure 1). The protein mobility differences are considered to be
186 primarily due to variation in the central repetitive region and therefore are not directly detectable with

187 short-read sequencing, though the variable central repeats are completely phased with diagnostic
188 haplotype variants within the terminal coding regions. Thus, we conclude that a sequence-based
189 resource such as this *Ae. tauschii* panel provides a superior tool for identification and tracking of unique
190 *Glu-D1* alleles in molecular breeding.

191

192 We also examined the connection between the glutenin protein mobility in *Ae. tauschii* compared to
193 hexaploid wheat. *Ae. tauschii* haplotype *Dx1a+Dy1a* matched with the wheat *2+12* haplotype and
194 exhibited the same SDS-PAGE mobility. Although we found an *Ae. tauschii* haplotype identical to the
195 wheat *2+12* allele haplotype, the exact wheat *5+10* haplotype was not detected in this panel, although a
196 very closely related Lineage 3 haplotype was found. Additionally, no *5+10* SDS-PAGE mobilities were
197 observed. This was a surprising observation given that previous studies reported *Ae. tauschii* alleles
198 with a *5+10* SDS-PAGE mobility¹⁸. However, Williams *et al.* (1993)¹⁸ did not reveal the identities of the
199 *Ae. tauschii* accessions with *5+10* SDS-PAGE mobility. Interestingly, the haplotype *Dx7a+Dy7a* in the
200 newly characterized Lineage 3²³ was most similar to *5+10* on the molecular level, however it carried
201 eight variant differences. This current panel, however, only has five unique accessions representing
202 Lineage 3. It is possible therefore that exploration of additional Lineage 3 accessions would reveal a
203 haplotype exactly matching the wheat *5+10* with the same mobility.

204

205 **Cryptic haplotypes**

206 One of the most valuable findings of this study was the high prevalence of cryptic molecular haplotypes
207 hidden within SDS-PAGE mobilities. Within every SDS-PAGE mobility pattern there were multiple
208 molecular haplotypes, often from very different subclades and occasionally from entirely different
209 clades (Figure 1). The cryptic SDS-PAGE haplotypes, accordingly, were geographically disperse (Figure
210 3). For example, within SDS-PAGE *2+12* were four haplotypes; one which was the same as wheat *2+12*
211 (*Dx1a + Dy1a*), another which was within the same subclade (*Dx1c+Dy1d*), and two from entirely
212 different major-clades (*Dx9a+Dy9b* and *Dx13b+Dy13a*). Also, within subclade 9 were the SDS-PAGE
213 mobilities *Dx2+Dy10* and *Dx2+Dy11*, and within subclade 13 were the SDS-PAGE mobilities *1t+12*,
214 *2.1*+12.1**, and *4+10* further supporting that these haplotypes are not all similar to the wheat *2+12*
215 haplotype at the molecular level. However, the proteins still migrate similarly on an SDS-PAGE. These
216 results suggest that SDS-PAGE alone is insufficient when characterizing HMW glutenin diversity in wild
217 relatives and will not be a suitable tool for tracking novel alleles in the hexaploid wheat germplasm.

218

219 While most molecular haplotypes delineated along the three *Ae. tauschii* lineages (Figure 1), a notable
220 exception was within the predominantly Lineage 1 major-clade, subclade 9, where the same three
221 haplotypes ($Dx9a+Dy9a$, $Dx9a+Dy9b$, and $Dx9a+Dy9c$) were observed in both Lineage 1 and Lineage 2
222 accessions. Interestingly, while there were three haplotypes at the y subunit, there was only a single x
223 haplotype associated with all three of these. The x subunit mobility was the same for all three
224 haplotypes, indicating that the x allele is in fact the same. However, the y subunit was differentiated
225 with the mobility $Dy9b$ was faster than that of $Dy9a$ and $Dy10c$ (Supplemental Figure S1).

226

227 **Recombinant haplotypes identified**

228 The close proximity of the glutenin genes results in such tight linkage that recombination is extremely
229 rare. To date, a recombination between the x and y subunit of any HMW-GS locus has yet to be verified.
230 Among the 242 *Ae. tauschii* accessions studied here, we found a clear example of a historical
231 recombination at *Glu-D1* in the accession TA1668 (Lineage 2). SDS-PAGE mobility of TA1668 matches
232 that of TA10081 ($Dx2+Dy10.2$), and though the y haplotype of TA1688 is the same as the y haplotype of
233 TA10081, the x subunit is very different and matching the Lineage 1 clade (Figure 4). Within this clade,
234 the subclade 9 contains both Lineage 1 and Lineage 2 accessions, indicating that there was incomplete
235 lineage sorting or admixture between the two lineages that lead to the introgression of a lineage *Glu-D1*
236 haplotype into the Lineage 2 population. In the presence of both haplotypes, it appears there was a rare
237 recombination between the Lineage 1 and Lineage 2 *Glu-D1* haplotypes, leading to the recombinant
238 haplotype $Dx9a+Dy5a$ found in TA1688.

239 The Lineage 3 accession TA2576 also appears to carry a recombinant haplotype ($Dx7b + Dy15b$) (Figure
240 4). However, our dataset did not contain the exact haplotypes involved in the recombination that led to
241 $Dx7b + Dy15b$. The closest x subunit haplotype is $Dx7a$, the only other Lineage 3 haplotype, from major-
242 clade III and the closest y subunit is the Lineage 2 haplotype $Dy15a$ from major-clade II (Lineage 2). We
243 therefore designated the x and y subunit haplotypes of TA2576 haplotypes within subclades 7 and 15.
244 Geographical analysis reveals that TA2576 was collected from a region shared with other Lineage 3
245 accessions. However, the accessions containing $Dy15a$ haplotype were not collected from a shared
246 region with the L3 accessions. Although not conclusive, the most parsimonious explanation is therefore
247 that $Dx7b + Dy15b$ represents a recombinant haplotype between the x and y subunits from two different
248 alleles. Within our current panel, however, we are unable to differentiate exactly which original
249 haplotypes gave rise to this recombinant haplotype.

250

251 **Conclusions**

252

253 **Importance of Glu-D1 diversity.** The *Glu-D1* locus of wheat provides the greatest contribution to gluten
254 strength, regardless of the allele present¹¹. The allelic diversity of *Glu-D1* in wheat is limited to two
255 predominant alleles *2+12* and *5+10*, and a few rare alleles (*3+12*, *4+10*) which are associated with similar
256 end-use quality as *2+12*^{25,26}. The unique *2.2+12* SDS-PAGE allele, which is found at high frequency in
257 Japanese wheat, was shown to be identical to the *2+12* haplotype with the exception of additional
258 repeats in the internal repeat domain of the x subunit^{25,27,28}. The x subunit protein from *5+10* has a
259 unique cysteine residue just within the central repeat domain which is suspected to increase disulfide
260 bonds in the forming dough. The early expression and greater transcription of this allele is also greater
261 than that of the other *Glu-D1* alleles, in particular *2+12*¹⁴. It is unclear which of these characteristics, or
262 the combination of the two, lend *5+10* the superior quality characteristics. The unique origin of *5+10*
263 from Lineage 3, however, further supports the important contributions of this lineage to the wheat D
264 genome consistent with the findings by Gaurav *et al.* (2021)²³.

265

266 **Unique and valuable sources of diversity.** Our haplotype analysis revealed that the x and y subunits
267 are strongly associated even in diverse germplasm and that the *Glu-D1* haplotypes were clustered to
268 specific geographic origins. Consistent with the findings of Gaurav *et al.* (2021)²³, we found evidence of
269 two lineages (Lineage 2 and Lineage 3) contributing to the D genome of wheat with the superior *5+10*
270 allele found associated with Lineage 3 accessions. Given the excellent end-use quality imparted by
271 *5+10*, understanding this unique origin of the wheat allele support the further exploration and
272 evaluation of novel *Glu-D1* alleles to further improve wheat quality. This also greatly supports the
273 potential of novel alleles and unique haplotypes from the breadth of *Ae. tauschii* diversity.

274

275 Wheat grain quality remains one of the most important targets for breeders to develop superior wheat
276 cultivars. Wild wheat relatives have been shown as a valuable resource for accessing novel genetic
277 diversity to improve a range of wheat breeding targets including yield and disease resistance. For
278 quality evaluation, however, the large quantity of grain needed for milling and baking and the
279 confounding morphological characteristic needed for quality evaluation, such as suitable seed size for
280 milling, make direct evaluation of various end-use quality traits intractable to phenotype directly these

281 wild relatives, including *Ae. tauschii*. In this work, we therefore took the first step in a reverse genetics
282 approach in *Ae. tauschii* by identifying and characterizing variants at the important *Glu-D1* locus. This
283 demonstrated the unique origin of the *Glu-D1* allele in wheat as well as uncovering novel allele variants
284 and haplotypes that can now be targeted for breeding. We also established the relation of wheat alleles
285 to those of *Ae. tauschii* and have shown that *Ae. tauschii* contains a trove of unique *Glu-D1* alleles very
286 unlike the alleles in current wheat germplasm. With accessible germplasm resources such as synthetic
287 hexaploids²³, the diagnostic variants will enable marked-assisted selection of novel *Ae. tauschii*
288 introgressions into wheat, characterization of their end-use quality, and utilization in wheat
289 improvement.

290

291

292 **Methods**

293

294 *Plant Material*

295 This study included 273 *Aegilops tauschii* accessions, of which 241 were from the Wheat Genetics
296 Resource Center (WGRC) collection at Kansas State University in Manhattan, KS, USA. Another 28 were
297 from the National Institute for Agricultural Botany (NIAB) in Cambridge, United Kingdom. An additional 2
298 were from the Commonwealth Scientific and Industrial Research Organisation (CSIRO) in Canberra,
299 Australia. The final accession, AL8/78, was obtained from the John Innes Center (JIC) in Norwich,
300 Norfolk, England. Data regarding original collection sites for the WGRC accessions is detailed in
301 Supplementary Data 1²⁴. *Aegilops tauschii* is divided into two subspecies, spp. *tauschii* (Lineage 1) and
302 the wheat D-genome donor spp. *strangulata* (Lineage 2). In this data set, 117 accessions were Lineage 1
303 and 143 Lineage 2. An additional eight accessions (five non-redundant) belonged to the newly described
304 Lineage 3²³.

305 *SDS-PAGE Analysis*

306 The sodium dodecyl sulphate polyacrylamide gel electrophoresis (SDS-PAGE) analysis of 72 of the *Ae.*
307 *tauschii* accessions, was conducted at the Wheat Chemistry and Quality Laboratory at the International
308 Maize and Wheat Improvement Center (CIMMYT) Texcoco, Mexico according to Singh *et al.* (1991)²⁹
309 with the following modifications. Specifically, 20 mg of whole meal flour were mixed at 1,400 rpm with
310 0.75 ml of 50% propanol (v/v) for 30 min at 65°C in a Thermomixer Comfort (Eppendorf). The tubes
311 were then centrifuged for 2 min at 10,000 rpm, and the supernatant containing the gliadins was
312 discarded. The pellet was then mixed with 0.1 ml of a 1.5% (w/v) DTT solution in a Thermomixer for 30
313 min at 65°C, 1,400 rpm, and centrifuged for 2 min at 10,000 rpm. A 0.1 ml volume of a 1.4% (v/v)
314 vinylpyridine solution was then added to the tube which was subsequently placed again in a
315 Thermomixer for 15 min at 65°C, 1,400 rpm, and centrifuged for 5 min at 13,000 rpm. The supernatant
316 was mixed with the same volume of sample buffer (2% SDS (w/v), 40% glycerol (w/v), and 0.02% (w/v)
317 bromophenol blue, pH 6.8) and incubated in the Thermomixer for 5 min at 90°C and 1,400 rpm. Tubes
318 were centrifuged for 5 min at 10,000 rpm, and 8 ml of the supernatant were used for the glutenin gel.
319 Glutenins were separated in polyacrylamide gels (15% or 13% T) prepared using 1 M Tris buffer, pH of
320 8.5. Gels were run at 12.5 mA for ~19 h. Alleles were identified using the nomenclatures proposed by
321 Payne and Lawrence (1983)¹⁷ for bread wheat high molecular weight glutenins and Lagudah and
322 Halloran (1988)²² for previously described *Ae. tauschii* high molecular weight glutenins.

323 *DNA Sequencing*

324 Whole genome Illumina paired-end sequencing to 10x coverage for most accessions, and 30x coverage
325 for select accessions, was obtained from TruSeq PCR-free libraries with 350 bp insert with Illumina
326 paired end sequencing of 150 bp according to manufactures recommendations. Sequence datasets are
327 detailed in Gaurav *et al.* (2021)²³.

328 *Variant Calling and Duplicated Accessions*

329 Paired-end reads of the *Ae. tauschii* samples were aligned to the *Ae. tauschii* AL8/78 genome assembly
330 (Aet v4.0; NCBI BioProject PRJNA341983, accession AL8/78) and hexaploid wheat samples aligned to an
331 *in silico* reference assembly including the hexaploid wheat A and B genomes from 'Jagger' ³⁰ (Aet v4.0;
332 NCBI BioProject PRJNA341983, accession AL8/78) combined with the Aetv4.0 D genome using HISAT2
333 version 2.1.0 with default parameters ³¹. Alignments were sorted and indexed using samtools v1.9 ³².
334 Variants for coding regions of the x and y subunits of *Glu-D1* were called using bcftools version 1.9 ³³
335 'mpileup' and 'call' commands with a minimum alignment quality of 20 (--q 20) ³⁴. Duplicated accessions
336 were identified as sharing greater than 99.8% variant calls.

337 *Molecular Haplotype Analysis*

338 *Ae. tauschii* and hexaploid wheat variant call format (vcf) files were merged in R and variant calls were
339 recoded to reference (-1) and alternate (1) alleles in R and heterozygous calls were set to missing.

340 Variants were filtered on the following criteria: a variant must be present in either hexaploid wheat or
341 *Ae. tauschii*, must have a quality score greater than 30 and be present in greater than 50% of samples.
342 Given that we expected novel alleles present in single accessions, no minimum minor allele frequency
343 was set. Samples sharing the same variants were considered to share the same molecular haplotype.

344 Genetic distances were calculated as the Euclidean distance on the A matrix of the variants in R. The A
345 matrix was calculated with 'A.mat()' from the rrBLUP package ³⁵ and Euclidean distances with 'dist()'.
346 Hierarchical clustering of the genetic distances were found using hclust() and converted to a
347 dendrogram object before plotting with the dendextend package ³⁶.

348 Molecular haplotypes were designated by the subclade number of the x and y subunits together and
349 then by the letter corresponding to the individual gene level haplotype within. For example, molecular
350 haplotype $x1a + y1b$ represents the a^{th} x haplotype and b^{th} y haplotype within the subclade 1. It should
351 be noted that letter designations across subclades have no correspondence. The ath x haplotype of
352 subclade 1 is different than that of subclade 2.

353 **Supplemental Materials**

354

355 Supplemental Table 1: *Glu-D1* gene positions in Aet v4 assembly

356 Supplemental Table 2: *Glu-D1x* Molecular haplotypes and associated SDS-PAGE alleles

357 Supplemental Table 3: *Glu-D1y* Molecular haplotypes and associated SDS-PAGE alleles

358 Supplemental Figure 1: SDS-PAGE images

359 Supplemental Data 1: Haplotypes and accession info (excel of TAs, SDS-PAGE allele, molecular allele,
360 passport information)

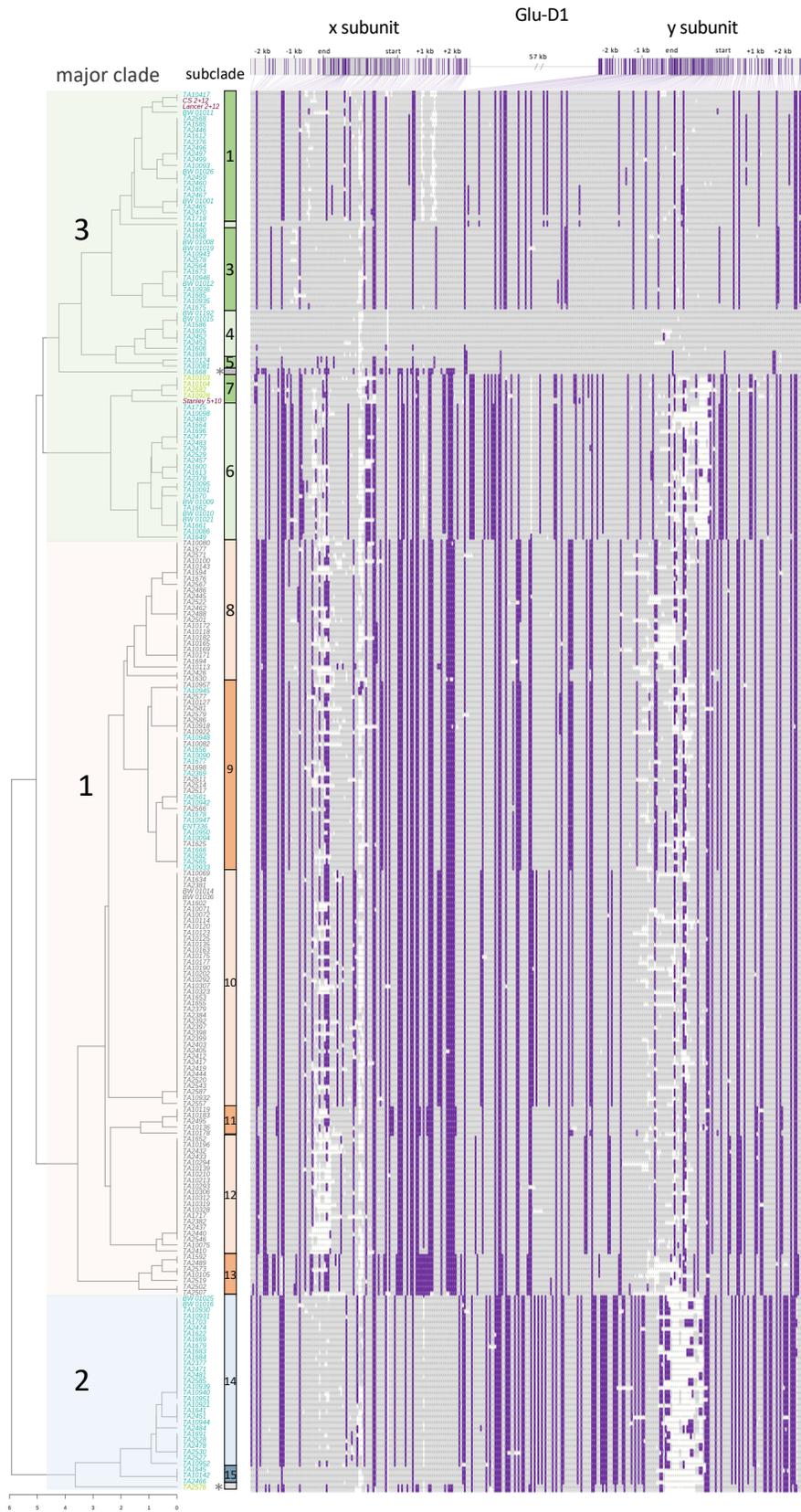
361 Supplemental Data 2: vcf of haplotypes

362

363 Code available in github

364

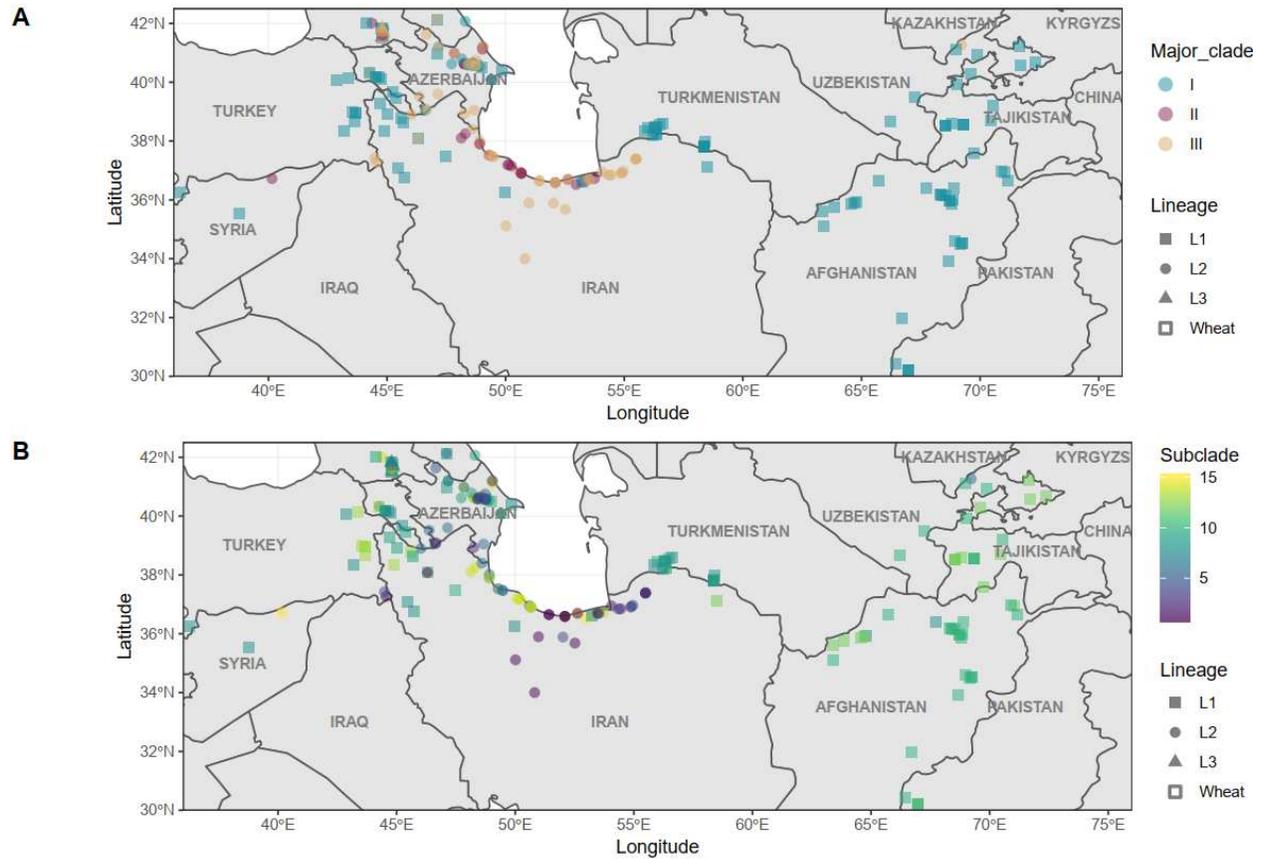
365



367 **Figure 1: Molecular haplotypes of *Glu-D1 Aegilops tauschii* accessions.** Variant positions within the x
368 and y subunit coding sequences and their 2.5 kb flanking sequences are marked with purple bars on
369 schematic of the genes. The molecular haplotypes with position as column and accession as row are
370 shown to the right of the dendrogram, reference allele is in gray and alternate allele is in purple.
371 Dendrogram of combined x and y subunit haplotypes is shown to the left. The corresponding lineage of
372 each accession is colored in blue (Lineage 1), red (Lineage 2) and green (Lineage 3). Haplotypes with
373 major clades and subclades are designated by numbers. The wheat alleles *2+12* and *5+10* are shown in
374 purple and recombinant haplotypes are shown in grey with asterisk.

375

376



378

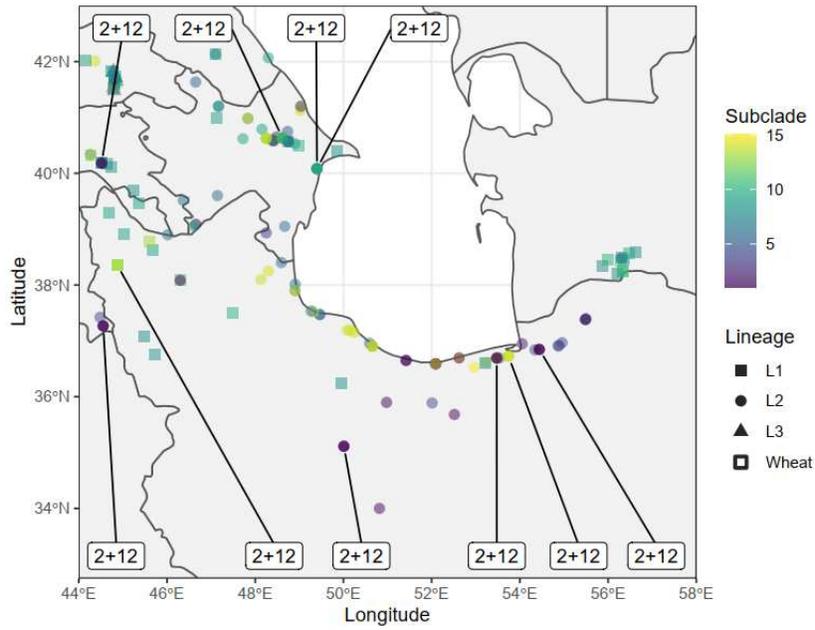
379 **Figure 2: Geographic distribution of *Glu-D1* haplotypes.** Molecular haplotypes for *Glu-D1* shown at the
 380 collection site for the given *Ae. tauschii* accession. **(a)** Distribution of accessions according to *Glu-D1*
 381 major clades with Clade I in blue, Clade II in red and Clade III in orange. **(b)** Distribution of accessions
 382 according to *Glu-D1* haplotype subclades. Haplotypes are shown on a scale from purple to yellow
 383 according to dendrogram order (Figure 1). Lineages are shown as circles for Lineage 1, triangles for
 384 Lineage 2 and squares for Lineage 3.

385

386

387

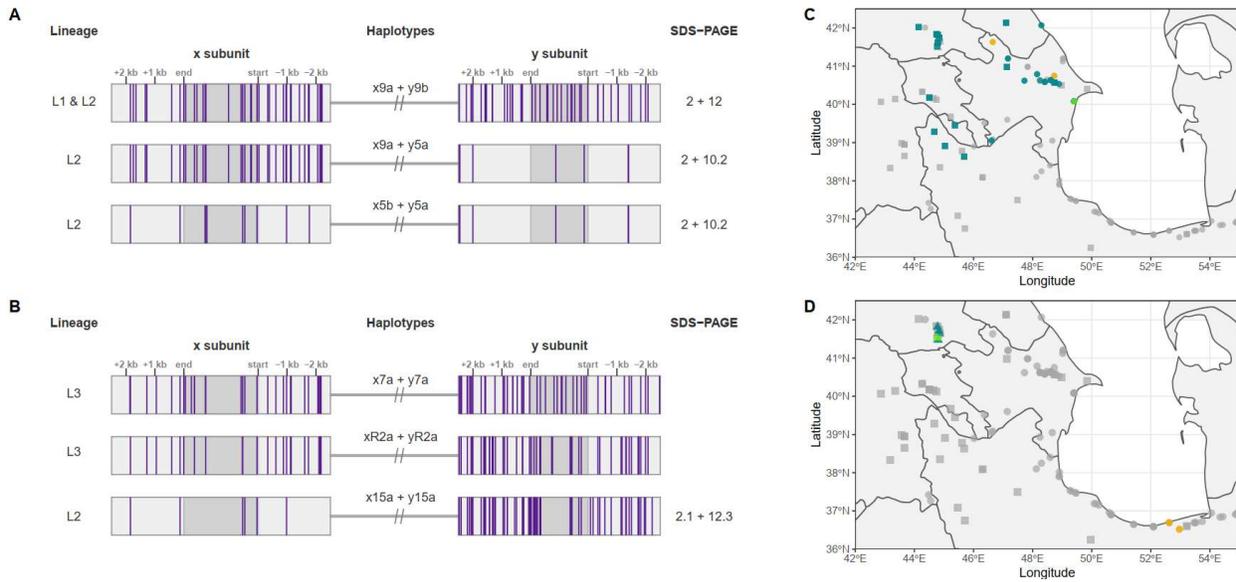
388



389

390 **Figure 3: Cryptic haplotypes within SDS-PAGE alleles.** The molecular haplotypes for *Ae. tauschii*
391 accessions at the sites where the accessions were collected. Corresponding SDS-PAGE allele is noted for
392 *2+12* mobility alleles.

393



395

396 **Figure 4: *Glu-D1* recombinants.** Molecular haplotype representation of recombinant *Glu-D1* haplotypes
 397 for accessions (a) TA1668 (Lineage 2) and (b) TA2576 (Lineage 3). Vertical purple bars represent the
 398 alternate allele variants as called against the AL8 7/8 genome assembly. SDS-PAGE allele for the given
 399 haplotypes are show to right of each gene model. Haplotype *Dx9a* is present in both Lineage 1 and
 400 Lineage 2 accessions. The closest potential x and y subunit haplotypes involved in the recombinant
 401 haplotype *xR2+yR2a* of TA2576 are *x7a* (Lineage 3) and *y15a* (Lineage 2). SDS-PAGE protein mobilities
 402 for *x7a + y7a* and *xR2+yR2a* were not analyzed. Geographical distribution of recombinant *Glu-D1*
 403 haplotypes for (c) TA1668 (Lineage 2) and (d) TA2576 (Lineage 3). Collection site of recombinant
 404 accessions is marked in lime green, whereas turquoise and orange designate the collection sites of
 405 accessions carrying x subunit and y subunit haplotypes, respectively. Accessions with unrelated
 406 haplotypes are in light gray. Lineages are shown in squares (Lineage 1), circles (Lineage 2) or triangles
 407 (Lineage 3).

408

409 **Acknowledgements**

410 ED was supported through the Monsanto's Beachell-Borlaug International Scholars Program. This
411 material is based upon work supported by the National Science Foundation under Award No. 1822162
412 "Phase II IUCRC at Kansas State University Center for Wheat Genetic Resources WGRC" and Award No.
413 1339389 "GPF-PG: Genome Structure and Diversity of Wheat and Its Wild Relatives". BW was supported
414 by the UK Biotechnology and Biological Sciences Research Council Designing Future Wheat Institute
415 Strategic Programme BB/P016855/1. Any opinions, findings, and conclusions or recommendations
416 expressed in this material are those of the author(s) and do not necessarily reflect the views of the
417 National Science Foundation.

418

419 Conflict of Interest: The authors declare no conflict of interest.

420

421

422 **References**

- 423 1 Salamini, F., Özkan, H., Brandolini, A., Schäfer-Pregl, R. & Martin, W. Genetics and geography of
424 wild cereal domestication in the near east. *Nature Reviews Genetics* **3**, 429-441,
425 doi:10.1038/nrg817 (2002).
- 426 2 Dubcovsky, J. & Dvorak, J. Genome plasticity a key factor in the success of polyploid wheat
427 under domestication. *Science (New York, N.Y.)* **316**, 1862-1866, doi:10.1126/science.1143986
428 (2007).
- 429 3 Wang, J. *et al.* Aegilops tauschii single nucleotide polymorphisms shed light on the origins of
430 wheat D-genome genetic diversity and pinpoint the geographic origin of hexaploid wheat. *The*
431 *New Phytologist* **198**, 925-937, doi:10.1111/nph.12164 (2013).
- 432 4 Zhou, Y. *et al.* Triticum population sequencing provides insights into wheat adaptation. *Nat*
433 *Genetics* **52**, 1412-1422, doi:10.1038/s41588-020-00722-w (2020).
- 434 5 Shewry, P. What Is Gluten—Why Is It Special? *Frontiers in Nutrition* **6**,
435 doi:10.3389/fnut.2019.00101 (2019).
- 436 6 Lutz, E., Wieser, H. & Koehler, P. Identification of disulfide bonds in wheat gluten proteins by
437 means of mass spectrometry/electron transfer dissociation. *Journal of Agricultural and Food*
438 *Chemistry* **60**, 3708-3716, doi:10.1021/jf204973u (2012).
- 439 7 Anderson, O., Rausch, C., Moullet, O. & Lagudah, E. The wheat D-genome HMW-glutenin locus:
440 BAC sequencing, gene distribution, and retrotransposon clusters. *Functional & Integrative*
441 *Genomics* **3**, 56-68, doi:10.1007/s10142-002-0069-z (2003).
- 442 8 Kong, X. Y., Gu, Y. Q., You, F. M., Dubcovsky, J. & Anderson, O. D. Dynamics of the evolution of
443 orthologous and paralogous portions of a complex locus region in two genomes of allopolyploid
444 wheat. *Plant Molecular Biology* **54**, 55-69, doi:10.1023/B:PLAN.0000028768.21587.dc (2004).
- 445 9 Gu, Y. Q. *et al.* Types and rates of sequence evolution at the high-molecular-weight glutenin
446 locus in hexaploid wheat and its ancestral genomes. *Genetics* **174**, 1493-1504,
447 doi:10.1534/genetics.106.060756 (2006).
- 448 10 Shewry, P. R., Halford, N. G., Belton, P. S. & Tatham, A. S. The structure and properties of gluten:
449 an elastic protein from wheat grain. *Philosophical Transactions of the Royal Society of London.*
450 *Series B, Biological Sciences* **357**, 133-142, doi:10.1098/rstb.2001.1024 (2002).
- 451 11 Wang, Z. *et al.* New insight into the function of wheat glutenin proteins as investigated with two
452 series of genetic mutants. *Scientific Reports* **7**, 3428, doi:10.1038/s41598-017-03393-6 (2017).

- 453 12 Rooke, L. *et al.* Overexpression of a Gluten Protein in Transgenic Wheat Results in Greatly
454 Increased Dough Strength. *Journal of Cereal Science* **30**, 115-120,
455 doi:<https://doi.org/10.1006/jcrs.1999.0265> (1999).
- 456 13 Lafiandra, D., D'Ovidio, R., Porceddu, E., Margiotta, B. & Colaprico, G. New Data Supporting High
457 Mr Glutenin Subunit 5 as the Determinant of Quality Differences among the Pairs 5 + 10 vs. 2 +
458 12. *Journal of Cereal Science* **18**, 197-205, doi:<https://doi.org/10.1006/jcrs.1993.1046> (1993).
- 459 14 Don, C., Lookhart, G., Naeem, H., MacRitchie, F. & Hamer, R. J. Heat stress and genotype affect
460 the glutenin particles of the glutenin macropolymer-gel fraction. *Journal of Cereal Science* **42**,
461 69-80, doi:<https://doi.org/10.1016/j.jcs.2005.01.005> (2005).
- 462 15 Zhang, P.-P., Ma, H.-X., Yao, J.-B. & He, Z.-H. Effect of Allelic Variation and Expression Quantity at
463 Glu-1 Loci on Size Distribution of Glutenin Polymer in Common Wheat. *Acta Agronomica Sinica*
464 **35**, 1606-1612, doi:[https://doi.org/10.1016/S1875-2780\(08\)60104-2](https://doi.org/10.1016/S1875-2780(08)60104-2) (2009).
- 465 16 Payne, P. I., Holt, L. M. & Law, C. N. Structural and genetical studies on the high-molecular-
466 weight subunits of wheat glutenin : Part 1: Allelic variation in subunits amongst varieties of
467 wheat (*Triticum aestivum*). *Theoretical and Applied Genetics* **60**, 229-236,
468 doi:10.1007/bf02342544 (1981).
- 469 17 Payne, P. I. & Lawrence, G. J. Catalogue of alleles for the complex gene loci, Glu-A1, Glu-B1, and
470 Glu-D1 which code for high-molecular-weight subunits of glutenin in hexaploid wheat. *Cereal*
471 *Research Communications* **11**, 29-35 (1983).
- 472 18 William, M. D. H. M., Peña, R. J. & Mujeeb-Kazi, A. Seed protein and isozyme variations in
473 *Triticum tauschii* (*Aegilops squarrosa*). *Theoretical and Applied Genetics* **87**, 257-263,
474 doi:10.1007/BF00223774 (1993).
- 475 19 Xu, S., Khan, K., Klindworth, D. & Nygard, G. Evaluation and characterization of high-molecular
476 weight 1D glutenin subunits from *Aegilops tauschii* in synthetic hexaploid wheats. *Journal of*
477 *Cereal Science* **52**, doi:10.1016/j.jcs.2010.05.004 (2010).
- 478 20 Mackie, A. M., Lagudah, E. S., Sharp, P. J. & Lafiandra, D. Molecular and biochemical
479 characterisation of HMW glutenin subunits from *T. tauschii* and the D genome of hexaploid
480 wheat. *Journal of Cereal Science* **23**, 213-225, doi:<https://doi.org/10.1006/jcrs.1996.0022>
481 (1996).
- 482 21 Gianibelli, M. C., Gupta, R. B., Lafiandra, D., Margiotta, B. & MacRitchie, F. Polymorphism of high
483 Mr glutenin subunits in *Triticum tauschii*: Characterisation by Chromatography and

484 Electrophoretic Methods. *Journal of Cereal Science* **33**, 39-52,
485 doi:<https://doi.org/10.1006/jcrs.2000.0328> (2001).

486 22 Lagudah, E. S. & Halloran, G. M. Phylogenetic relationships of *Triticum tauschii* the D genome
487 donor to hexaploid wheat. *Theoretical and Applied Genetics* **75**, 592-598,
488 doi:[10.1007/BF00289125](https://doi.org/10.1007/BF00289125) (1988).

489 23 Gaurav, K. *et al.* Evolution of the bread wheat D-subgenome and enriching it with diversity from
490 *Aegilops tauschii*. *bioRxiv*, 2021.2001.2031.428788, doi:[10.1101/2021.01.31.428788](https://doi.org/10.1101/2021.01.31.428788) (2021).

491 24 Singh, N. *et al.* Genomic analysis confirms population structure and identifies inter-Lineage
492 Hybrids in *Aegilops tauschii*. *Frontiers in Plant Science* **10**, doi:[10.3389/fpls.2019.00009](https://doi.org/10.3389/fpls.2019.00009) (2019).

493 25 Wan, Y. *et al.* Comparative analysis of the D genome-encoded high-molecular weight subunits of
494 glutenin. *Theoretical and Applied Genetics* **111**, 1183-1190, doi:[10.1007/s00122-005-0051-y](https://doi.org/10.1007/s00122-005-0051-y)
495 (2005).

496 26 Dong, Z. *et al.* Haplotype Variation of Glu-D1 Locus and the Origin of Glu-D1d Allele Conferring
497 Superior End-Use Qualities in Common Wheat. *PLoS ONE* **8**, e74859,
498 doi:[10.1371/journal.pone.0074859](https://doi.org/10.1371/journal.pone.0074859) (2013).

499 27 Payne, P. I., Holt, L. M. & Lawrence, G. J. Detection of a novel high molecular weight subunit of
500 glutenin in some Japanese hexaploid wheats. *Journal of Cereal Science* **1**, 3-8,
501 doi:[https://doi.org/10.1016/S0733-5210\(83\)80003-4](https://doi.org/10.1016/S0733-5210(83)80003-4) (1983).

502 28 Shimizu, K. K. *et al.* De Novo Genome Assembly of the Japanese Wheat Cultivar Norin 61
503 Highlights Functional Variation in Flowering Time and Fusarium Resistance Genes in East Asian
504 Genotypes. *Plant and Cell Physiology*, doi:[10.1093/pcp/pcaa152](https://doi.org/10.1093/pcp/pcaa152) (2020).

505 29 Singh, N. K., Shepherd, K. W. & Cornish, G. B. A simplified SDS-PAGE procedure for separating
506 LMW subunits of glutenin. *Journal of Cereal Science* **14**, 203-208,
507 doi:[https://doi.org/10.1016/S0733-5210\(09\)80039-8](https://doi.org/10.1016/S0733-5210(09)80039-8) (1991).

508 30 Walkowiak, S. *et al.* Multiple wheat genomes reveal global variation in modern breeding. *Nature*
509 **588**, 277-283, doi:[10.1038/s41586-020-2961-x](https://doi.org/10.1038/s41586-020-2961-x) (2020).

510 31 Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome alignment and
511 genotyping with HISAT2 and HISAT-genotype. *Nature Biotechnology* **37**, 907-915,
512 doi:[10.1038/s41587-019-0201-4](https://doi.org/10.1038/s41587-019-0201-4) (2019).

513 32 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford,*
514 *England)* **25**, 2078-2079, doi:[10.1093/bioinformatics/btp352](https://doi.org/10.1093/bioinformatics/btp352) (2009).

515 33 Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and
516 population genetical parameter estimation from sequencing data. *Bioinformatics (Oxford,
517 England)* **27**, 2987-2993, doi:10.1093/bioinformatics/btr509 (2011).

518 34 Gao, L. (Zenodo, 2020).

519 35 Endelman, J. B. Ridge Regression and Other Kernels for Genomic Selection with R Package
520 rrBLUP. *The Plant Genome* **4**, doi:https://doi.org/10.3835/plantgenome2011.08.0024 (2011).

521 36 Galili, T. dendextend: an R package for visualizing, adjusting and comparing trees of hierarchical
522 clustering. *Bioinformatics* **31**, 3718-3720, doi:10.1093/bioinformatics/btv428 (2015).

523

Figures

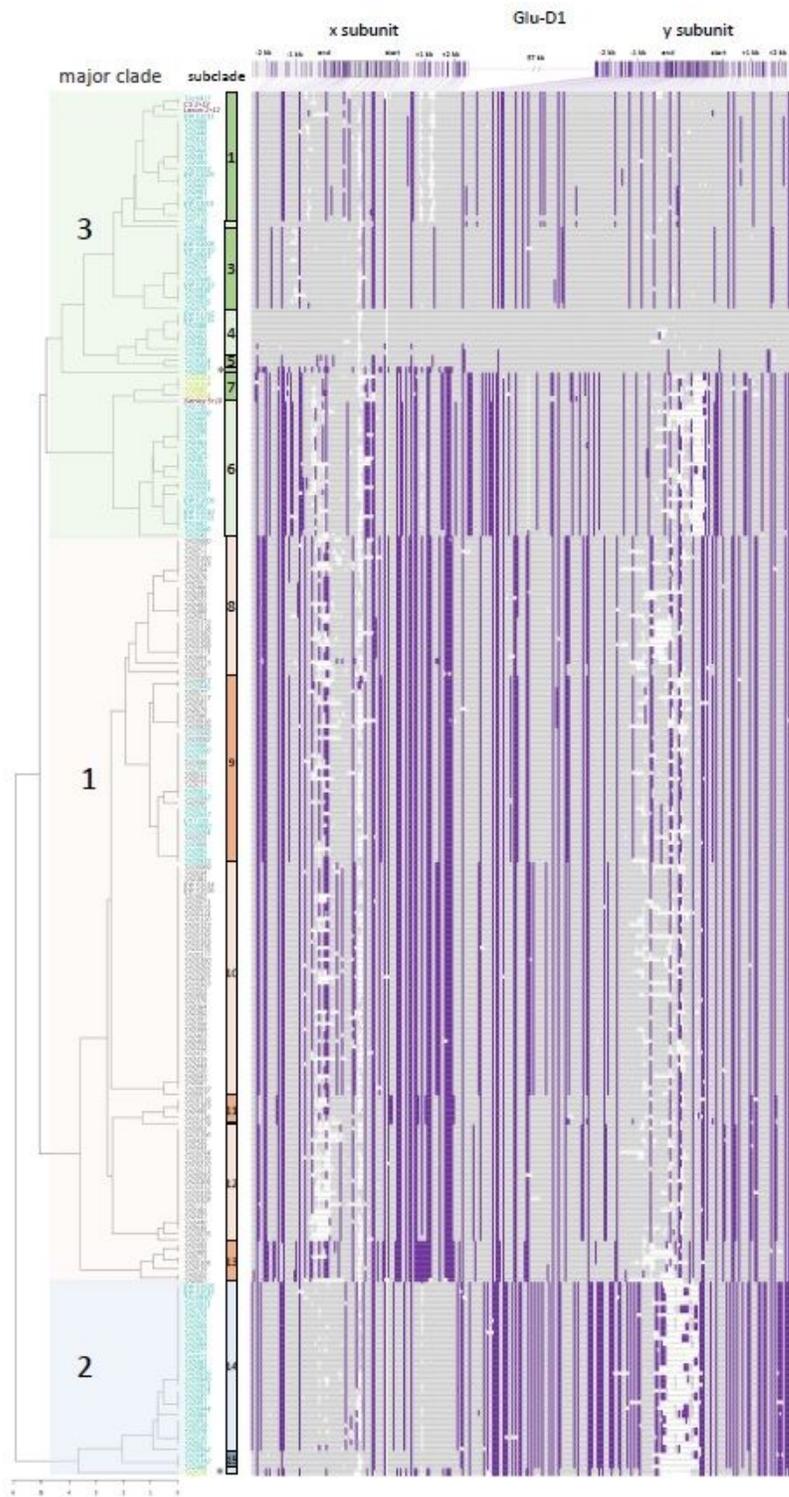


Figure 1

Molecular haplotypes of Glu-D1 *Aegilops tauschii* accessions. Variant positions within the x and y subunit coding sequences and their 2.5 kb flanking sequences are marked with purple bars on schematic of the genes. The molecular haplotypes with position as column and accession as row are shown to the

right of the dendrogram, reference allele is in gray and alternate allele is in purple. Dendrogram of combined x and y subunit haplotypes is shown to the left. The corresponding lineage of each accession is colored in blue (Lineage 1), red (Lineage 2) and green (Lineage 3). Haplotypes with major clades and subclades are designated by numbers. The wheat alleles purple and recombinant haplotypes are shown in grey with asterisk.

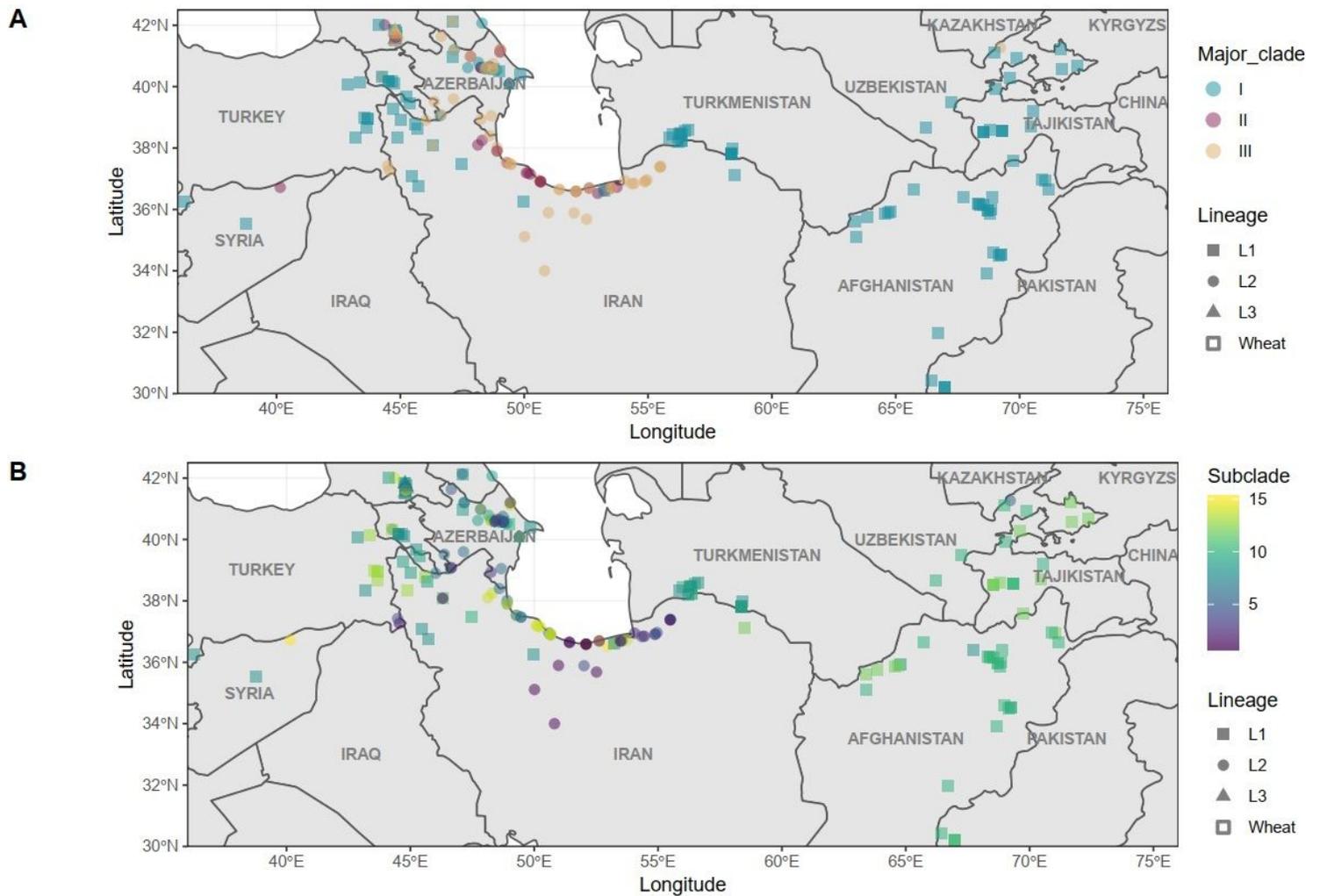


Figure 2

Geographic distribution of Glu-D1 haplotypes. Molecular haplotypes for Glu-D1 shown at the collection site for the given *Ae. tauschii* accession. (a) Distribution of accessions according to Glu-D1 major clades with Clade I in blue, Clade II in red and Clade III in orange. (b) Distribution of accessions according to Glu-D1 haplotype subclades. Haplotypes are shown on a scale from purple to yellow according to dendrogram order (Figure 1). Lineages are shown as circles for Lineage 1, triangles for Lineage 2 and squares for Lineage 3.

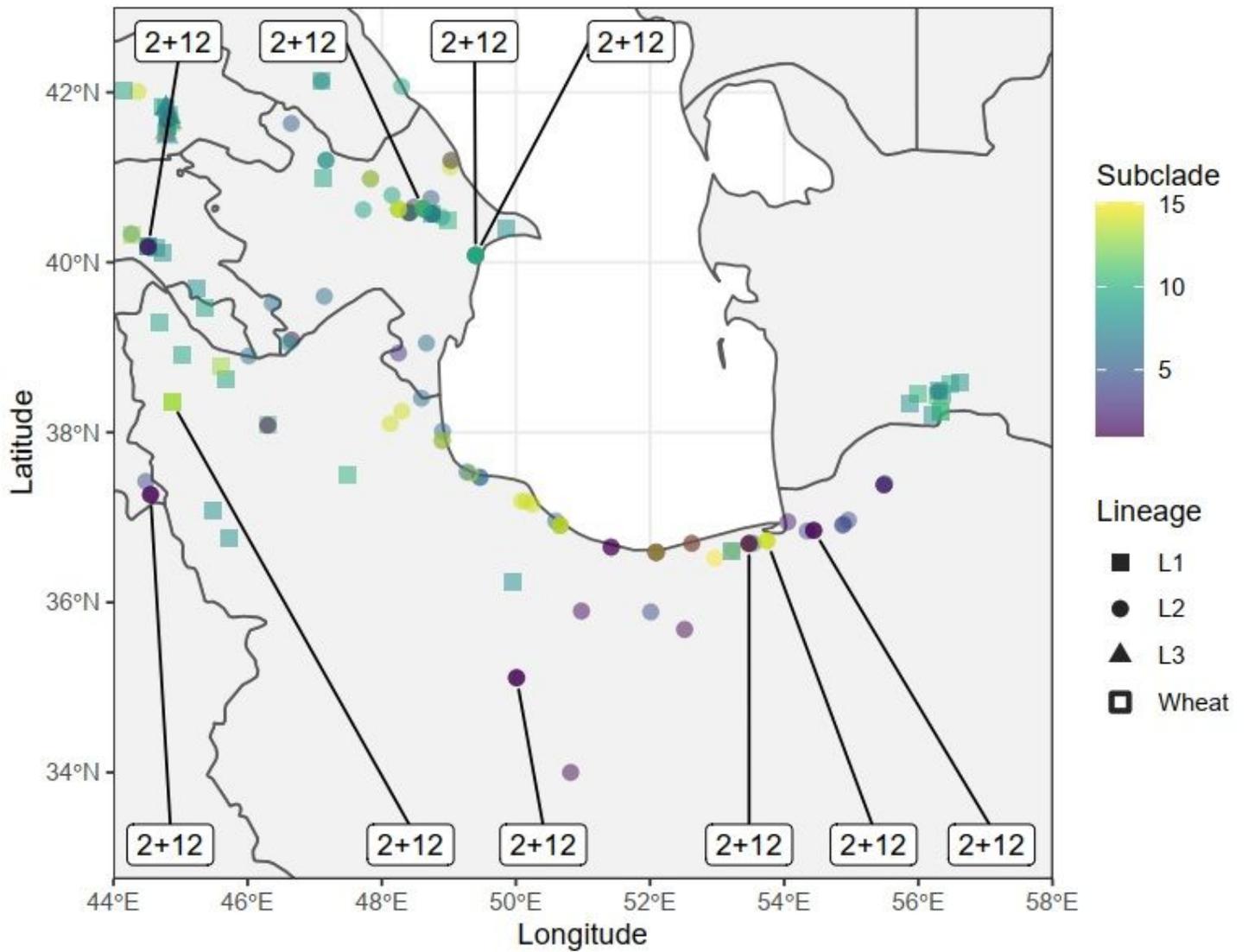


Figure 3

Cryptic haplotypes within SDS-PAGE alleles. The molecular haplotypes for *Ae. tauschii* accessions at the sites where the accessions were collected. Corresponding SDS-PAGE allele is noted for 2+12 mobility alleles.

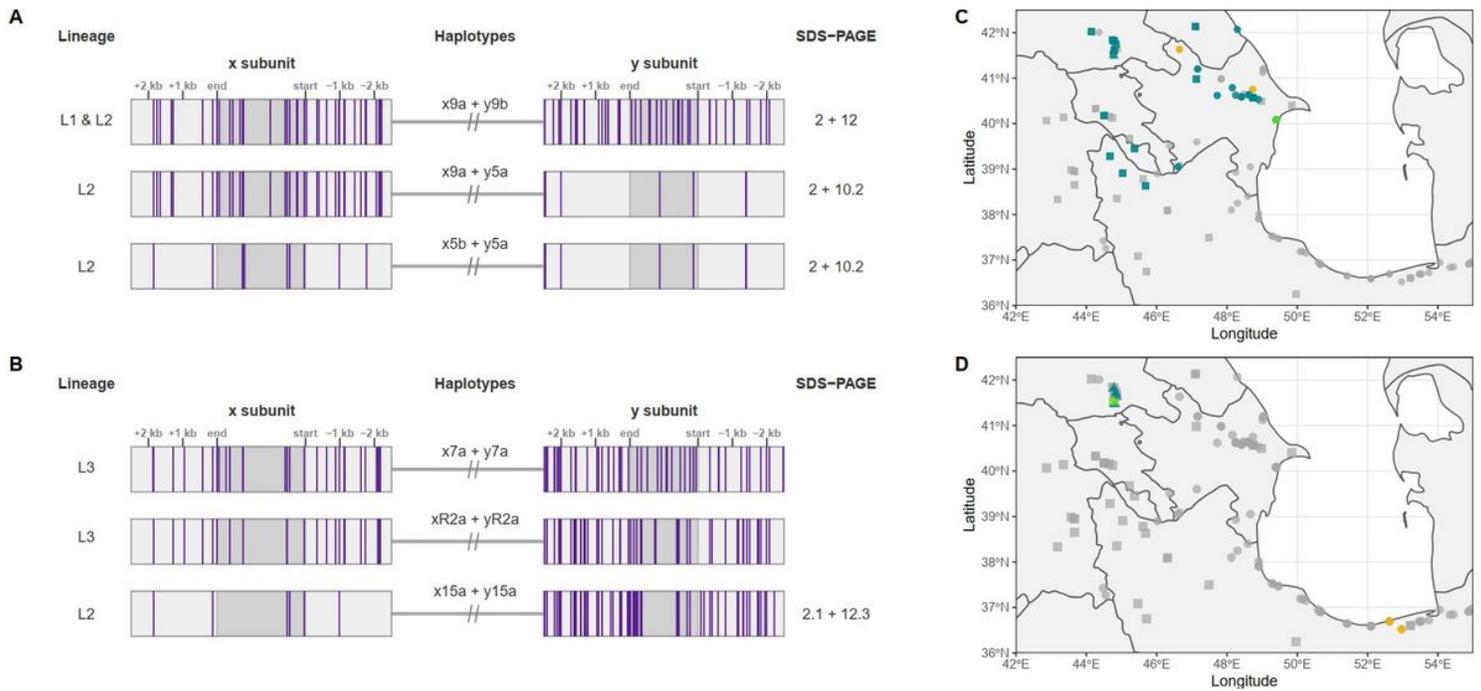


Figure 4

Glu-D1 recombinants. Molecular haplotype representation of recombinant Glu-D1 haplotypes for accessions (a) TA1668 (Lineage 2) and (b) TA2576 (Lineage 3). Vertical purple bars represent the alternate allele variants as called against the AL8 7/8 genome assembly. SDS-PAGE allele for the given haplotypes are show to right of each gene model. Haplotype Dx9a is present in both Lineage 1 and Lineage 2 accessions. The closest potential x and y subunit haplotypes involved in the recombinant haplotype xR2+yR2a of TA2576 are x7a (Lineage 3) and y15a (Lineage 2). SDS-PAGE protein mobilities for x7a + y7a and xR2+yR2a were not analyzed. Geographical distribution of recombinant Glu-D1 haplotypes for (c) TA1668 (Lineage 2) and (d) TA2576 (Lineage 3). Collection site of recombinant accessions is marked in lime green, whereas turquoise and orange designate the collection sites of accessions carrying x subunit and y subunit haplotypes, respectively. Accessions with unrelated haplotypes are in light gray. Lineages are shown in squares (Lineage 1), circles (Lineage 2) or triangles (Lineage 3).

Supplementary Files

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

- [SupplementalMaterials.pdf](#)
- [SupplementalData1popmap.csv](#)
- [SupplementalData2GluD1.txt](#)