

1 **Novel design of imputation-enabled SNP arrays for breeding and research applications supporting**
2 **multi-species hybridisation**

3

4 *Running Title: Illumina Infinium Wheat Barley 40K SNP array*

5 Keeble-Gagnère G¹, Pasam R¹, Forrest KL¹, Wong D¹, Robinson H², Godoy J², Rattey A², Moody D²,
6 Mullan D², Walmsley T², Daetwyler HD^{1,3}, Tibbits J¹, Hayden MJ^{1,3,4}

7 ¹Agriculture Victoria, AgriBio, Centre for AgriBioscience, Bundoora, Victoria 3083, Australia

8 ²InterGrain, 19 Ambitious Link, Bibra Lake, WA 6163, Australia

9 ³School of Applied Systems Biology, La Trobe University, Bundoora, Victoria 3083, Australia

10 ⁴Corresponding Author

11

12 Word Count: 8,563 (includes all sections and figure and table legends but excludes references and supplementary data)

13

14 **Abstract**

15 Array-based SNP genotyping platforms have low genotype error and missing data rates compared to
16 genotyping-by-sequencing technologies. However, design decisions used to create array-based SNP
17 genotyping assays for both research and breeding applications are critical to their success. We
18 describe a novel approach applicable to any animal or plant species for the design of cost-effective
19 imputation-enabled SNP genotyping arrays with broad utility and demonstrate its application through
20 the development of the Infinium Wheat Barley 40K SNP array. We show the approach delivers high-
21 quality and high-resolution data for wheat and barley, including when samples are jointly hybridised.
22 The new array aims to maximally capture haplotypic diversity in globally diverse wheat and barley
23 germplasm while minimising ascertainment bias. Comprising mostly biallelic markers designed to be
24 species-specific and single-copy, it permits highly accurate imputation in diverse germplasm to
25 improve statistical power for GWAS and genomic selection. The SNP content captures tetraploid
26 wheat (A- and B-genome) and *Ae. tauschii* (D-genome) diversity and delineates synthetic and
27 tetraploid wheat from other wheats, as well as tetraploid species and subgroups. The content includes
28 SNP tagging key trait loci in wheat and barley and that directly connect to other genotyping platforms
29 and legacy datasets. The utility of the array is enhanced through the web-based tool *Pretzel*
30 (<https://plantinformatics.io/>) which enables the array's content to be visualised and interrogated
31 interactively in the context of numerous genetic and genomic resources to more seamlessly connect
32 research and breeding. The array is available for use by the international wheat and barley community.

33 *(248 words)*

34

35 **Short summary**

36 Designing SNP genotyping arrays for closely related species with broad applicability in both research
37 and breeding is challenging. Here we describe a novel generic approach to select SNP content for
38 such arrays and demonstrate its utility in wheat and barley to:

- 39 • capture haplotypic diversity while minimising ascertainment bias;
- 40 • accurately impute to high SNP density in diverse germplasm;
- 41 • generate high-quality high-resolution genotypic data; and
- 42 • jointly hybridise samples to the same bead chip array.

43

44 **Keywords**

45 Triticum aestivum, wheat, Hordeum vulgare, barley, SNP genotyping array, imputation, GWAS,
46 genomic selection, dual sample hybridisation, molecular breeding

47

48 **Introduction**

49 High-density genotyping arrays that simultaneously interrogate thousands of single nucleotide
50 polymorphisms (SNP) have proven a powerful tool in genetic studies. The first generation of these
51 have been widely used in wheat and barley for various applications including genome-wide association
52 studies (GWAS), characterization of genetic resources, marker-assisted breeding and genomic
53 selection (Pasam *et al.* 2017, Joukhadar *et al.* 2017, Balfourier *et al.* 2019). Continued advances in
54 genome assembly and genotyping technologies present powerful new opportunities to continue the
55 integration of genomics information into operational plant breeding systems and extend the potential
56 of more academic research applications; e.g. studying genomic patterns of diversity, inferring
57 ancestral relationships between individuals in populations and studying marker-trait associations in
58 mapping experiments.

59 The publication of chromosome-scale genome assemblies are becoming available for more and more
60 species and this availability is expected to accelerate with international projects such as the Earth
61 BioGenome project (<https://www.earthbiogenome.org/>) which aims to sequence, catalog and
62 characterize the genomes of all of the earth's eukaryotic biodiversity over the next ten years. High
63 quality assemblies are already available in cereal crop species such as barley (Mascher *et al.* 2017,
64 Monat *et al.* 2019), emmer wheat (Avni *et al.* 2017), durum wheat (Maccaferri *et al.* 2019) and bread
65 wheat (IWGSC 2018), as well as for the diploid ancestors of wheat (Luo *et al.* 2017, Ling *et al.* 2018).
66 These assemblies have accelerated SNP discovery and our understanding of the breeding history of
67 wheat and patterns of genome-wide linkage disequilibrium (LD) in different germplasm pools. For
68 example, He *et al.* (2019) used an exome capture array in 890 globally diverse hexaploid and tetraploid
69 wheat accessions to discover 7.3M varietal SNP and investigate the role of wild relative introgressions
70 in shaping wheat improvement and environmental adaptation. Pont *et al.* (2019) exome sequenced a
71 worldwide panel of 487 accessions selected from across the geographical range of the wheat species
72 complex to explore how 10,000 years of hybridisation, selection, adaptation and plant breeding has
73 shaped the genetic makeup of modern bread wheats. Similarly, Mascher *et al.* (2019) discovered
74 almost 15M varietal SNP from exome sequence generated for 96 two-row spring and winter barley
75 accessions, a subset of which was used to investigate the extent and partitioning of molecular
76 variation within and between the two groups.

77 While SNP discovery using whole genome sequence data is currently limited to a relatively small
78 number of wheat and barley accessions, this situation is expected to rapidly change as sequencing
79 costs continue to decrease. For example, Lai *et al.* (2015) and Montenegro *et al.* (2017) used whole
80 genome sequence data from 16 and 18 bread wheat accessions to identify more than 4M and 36M
81 SNP on group 7 chromosomes and at the whole genome level, respectively. The more recent

82 publication of whole genome sequence assemblies for 14 modern bread wheat varieties from global
83 breeding programs (Walkowiak *et al.* 2020) provides additional new resources for *de novo* whole
84 genome SNP discovery and to investigate structural variation within the wheat genome. In barley, Hill
85 *et al.* (2020) used a combination of data sources including low coverage whole genome sequence of
86 632 genotypes representing major global barley breeding programs to investigate genomic selection
87 signatures of breeding in modern varieties.

88 Increasing genomic resources and increased understanding of global and local population structure
89 (Joukhadar *et al.* 2017) is enabling a shift from high to lower-density genotyping assays as a basis for
90 undertaking genetic analyses for trait dissection and mapping. Where high-density data is still
91 required, imputation can be effective to accurately infer higher marker density. Imputation uses
92 statistical approaches to fill missing genotype data and increase low-density genotype data to
93 genome-wide high-density data (Money *et al.* 2015). Imputation has been shown to increase power
94 for the detection of marker-trait associations in GWAS (Jordan *et al.* 2015, Fikere *et al.* 2020) and
95 genomic selection (Nyine *et al.* 2019). Currently, hybridisation-based SNP arrays are better suited for
96 imputation, compared to genotyping-by-sequencing (GBS) approaches, due to their lower missing
97 data rates and higher genotype calling accuracies (Rasheed *et al.* 2017, Elbasyoni *et al.* 2018).

98 To date, several hybridisation-based SNP genotyping arrays providing genome-wide coverage have
99 been developed for wheat and barley. Cavanagh *et al.* (2013) developed an Illumina iSelect array that
100 genotyped 9,000 SNP. The same technology was used a year later to design an array that assayed
101 90,000 SNP (Wang *et al.* 2014), which was subsequently used to derive a breeder-oriented Infinium
102 15K array (Soleimani *et al.* 2020). Winfield *et al.* (2016) reported an Affymetrix Axiom 820K SNP array,
103 which was also subsequently used to derive an Axiom 35K Wheat Breeders' array that targeted
104 applications in elite wheat germplasm (Allen *et al.* 2015). These genotyping arrays were largely based
105 on genome sequence fragments from early Roche 454 and Illumina assemblies, or from exome capture
106 sequence, and were generally enriched for gene-associated SNP. More recently, Rimbart *et al.* (2018)
107 reported an Axiom 280K SNP array based on content derived from the intergenic fraction of the wheat
108 genome, which to date has been poorly exploited for SNP, while Sun *et al.* (2020) described an Axiom
109 660K array based on genome-specific markers from hexaploid and tetraploid wheat, emmer wheat
110 and *Ae. tauschii*. In barley, two Infinium iSelect genotyping arrays comprising 9K and 50K SNP have
111 been reported (Comadran *et al.* 2012, Bayer *et al.* 2019).

112 While SNP genotyping arrays provide robust allele calling with high call rates and fast sample turn
113 around (typically about 3 days), they have high set up costs. The latter has presented significant
114 challenges for the development of SNP arrays that can comprehensively serve both research and
115 breeding applications; researchers have traditionally preferred high SNP density (which creates a high
116 genotyping cost per sample but low cost per datapoint), while breeders typically only want a minimally
117 sufficient marker density. This challenge drove us to develop a general approach to SNP array design
118 that specifically takes into consideration the need for low-cost genotyping across a wide range of
119 research and breeding applications, with the aim to seamlessly connect research to breeding.

120 Here, we present the design methodology and an example of its implementation in the Infinium
121 Wheat Barley 40K SNP array Version 1.0, a new and highly optimised genotyping platform containing
122 25,363 wheat-specific and 14,261 barley-specific SNP, the vast majority of which behave as easily
123 scored, single-copy biallelic markers. The SNP content was carefully selected to enable accurate
124 imputation to high SNP density in globally diverse wheat and barley germplasm, as well as within the
125 more restricted germplasm pools of breeding programs. The array is well connected to markers on
126 other commonly used SNP arrays, as well as to many existing genomic resources, and provides high
127 utility in research and breeding from germplasm resource characterisation, GWAS and genetic

128 mapping to tracking introgressions from different sources, marker-assisted breeding and genomic
129 selection. In addition, the SNP have been selected to enable joint hybridisation of wheat and barley
130 samples in the same assay, potentially halving costs for large scale deployment. The array is available
131 for use by the international wheat and barley community and is supported by the web-tool *Pretzel*
132 (Keeble-Gagnère *et al.* 2019, <https://plantinformatics.io/>).

133

134 **Materials and Methods**

135 *Germplasm and genomic resources*

136 SNP genotypes for 1,041 exome sequenced bread wheat accessions were used to select content for
137 the Infinium Wheat Barley 40K SNP array. The accessions included 790 previously reported in He *et al.*
138 (2019) to capture global wheat diversity, an additional 149 accessions selected from the global
139 collection contained in the associated VCF file (<http://wheatgenomics.plantpath.ksu.edu/1000EC/>) to
140 expand the diversity captured and 102 historical breeding lines from the InterGrain commercial wheat
141 breeding program (www.intergrain.com). The first two sets of accessions maximally captured genetic
142 diversity among 6,087 globally diverse wheat accessions comprising landraces, varieties, synthetic
143 derivatives and novel trait donor lines (He *et al.* 2019). The additional 149 accessions were selected to
144 capture genetic diversity within synthetic derivative germplasm derived from crossing 100 primary
145 synthetics (derived from interspecific hybridisation of durum wheat with *Ae. tauschii*) to three
146 Australian varieties: Yitpi, Annuello and Correll (Ogbonnaya *et al.* 2007). The latter two sets of
147 accessions were exome capture sequenced as described in He *et al.* (2019). SNP discovery was
148 performed using the first two sets of accessions and the resulting SNP list was used to call SNP
149 genotypes across all accessions.

150 The Infinium 90K wheat SNP genotypes reported in Maccaferri *et al.* (2019) for a globally diverse
151 tetraploid wheat collection of 1,856 accessions comprising wild emmer (*Triticum turgidum* ssp.
152 *dicoccoides*), domesticated emmer (*T. turgidum* ssp. *dicoccum*) and *T. turgidum* genotypes including
153 durum landraces and cultivars were used to select tetraploid wheat specific SNP.

154 A georeferenced landrace collection of 267 exome sequenced barley accessions, including 2- and 6-
155 rowed *Hordeum vulgare* landraces as well as *Hordeum spontaneum* (Russell *et al.* 2016), and 117
156 whole genome sequenced accessions representing historical breeding lines from the InterGrain
157 commercial barley breeding program were used to select content for the SNP array.

158 *SNP discovery*

159 In wheat, SNP discovery and genotype calling were performed as described in He *et al.* (2019), against
160 IWGSC RefSeq v1.0 (IWGSC 2018). After filtering for >40% call rate and >1% minor allele frequency
161 (MAF), 2.04M SNP were used for LD analysis. To filter for nucleotide variation originating from *Ae.*
162 *tauschii*, D-genome-specific SNP that had a MAF >0.1 in the synthetic derivative wheat and MAF <0.1
163 in the globally diverse wheat collection were identified. In addition, the top 2% of D-genome SNP that
164 showed differential allele frequencies between these two groups based on *Fst* values (Weir and
165 Cockerham 1984) were selected. From these two SNP sets, SNP uniformly distributed across the D-
166 genome were selected for inclusion as SNP content.

167 In barley, SNP discovery was performed as described in He *et al.* (2019) using the exome sequence
168 data published in Russell *et al.* (2016), against Morex v1.0 (Mascher *et al.* 2017). Following removal of
169 *Hordeum spontaneum*-like accessions based on PCA clustering (which left 157 *Hordeum vulgare*-like
170 accessions), the resulting SNP list was used to call SNP genotypes in the 120 InterGrain historical

171 breeding lines. After filtering for >40% call rate and >5% MAF (a higher cut-off was used in barley due
172 to the smaller reference population), 932,098 SNP were used for LD analysis.

173 *Linkage disequilibrium analysis*

174 LD analysis for the filtered SNP was performed using PLINK (Purcell *et al.* 2007) at the chromosome
175 level within each species with a maximum window size of 2 Mb; i.e. all the SNP in a tag SNP set had to
176 be within a 2 Mb window. The squared correlation coefficient (r^2) based on the allele frequency in the
177 global barley or wheat diversity panel (excluding the synthetic derivatives) between two SNP was
178 considered as a measure of LD.

179 *Choice of SNP probe designs*

180 To maximise the number of SNP assayed for a given number of probes on the bead chip array, A/T and
181 C/G variants (Infinium Type I SNP which require two probes) were avoided. To maximise SNP
182 scorability and genotype calling accuracy, polymorphism underlying the 50-mer oligonucleotide SNP
183 probe sequences was also avoided as they are known to cause shifts in SNP cluster position (Wang *et*
184 *al.* 2014). For tSNPs, the probe sequences were required to align uniquely to the target genome and
185 not align to the other genome; i.e. a wheat SNP probe had to align uniquely to the wheat genome and
186 not to the barley genome, and vice versa. Finally, an Illumina Design Tool score of ≥ 0.6 was required
187 for a probe to be included as array content. A relaxed set of criteria was also used (to tag SNP sets
188 otherwise missed) which allowed up to 3 alignments to the target genome.

189 *Selection of tagging SNP (tSNP) for imputation*

190 A custom algorithm was used to select tSNP tagging LD blocks in each of the global collections and to
191 facilitate imputation from the density of the SNP array. In brief, for each chromosome the algorithm
192 iteratively selected the most informative tSNP passing all filters (based its r^2 value from the LD
193 analysis), removed all SNPs linked to the selected tSNP from the remaining list of SNPs, as well as all
194 SNP linked to any SNP in the selected tSNP set to avoid directly tagging any SNP at $r^2 \geq 0.9$ more than
195 once, before repeating the process until a target number of tSNP was reached. This process ensured
196 the set of tSNP selected was the minimum set required to tag the most SNPs at $r^2 \geq 0.90$. Specifically,
197 for a given a set of SNPs $S = \{s_1, s_2, \dots\}$ and function $r^2(s_i, s_j)$ defining the *Pearson correlation*
198 *squared* $\forall s_i, s_j \in S$, we defined the tSNP set for s_i at q to be:

$$199 \quad T_{s_i}^q = \{s_j \in S \mid r^2(s_i, s_j) \geq q\}.$$

200 Rename the $T_{s_i}^q$ and define $T_{sorted}^q = (T_{s_j}^q)_{j=1}^n = T_{s_1}^q, T_{s_2}^q, T_{s_3}^q, \dots$ where $i \geq j \Rightarrow |T_{s_i}^q| \geq |T_{s_j}^q|$.

201 In other words, T_{sorted}^q is an ordering of tSNP sets, monotonically decreasing in size.

202 Let $F \subset S$ be a subset of filtered SNPs. Define $F(T_{sorted}^q) = \{T_{s_j}^q \mid s_j \in F\}$.

203 We define $T_{sorted}^q - T_{s_i}^q = \{(T_{s_j}^q)_{s_j \in S} \mid T_{s_i}^q \cap T_{s_j}^q = \emptyset\}$, and $head(L)$ to be the first element of the
204 ordered sequence L .

205 The algorithm is then:

206 $S \leftarrow \emptyset$

207 $T \leftarrow head(F(T_{sorted}^q))$

208 while $|T| \geq m$:

$$\begin{aligned} 209 \quad & S \leftarrow S \cup T \\ 210 \quad & T_{sorted}^q \leftarrow T_{sorted}^q - T \\ 211 \quad & T \leftarrow head(F(T_{sorted}^q)) \end{aligned}$$

212 The above is applied with $q = 0.9, m = 10$ to define the imputation set S .

213 To guard against possible loss of imputation accuracy due to SNP assays failing to provide reliable
214 genotypes calls, a level of redundancy was included in the tSNP sets for wheat and barley. Specifically,
215 three tSNP were chosen when the number of SNP tagged was ≥ 50 and two tSNP were selected when
216 the number of SNP tagged was ≥ 20 . Single tSNP were included as array content when they tagged at
217 least 10 SNP. Some tag SNP sets could not be tagged because no probe passed all filters; in this case
218 we ran the algorithm on the remaining sets allowing SNP passing relaxed filters (up to 3 hits to target
219 genome were allowed). In addition, tSNP were selected to tag genomic regions that had sparse SNP
220 coverage but high LD; i.e. tagging < 10 SNP within windows larger than 500Kb in wheat and 1Mb in
221 barley. Finally, SNP were selected in regions still lacking SNP after the previous steps.

222 *Optimisation of SNP content*

223 To ensure broad applicability of the SNP array in research and breeding, the content included SNP
224 selected to specifically interlink germplasm resources such as the 19,778 domesticated barley
225 accessions with GBS genotypes described in Milner *et al.* (2019). It also included SNP probes designed
226 to interrogate published trait-linked markers in wheat and barley. Designs for these markers were
227 based directly on published sequence or from alignment of published primers or flanking sequences
228 and inference of the targeted nucleotide variation. For all trait-linked markers, the best probe design
229 was selected based solely on the Illumina quality score. Due to difficulty for designing SNP probes
230 targeting known alleles of phenology genes, we selected 293 exome SNPs around the genes reported
231 in Shi *et al.* (2019).

232 *Imputation*

233 The wheat and barley global diversity sets were used as reference haplotypes for imputation. For
234 wheat, accessions clustering with the synthetic derivatives in a PCA analysis were excluded. For barley,
235 only samples with $< 20\%$ missing data were used. In both species, missing data was filled in using
236 Beagle (Browning *et al.* 2007) and phased with Eagle (Loh *et al.* 2016). In total, 868 and 155 wheat and
237 barley lines were used as reference haplotypes.

238 In wheat, SNP coordinates were converted to IWGSC v2.0 pseudomolecules
239 (https://urgi.versailles.inra.fr/download/iwgsc/IWGSC_RefSeq_Assemblies/v2.0/, Zhu *et al.* 2021) before
240 imputation. After transfer into the v2.0 assembly, there were 18,521 SNP before imputation, with
241 630,058, 549,003 and 352,947 tagged at $r^2 \geq 0.50, 0.70$ and 0.90 , respectively.

242 To assess the accuracy for imputation into globally diverse germplasm, 100-fold cross validation was
243 performed. A random subset of 100 wheat (or 10 barley) lines had their true genotypes masked,
244 leaving only the tSNP. The remaining lines were then used as the reference population with Minimac3
245 software (Das *et al.* 2016) to impute back the missing genotypes for three different target SNP sets:
246 the set of SNP tagged at $r^2 \geq 0.50, 0.70$ and 0.90 , respectively. The imputation accuracy for each line,
247 measured as both correlation and concordance between the actual and imputed genotypes, was
248 calculated from 100 repetitions of this process in each of wheat and barley. Correlation was measured
249 as the Pearson r^2 between SNP called in both genotypes being compared, while concordance was

250 measured as the fraction of SNP in agreement between those called in both genotypes being
251 compared.

252 *SNP assay and genotype calling*

253 Samples were assayed following the protocol for Infinium XT bead chip technology (Illumina Ltd). SNP
254 clustering and allele calling was performed using GenomeStudio Polyploid software (Illumina Ltd)
255 using the Illumina-supplied wheat or barley SNP manifest file. The custom genotype calling pipeline
256 described in Maccaferri *et al.* (2019) was also used.

257

258 **Results**

259 *Overview of design approach*

260 The central idea of the design concept is to exploit LD using the r^2 measure to define sets of SNP that
261 can be considered equivalent: for a given SNP, we define its tag SNP set as the set of SNP with $r^2 \geq 0.9$
262 (the set of SNP in this set are referred to as tSNP). This metric provides a measure of equivalence as
263 well as a natural ranking of SNP by their informativeness, as defined by the size of the tSNP set to
264 which they belong. We assume the relationship is symmetrical; i.e. if SNP A is in SNP B's tSNP set, then
265 SNP B should be in SNP A's tSNP set. The original set of SNP is then filtered using technology and
266 application-specific criteria (see Materials and Methods) while maintaining connectivity to SNP that
267 fail the filters via the tSNP sets of SNP that pass the filters.

268 To design a genotyping array that has broad applicability in research and breeding, the SNP should be
269 discovered in diverse germplasm to avoid ascertainment bias (since LD is population dependent) and
270 with sufficient density to produce large tSNP sets. The latter helps ensure at least one SNP in a tSNP
271 set will pass all the design filters in most instances. Here, we used a globally diverse set of barley
272 landrace accessions and a globally diverse set of wheat accessions that included landraces, varieties,
273 novel trait donor and historical breeding lines (Figure 1). For array designs focused only on breeding
274 applications, SNP discovery should aim to capture the genetic diversity within the breeding germplasm
275 pool.

276 A novel selection algorithm (described in Materials and Methods) is then used to select SNP which
277 maximise LD capture, while minimising the number of SNP assayed on the array, using only SNP that
278 pass the design filters.

279 The design concept can be applied to any animal or plant species. In addition to this set of SNP, utility
280 in research and breeding can be further enhanced by including context-relevant SNP such as trait-
281 linked markers and markers that link germplasm resources across different genotyping technologies.

282 *SNP discovery and filtering*

283 Filtering for a minimum minor allele frequency (MAF) of 1% and maximum missing rate of 60% using
284 the 8,869,370 wheat SNP published in He *et al.* (2019) resulted in 2,037,434 high quality SNP for
285 downstream analysis. Of these, 122,799 SNP had at least one array probe that passed all design filters.
286 In barley, filtering of the 1,843,823 SNP identified from our processing of exome capture sequence
287 from the accessions from Russel *et al.* (2016) for MAF >5% and missing rate <60% resulted in 932,098
288 high quality SNP for downstream analysis, of which 119,633 SNP had at least one array probe passing
289 all filters. The filtered SNP matrices used in subsequent analysis are available at
290 https://dataverse.harvard.edu/dataverse/WheatBarley40k_v1.

291 *LD analysis and selection of tagging SNP for imputation*

292 Based on LD values of $r^2 \geq 0.9$, a total of 1.07M wheat and 413,508 barley high quality SNP were
293 singletons; i.e. had no SNP within 1Mb up- and downstream with $r^2 \geq 0.9$. These SNP were either
294 genuine singletons or categorised as singletons due to the absence of additional SNP within the
295 surrounding 2Mb region. As singleton SNP can only be tagged directly, which is not feasible on a low-
296 density array, these SNP were not considered further for inclusion on the array.

297 The custom selection algorithm grouped the 122,799 non-singleton wheat SNP passing all design
298 filters into 11,076 tSNP tagging SNP sets containing ≥ 10 SNP within a 2 Mb window. These tSNP tagged
299 317,599, 538,326 and 652,476 SNP at $r^2 \geq 0.9$, 0.7 and 0.5, respectively. Of the 119,633 non-singleton
300 barley SNP passing all filters, the selection algorithm identified 7,316 tSNP which tagged a total of
301 150,096, 294,659 and 390,844 SNP at $r^2 \geq 0.9$, 0.7 and 0.5, respectively. At the genome level, the rate
302 of return per tSNP was surprising similar for wheat and barley and plateaued at about 15,000 tSNP at
303 $r^2 \geq 0.9$ (Figure 2). However, the rate of return per tSNP varied at the chromosome level (Figure S1).

304 In total 21,012 wheat and 13,469 barley tSNP were included as content on the array. This tally includes
305 redundant SNP selected to guard against possible loss of imputation accuracy due to SNP assays that
306 might fail; SNP passing a relaxed set of filters (allowing up to 3 alignments to the target genome) and
307 tagging SNP sets untaggable with the stricter filtered SNP; and SNP to tag genomic regions that had
308 sparse SNP coverage but high LD; i.e. tagging < 10 SNPs within windows larger than 500Kb in wheat
309 and 1Mb in barley. The latter SNP are expected to support increased imputation density in these
310 regions as higher density SNP datasets become available into the future. The wheat tSNP tagged a
311 total of 394,034, 636,641 and 758,452 SNP at $r^2 \geq 0.9$, 0.70 and 0.50 respectively, while the barley tSNP
312 tagged a total of 187,412, 361,012 and 471,645 SNP, respectively. Importantly the MAF distributions
313 for the tSNP, tagged SNP and filtered SNP from the globally diverse wheat and barley collections
314 closely matched one another, respectively (Figure 3).

315 *Accuracy for imputing into globally diverse germplasm*

316 The ability to impute from the tSNP on the array to the sets of SNP tagged at $r^2 \geq 0.50$, 0.70 and 0.90
317 respectively in globally diverse wheat and barley germplasm was assessed using 100-fold cross
318 validation. Accuracy was determined from the correlation and concordance between the imputed and
319 actual genotypes for each wheat or barley line averaged over the occurrences of that sample within
320 the 100 iterations.

321 As expected, all metrics were highest when imputing to the set of SNP tagged at $r^2 \geq 0.90$ and lowest
322 for those tagged at $r^2 \geq 0.50$ (Table 1). In wheat, only a small decrease in accuracy was observed for
323 most accessions as the size of the tagged SNP set increased (i.e. r^2 decreased), with reduced accuracy
324 most evident in the bottom 50 accessions (Figure 4). For these accessions, the difference in accuracy
325 (both correlation and concordance) between comparisons including and excluding heterozygous
326 genotype calls was almost 10%, suggesting the possibility of high error rates in the heterozygous
327 exome SNP calls for these accessions. 768 (88.5%) of the wheat accessions had accuracies $\geq 90\%$ with
328 the strictest correlation metric (which included heterozygous calls) for the set of SNP tagged at r^2
329 ≥ 0.50 . When comparing only non-heterozygous calls, the number of lines above this threshold rose to
330 866 (99.8%) (Figure 4).

331 Reduced accuracy when imputing to higher tagged SNP numbers was more pronounced in barley. A
332 difference of 10.8% (from 96.8% to 86%) was observed between the average correlation (which
333 included heterozygous calls) for the set of SNP tagged at $r^2 \geq 0.90$, compared to those tagged at r^2
334 ≥ 0.50 (Table 1). As observed in wheat, the inclusion of heterozygous calls reduced the accuracy,

335 particularly when imputing to the set of SNP tagged at $r^2 \geq 0.50$, again suggesting possible erroneous
336 heterozygous calls in the sequence genotypes (Figure 4). The reduced accuracies observed in barley
337 compared to wheat are also likely partly due to the reduced size of reference haplotypes (155 versus
338 868). Accuracies in barley would likely improve if the reference haplotype set was expanded.

339 *Wheat-barley 40K SNP array content*

340 The final array design comprised 34,481 imputation SNP and two additional categories of context-
341 specific SNP (content summarised in Table 2, full details are in Table S1).

342 The first context-specific category included 2,609 SNP from the Infinium wheat 90K SNP array (Wang
343 *et al.* 2014) that were selected based on allele differentiation to tag tetraploid wheat (A- and B-
344 genome) diversity and to clearly delineate tetraploid wheat from other types of wheat, as well as
345 distinguish tetraploid species and subgroups from one another. The SNP comprised four classes: 1)
346 differentiating SNP that represent the top 2% F_{st} values in Maccaferri *et al.* (2019) between the four
347 subgroups of tetraploid species: wild emmer, domesticated emmer, domesticated wild emmer, durum
348 landraces and durum cultivars; 2) subgroup-specific private SNP that showed a $MAF \geq 0.1$ in one of the
349 subgroups and were either monomorphic or showed a $MAF < 0.05$ in the other subgroups; 3)
350 subgroup-specific high MAF SNP that were present at ≥ 0.3 MAF in any one of the subgroups; and 4)
351 neutral SNP that did not show any signatures of selection, were polymorphic in all subgroups and
352 showed an overall MAF of ≥ 0.4 . The ability of these SNP to reliably differentiate the tetraploid species
353 subgroups as efficiently as the Infinium wheat 90K array is shown in Figure S2.

354 The second category included 1,206 exome SNP tagging *Ae. tauschii* (D-genome) diversity present in
355 backcross synthetic derivatives that originated from crosses involving 100 primary synthetic parents,
356 which were selected for phenotypic and genetic diversity among about 400 primary synthetics
357 developed at CYMMIT and imported into Australia in 2001. Each of the 100 primary synthetic parents
358 was derived from a different *Ae. tauschii* accession. The SNP were selected to provide high D-genome
359 coverage, enriched density in highly recombining chromosomal regions and to clearly delineate bread
360 wheat from other types of wheat, as well as tag diversity in synthetic wheat and their derivatives and
361 *Ae. tauschii*. The SNP comprised two classes: 1) differentiating SNP that represent the top 2% F_{st}
362 values between the global diversity wheat and synthetic derivative collections; and 2) D-genome
363 diversity from *Ae. tauschii* that showed a $MAF \geq 0.1$ in the synthetic derivative collection and $MAF \leq 0.1$
364 in the global diversity wheat collection. The ability of these SNP to reliably differentiate synthetic
365 wheat from common wheat as efficiently as the Infinium wheat 90K array is shown in Figure S3.

366 The final category included linked SNP for key breeding traits and SNP linking major germplasm
367 resources genotyped with different technologies. In total, 457 wheat and 178 barley SNP
368 corresponded to published trait-linked markers with 109 SNP associated with agronomically
369 important genes (Table S1). Another 614 SNP provide a direct link to 19,778 GBS genotyped
370 domesticated barley accessions (Milner *et al.* 2019).

371 *Assay performance – Single sample hybridisations*

372 A limitation of hybridisation-based genotyping arrays is that their oligonucleotide probes hybridise
373 both to the targeted locus and its homoeologues and paralogues if present (Cavanagh *et al.* 2013;
374 Wang *et al.* 2014). Consequently, the ratio of allele-specific fluorescent signals observed for an assay
375 depends on the locus copy number in the genome, with increasing copy number reducing the allele-
376 specific fluorescent signal ratio and separation of SNP allele clusters. Further, SNP assay scorability
377 and genotype calling can be confounded by the presence of mutations that modify oligonucleotide
378 annealing such that different cluster patterns are observed across germplasm (Wang *et al.* 2014). An

379 ideal assay design for a hybridisation-based genotyping array is therefore an oligonucleotide probe
380 that binds at only one locus in the genome and has no known nucleotide variation underlying the
381 probe hybridisation site. Theoretically this should ensure three distinct clusters corresponding to the
382 genotypic states (REF, HET and ALT) expected of a single copy biallelic SNP. The increasing availability
383 of genomic resources is now allowing this historical problem to be addressed. Hence, we used the
384 combination of reference genome assemblies and genotypic data for large globally diverse wheat and
385 barley collections to specifically target the design of single copy biallelic SNP assays.

386 For the purpose of evaluating the performance of the array, the wheat and barley diversity
387 populations were used to define cluster positions for SNP genotype calling. The vast majority (98%) of
388 the 39,654 SNP assays on the array produced scorable cluster patterns when hybridised with a barley
389 or wheat sample; 91% (12,949/14,261) of the barley and 83% (20,090/24,598) of the wheat SNP assays
390 could be reliably scored as single-copy biallelic markers, with the REF and ALT clusters having Theta
391 values close to 0 and 1 in GenomeStudio SNP plots (Figure 5). While the remaining SNP could typically
392 be reliably scored as biallelic markers, they showed cluster compression indicative of multiple loci.
393 Few assays showed complex clustering patterns indicating the success of designing probes without
394 underlying polymorphism. Five and 7% of wheat and barley assays showed a clustering pattern typical
395 for the presence of a null allele. The occurrence of assays not behaving as single-copy biallelic markers
396 reflects current knowledge gaps for structural variation in the genomes of wheat and barley including
397 both copy number variation and presence-absence variation (Wang *et al.* 2014, Balfouier *et al.* 2019,
398 Walkowiak *et al.* 2020).

399 The concordance between called and actual genotypes was exceptionally high for both wheat and
400 barley. The genotype concordance and correlation were 99.5 and 98.1% respectively in wheat when
401 heterozygous genotype calls were excluded, and 97.6 and 95.7% when heterozygous calls were
402 included. Similarly, 99.8% concordance and 99.2% correlation were observed in barley when
403 heterozygous calls were excluded, and 98.2 and 97.2% was observed with heterozygous calls included.
404 The average missing data rates was 4.8 and 3.8% in wheat and barley, respectively.

405 *Assay performance – Dual sample hybridisations*

406 The design process specifically aimed to select species-specific SNP probes and thus it should be
407 theoretically possible to jointly hybridise a wheat and barley sample to the same bead chip array (dual
408 hybridisation) without loss of genotype calling accuracy. Cross-hybridisation between species is
409 expected to confound genotype calling accuracy by creating shifts in SNP cluster positions and/or
410 complex clustering patterns that cannot be easily scored.

411 To evaluate assay performance of a dual hybridisation, samples from the InterGrain commercial barley
412 and wheat breeding programs were used to define cluster positions and call SNP genotypes for 576
413 dual hybridisation assays. The same samples were also assayed in single sample hybridisation assays
414 to enable genotype calling accuracy between dual and single hybridisation assays to be directly
415 compared.

416 Most of the barley and wheat SNP in dual hybridisation assays produced scorable cluster patterns.
417 Shifts in cluster positions were observed, which indicated either that some oligonucleotide probes
418 showed a degree of cross-species hybridisation or that deviation from the standard amount of sample
419 DNA (200 ng per sample) recommended for the bead chip assay affected signal-to-noise. Through
420 empirical testing, we found the quantity of genomic DNA per sample was a major factor causing shifts
421 in cluster position (data not shown) and could be minimised by adjusting the input DNA for each

422 sample to match the ratio of the genome size for each species; e.g. 200 ng barley DNA and 600 ng
423 wheat DNA; the bread wheat genome is about three times larger than that of barley.

424 For the purpose of assessing genotype calling accuracy for dual hybridisation assays, only SNP that
425 revealed polymorphism among the 576 wheat and barley samples assayed were considered. Of the
426 9,826 barley and 9,118 wheat SNP showing polymorphism, the vast majority were easily scored as
427 biallelic markers and had good cluster separation, indicating that oligonucleotide probe cross-species
428 hybridisation was minimal (Figure 6). The average concordance between genotype calls for the same
429 wheat and barley samples in single and dual sample hybridisation assays was 99.9, 96.7, and 99.8%
430 for the REF, HET and ALT alleles, respectively. The average missing data rate across the wheat and
431 barley samples was similar for both assay types, with 4.7 and 2.0% in dual and single hybridisation
432 assays, respectively.

433

434 **Discussion**

435 High-throughput, low-cost and flexible genotyping platforms are required for both research and
436 breeding applications. Compared to GBS and PCR-based marker systems, array-based genotyping
437 platforms are highly commercialised and highly customisable, both for the number of markers and
438 samples assayed. They also have low genotype error and missing data rates compared to GBS
439 technologies (Rasheed *et al.* 2017). Consequently, SNP arrays are widely utilised and several low-
440 density SNP genotyping arrays have been developed for wheat and barley. Here, we described a novel
441 approach that is applicable to any animal or plant species for the design of cost-effective, imputation-
442 based SNP genotyping arrays with broad utility and that support the hybridisation of multiple samples
443 to the same SNP array. The utility of the approach was demonstrated through the development of the
444 Infinium Wheat Barley 40K SNP array.

445 The key difference between Infinium Wheat Barley 40K SNP array and previously reported array-based
446 genotyping assays is a paradigm shift in the logic underpinning its design. To date, commonly used
447 low-density genotyping arrays are comprised of the most scorable and informative markers from
448 higher density arrays. For example, the Infinium Wheat 15K SNP array (Soleimani *et al.* 2020) and
449 Axiom Wheat Breeders' 35K SNP array (Allen *et al.* 2015) are derived from the Infinium Wheat 90K
450 SNP array (Wang *et al.* 2014) and Axiom Wheat 820K SNP array (Winfield *et al.*, 2016). SNP on the
451 Infinium 90K SNP array were derived from transcriptome sequence of 26 bread wheat accessions,
452 while those on the Axiom 820K array were based on exome capture sequence from 43 bread wheat
453 and wild species accessions representing the primary, secondary and tertiary gene pools. While these
454 derived low-density arrays are affordable for routine deployment in breeding and research, their
455 content is breeder-oriented and has limited utility outside the primary gene pool of hexaploid wheat.

456 The design implemented in the Infinium Wheat Barley 40K SNP array is based on the hugely expanded
457 genotypic and genomic resources now available for wheat and barley. By using these resources, we
458 were able to identify species-specific single-copy tSNP that capture a large proportion of the
459 haplotypic diversity in globally diverse germplasm, are highly scorable for accurate genotype calling,
460 minimise ascertainment bias and enable accurate imputation to high SNP density. In the case of
461 wheat, this included the use of 2.04M SNP identified from exome sequence data of 1,041 accessions
462 selected to maximally capture genetic diversity among a global collection of 6,700 accessions
463 genotyped using the Infinium 90K SNP array (He *et al.* 2019; Figure 1a). The global collection included
464 landraces, released varieties, synthetic derivatives, and novel trait donor and historical breeding lines.
465 For barley, this included 932,098 SNP identified from exome sequence data of 267 accessions selected

466 to maximally capture geographic diversity among landraces (Russell *et al.* 2016; Figure 1b), as well as
467 SNP identified from target capture sequencing of 174 flowering time-related genes performed in 895
468 worldwide accessions (Hill *et al.* 2019). The latter dataset included global diverse cultivated and
469 landrace germplasm.

470 By selecting tSNP enabling accurate imputation of common haplotype block diversity in globally
471 diverse germplasm, the Infinium Wheat Barley 40K array is expected to maintain power for GWAS,
472 genetic mapping and genomic selection (Jordan *et al.* 2015, He *et al.* 2015, Negro *et al.* 2019, Nyine *et al.*
473 *et al.* 2019). Haplotype blocks are essentially fixed stretches of DNA sequence that show little historical
474 evidence of recombination and are effectively inherited as genetic units that are shuffled and
475 assembled during breeding. The univariate LD metric r^2 has been used in many tSNP algorithms as it
476 is a major determinant of imputation accuracy and has a simple inverse relationship with the sample
477 size required to detect associations in GWAS (Carlson *et al.* 2004, Ding and Kullo 2007). By selecting
478 tSNP with an $r^2 \geq 0.9$ cut-off, we aimed to retain most of the information content in the original SNP
479 set and to balance the power loss with the effort needed to compensate with increased sample
480 numbers in downstream GWAS (~11%; i.e. $1/0.9$). A significant advantage for using r^2 is that it allows
481 a high degree of flexibility in the composition of the final tSNP set, thereby enabling other design
482 criteria to be applied without compromising overall tagging efficiency. This was especially important
483 for implementing array design principles such as for selecting species-specific single-copy SNP targets
484 that had no nucleotide variation underlying the probes to both maximise SNP scorability and support
485 dual sample hybridisation assays. The success of our approach was confirmed by >97% accuracy (as
486 measured by both correlation and concordance between the imputed and actual SNP genotypes) for
487 imputing the set of SNP tagged at $r^2 \geq 0.9$ (inclusive of heterozygous calls) in both wheat and barley.
488 Importantly, imputation accuracy was also high for the set of SNP tagged at $r^2 \geq 0.5$ (Table 1). To
489 futureproof the array design, we added tSNP tagging genomic regions in wheat and barley that had
490 sparse exome SNP coverage but high LD. We expect this content will similarly support accurate
491 imputation to whole genome sequence once genomic resources needed to achieve this are available.

492 In emphasising the design focus on selecting tSNP for imputation, we also point out the limitations it
493 has for fully capturing haplotype diversity in global wheat and barley germplasm. First, we did not tag
494 LD blocks comprised of fewer than 10 SNP since this would have required an order of magnitude more
495 SNP assays on the array; about 30,000 tSNP per species was required to tag about half of the non-
496 singleton exome SNP at $r^2 \geq 0.9$ in each of wheat and barley (Figure 2). This presents a limitation for
497 trait mapping using GWAS (but not genetic mapping) since trait loci located in untagged LD blocks will
498 become increasingly harder to detect as their LD with a SNP on the array decreases. This limitation
499 can be partly overcome by increasing sample size but is an unavoidable consequence of low-density
500 arrays, despite our tSNP selection algorithm ensuring that we maximised the number of SNP tagged
501 in LD. And second, in wheat the set of SNPs and LD relationships between them is still limited by the
502 data currently available. As exome capture sequencing assays only 2-3% of the genome, the SNP
503 discovered represent just a fraction of the true SNP density. It is therefore possible that SNP were not
504 selected simply because the haplotype they represent was only sampled by a small number of SNP in
505 that region and was below our selection thresholds. This limitation will only be overcome by large-
506 scale whole genome sequencing efforts which are just beginning to become affordable for large
507 genome-sized species. It should be noted that the LD patterns detected in this study will remain valid
508 even with higher density sequencing and that the majority of the tagged LD haplotypes span across
509 capture regions and so the number of SNP in high LD with the selected tSNP will only increase as higher
510 density SNP data becomes available.

511 An argued advantage for GBS assays is that they are ascertainment bias free. Ascertainment bias can
512 result in rare alleles being missed and genetic diversity being underestimated in non-ascertained
513 populations (Clark *et al.* 2005), with its impact dependent on the study being undertaken. Increasing
514 marker density and including low MAF markers in GWAS boosts power for QTL detection (Negro *et al.*
515 2019, Fikere *et al.* 2020). Chu *et al.* (2020) reported that very low frequency markers (MAF <0.05)
516 contributed to an improvement of genomic prediction accuracy in 378 winter bread wheat genotypes,
517 and combined with the expectation that valuable novel diversity is most likely rare (Mascher *et al.*
518 2019), suggests that rare markers deserve careful consideration. Our tSNP selection algorithm
519 prioritises haplotypes that diverge significantly from the reference genome used for SNP discovery in
520 order to maximise the number of SNP tagged in LD; it is agnostic to the MAF of individual SNP (beyond
521 the MAF cut-offs of 1% and 5% in wheat and barley, respectively). Consequently, the MAF spectrum
522 of the wheat and barley tSNP closely resembled that observed for both the sets of tagged SNP and the
523 filtered SNP in the globally diverse collections (Figure 3). Hence, we suggest the Infinium Wheat Barley
524 40K array has minimal ascertainment bias. Since tagging all minor variants is not feasible using low-
525 density arrays, a better solution is to add minor variants into future versions of the array as trait
526 associations are discovered, essentially as we have currently done for published trait linked markers.

527 To drive efficiencies for large-scale genotyping in commercial breeding programs, we explored the
528 limits of the Infinium bead chip technology. One advantage of this technology is that each
529 oligonucleotide assay probe has a unique physical position on the bead chip. This allows SNP arrays to
530 be designed to genotype multiple crop species, with a user-defined number of SNP assigned to each
531 species. The Infinium Wheat Barley 40K array assays 25,393 SNP in wheat and 14,261 SNP in barley.
532 To the best of our knowledge, multispecies SNP arrays have only been used to assay a single sample
533 at the time. Here, we demonstrated that through careful selection of species-specific oligonucleotide
534 probes it is possible to jointly hybridise a wheat and barley sample to the same bead chip array,
535 without substantial loss of genotype calling accuracy (Figure 6). The selection of such probes is
536 facilitated by our design concept which exploits LD to identify SNP that can be considered equivalent
537 for the purpose of genotyping. From a deployment perspective in a commercial breeding program,
538 dual hybridisation doubles genotyping throughput, since twice as many samples can be processed
539 given the same amount of time and resource. Dual hybridisation genotyping is potentially a game
540 changing option for the adoption of genomics technologies by breeding companies that have large
541 numbers of samples that can be co-ordinated into genotyping.

542 To ensure broad utility in research and breeding, we added SNP content capturing genetic diversity in
543 the secondary and tertiary gene pools of wheat. This included 2,609 SNP from the Infinium 90K SNP
544 array (Wang *et al.* 2014) tagging tetraploid wheat (A- and B-genome) diversity and clearly delineating
545 tetraploid wheat from other types of wheat, as well as tetraploid species and subgroups from one
546 another. Each SNP is single copy in tetraploid wheat and has been genetically and physically mapped
547 (Maccaferri *et al.* 2019). It also included 1,206 single-copy SNP tagging *Ae. tauschii* (D-genome)
548 diversity represented in 100 primary synthetic wheats, where each primary synthetic was derived
549 from a different *Ae. tauschii* accession. Collectively, these SNP provide broad utility ranging from the
550 differentiation and genetic characterisation of tetraploid and synthetic wheat (as well as other
551 secondary and tertiary gene pools of wheat) to the tracking of introgressed genomic segments during
552 breeding. Also included are SNP that directly link to the Infinium 90K (Wang *et al.* 2014) and 15K
553 (Soleimani *et al.* 2020) wheat arrays to ensure connectivity with legacy genotypic datasets and
554 research. For barley, we included 685 SNP that overlap with SNP reported for 19,778 GBS genotyped
555 accessions from the IPK Genebank (Milner *et al.* 2019) to provide a direct anchor to that resource, and
556 1,239 SNP that overlap with the Infinium 50K barley SNP array (Bayer *et al.* 2017) which link to 21,606

557 common SNP following imputation. Finally, we included trait-linked SNP and SNP tagging GWAS signals
558 for key breeding and research targets reported in the published literature.

559 The overall array design makes it ideal for a wide range of research and breeding applications, from
560 germplasm resource characterisation, GWAS and genetic mapping to tracking introgressions from
561 different sources, marker-assisted breeding and genomic selection. Its utility is further enhanced
562 through the web-based tool *Pretzel* (Keeble-Gagnère *et al.* 2019; <https://plantinformatics.io/>) which
563 enables the array's content to be visualised and interrogated in real-time in the context of numerous
564 genetic and genomic resources. For example, the SNP can be visualised relative to the genetic and
565 physical positions of other DNA marker types (e.g. SSRs, DArT), SNP on other genotyping arrays, trait
566 loci, annotated genes and syntenic positions in the genomes of other crops and model species. The
567 ability to upload and visualise data in *Pretzel* allows breeders and researchers to seamlessly link and
568 interrogate their own data in the context of publicly available datasets hosted in *Pretzel*. Combined,
569 the Infinium Wheat Barley 40K SNP array and *Pretzel* enable legacy and current research to seamlessly
570 connect to breeding.

571 In conclusion, we have described a novel approach applicable to any animal or plant species for
572 designing cost-effective imputation-enabled SNP genotyping arrays which have broad applicability in
573 research and industry applications (e.g. GWAS, genomic prediction and operational breeding) and
574 support the hybridisation of multiple samples to the same array. The utility of this design approach
575 was demonstrated through its implementation to develop a new Infinium Wheat Barley 40K SNP array.
576 In addition, to supporting broad utility in research and breeding, this array can be used as a resource
577 to connect genetic and genomic datasets generated across germplasm pools and time. The array is
578 further supported by the publicly available web-tool *Pretzel* and is available for purchase by the
579 international wheat and barley community from Illumina Ltd, the manufacturer of the Infinium bead
580 chip technology.

581

582 **Data Statement**

583 Exome data used from Russel *et al.* 2016 and He *et al.* 2019 are accessible under EBI ENA project
584 accession numbers PRJEB8044 and PRJEB31218, respectively. The filtered set of exome genotype calls
585 for accessions and SNP underpinning the LD analysis and tag SNP selection for wheat
586 (<https://doi.org/10.7910/DVN/5LVYI1>) and barley (<https://doi.org/10.7910/DVN/CUPAXD>) as well as
587 the D-genome synthetic derivative-enriched SNP matrix (<https://doi.org/10.7910/DVN/0QEASE>) are
588 available through Dataverse at https://dataverse.harvard.edu/dataverse/WheatBarley40k_v1.
589 Information about the status of each SNP, including tag SNP set ID and whether the SNP passed design
590 filters, is included in the INFO column. Illumina 90k iSelect genotypes for the accessions used to select
591 tetraploid-specific content is available at
592 https://figshare.com/articles/dataset/Durum_Wheat_cv_Svevo_annotation/6984035 (Maccaferri *et*
593 *al.* 2019).

594

595 **Author Contributions**

596 R.P. performed LD analysis. G.K-G selected tagging SNP, performed imputation analyses and produced
597 the final designs. K.F. and D.W. performed exome and whole genome sequencing, Infinium Wheat
598 Barley 40K assays and genotype calling. J.T. performed sequence alignments and genotype calling.

599 H.R., J.G., A.R., D.M., D.M., selected non-tagging SNP and provided wheat and barley germplasm. T.W.,
600 H.D, J.T. and M.H. conceived the project. G.K-G and M.H. wrote the manuscript.

601

602 **Conflict of Interest**

603 The authors declare that they have no conflict of interest.

604

605 **References**

606 Avni R, Nave M, Barad O, Baruch K, Twardziok SO, Gundlach H, *et al.* (2017) Wild emmer genome
607 architecture and diversity elucidate wheat evolution and domestication. *Science* 357, 93-97

608 Allen AM, Winfield MO, BurrIDGE AJ, Downie RC, Benbow HR, Barker GLA, *et al.* (2015) Characterization
609 of a Wheat Breeders' Array suitable for high-throughput SNP genotyping of global accessions of
610 hexaploidy bread wheat (*Triticum aestivum*). *Plant Biotechnology Journal* 15, 390-401 DOI:
611 [10.1111/pbi.12635](https://doi.org/10.1111/pbi.12635)

612 Balfourier F, Bouchet S, Robert S, De Oliveira R, Rimbart H, Kitt J, *et al.* (2019) Worldwide
613 phylogeography and history of wheat genetic diversity. *Science Advances* 5, eaav0536 DOI:
614 [10.1126/sciadv.aav0536](https://doi.org/10.1126/sciadv.aav0536)

615 Bayer MM, Rapazote-Flores P, Ganal M, Hedley PE, Macaulay M, Plieske J, *et al.* (2019) Development
616 and evaluation of a barley 50k iSelect SNP array. *Frontiers in Plant Sciences* 8, 1792 DOI:
617 [10.3389/fpls.2017.01792](https://doi.org/10.3389/fpls.2017.01792)

618 Browning SR, Browning BL (2007) Rapid and accurate haplotype phasing and missing-data inference
619 for whole-genome association studies by use of localized haplotype clustering. *American Journal of*
620 *Human Genetics* 81, 1084-1097 DOI: [10.1086/521987](https://doi.org/10.1086/521987)

621 Carlson CS, Eberle MA, Rieder MJ, Yi Q, Kruglyak L, Nickerson DA (2004) Selecting a maximally
622 informative set of single-nucleotide polymorphisms for association analyses using linkage
623 disequilibrium. *Am J Hum Genet* 74, 106-120 DOI: [10.1086/381000](https://doi.org/10.1086/381000)

624 Cavanagh C, Chao S, Wang S, Huang BE, Stephen S, Kianic S, *et al.* (2013) Genome-wide comparative
625 diversity uncovers multiple targets of selection for improvement in hexaploid wheat landraces and
626 cultivars. *Proceedings of National Academy of Science, USA* 110, 8057-8062 DOI:
627 [10.1073/pnas.1217133110](https://doi.org/10.1073/pnas.1217133110)

628 Chu J, Zhao Y, Beier S, Schulthess AW, Stein N, Philipp N, *et al.* (2020) Suitability of single-nucleotide
629 polymorphism arrays versus genotyping-by-sequencing for genebank genomics in wheat. *Frontiers in*
630 *Plant Science* 14, 11-42 DOI: [10.3389/fpls.2020.00042](https://doi.org/10.3389/fpls.2020.00042)

631 Clark AG, Hubisz MJ, Bustamante CD, Williamson SH, Nielsen R (2005) Ascertainment bias in studies
632 of human genome-wide polymorphism. *Genome Research* 15, 1496-1502 DOI: [10.1101/gr.4107905](https://doi.org/10.1101/gr.4107905)

633 Comadran J, Kilian B, Russell J, Ramsay L, Stein N, Ganal M, *et al.* (2012) Natural variation in a homolog
634 of *Antirrhinum CENTRORADIALIS* contributed to spring growth habit and environmental adaptation in
635 cultivated barley. *Nature Genetics* 44, 1388-1392

636 Ding K, Kullo IJ (2007) Methods for the selection of tagging SNPs: a comparison of tagging efficiency
637 and performance. *European Journal of Human Genetics* 15, 228-236 DOI: [10.1038/sj.ejhg.5201755](https://doi.org/10.1038/sj.ejhg.5201755)

- 638 Das S, Forer L, Schönherr S, Sidore C, Locke A, *et al.* (2016) Next-generation genotype imputation
639 service and methods. *Nature Genetics* 48, 1284-1287 DOI: 10.1038/ng.3656
- 640 Elbasyoni IS, Lorenz AJ, Guttieri M, Frels K, Baenziger PS, Poland J, *et al.* (2018) A comparison between
641 genotyping-by-sequencing and array-based scoring of SNPs for genomic prediction accuracy in winter
642 wheat. *Plant Science Journal* 270, 123-130 DOI: 10.1016/j.plantsci.2018.02.019
- 643 Fikere M, Barbulescu DM, Malmberg MM, Spangenberg GC, Cogan NOI, Daetwyler HD (2020) Meta-
644 analysis of GWAS in canola blackleg (*Leptosphaeria maculans*) disease traits demonstrates increased
645 power from imputed whole-genome sequence. *Scientific Reports* 10, 14300 DOI: 10.1038/s41598-
646 020-71274-6
- 647 He F, Pasam R, Shi F, Kant S, Keeble-Gagnere G, Kay P, *et al.* (2019) Exome sequencing highlights the
648 role of wild-relative introgression in shaping the adaptive landscape of the wheat genome. *Nature*
649 *Genetics* 51, 896-904 DOI: 10.1038/s41588-019-0382-2
- 650 Hill CB, Angessa T, McFawn L-A, Wong D, Tibbits J, Zhang X-Q, *et al.* (2019) Hybridisation-based target
651 enrichment of phenology genes to dissect the genetic basis of yield and adaptation in barley. *Plant*
652 *Biotechnology Journal* 17, 932-944 DOI: 10.1111/pbi.13029
- 653 Hill CB, Angessa TT, Zhang X-Q, Chen K, Zhou G, Tan C, *et al.* (2020) A global barley panel revealing
654 genomic signatures of breeding in modern cultivars. *bioRxiv* DOI: 10.1101/2020.03.04.976324
- 655 Jordan KW, Wang S, Lun Y, Gardiner L-J, MacLauchlan R, Hucl P, *et al.* (2015) A haplotype map of
656 allohexaploid wheat reveals distinct patterns of selection on homoeologous genomes. *Genome*
657 *Biology* 16,48 DOI: 10.1186/s13059-015-0606-4
- 658 Joukhadar R, Daetwyler HD, Bansal UK, Gendall AR, Hayden MJ (2017) Genetic diversity, population
659 structure and ancestral origin of Australian wheat. *Frontiers in Plant Science* 8, 2115 DOI:
660 10.3389/fpls.2017.02115
- 661 Keeble-Gagnère G, Isdale D, Suchecki R, Kruger A, Lomas K, Carroll D, *et al.* (2019) Integrating past,
662 present and future wheat research with Pretzel. *bioRxiv* DOI: 10.1101/517953
- 663 Lai K, Lorenc MT, Lee HC, Berkman PJ, Bayer PE, Visendi P, *et al.* (2015) Identification and
664 characterization of more than 4 million intervarietal SNPs across the group 7 chromosomes of bread
665 wheat. *Plant Biotechnology Journal* 13, 97-104. DOI: 10.1111/pbi.12240
- 666 Ling HQ, Ma B, Shi X, Liu H, Dong L, Sun H, *et al.* (2018) Genome sequence of the progenitor of wheat
667 A subgenome *Triticum urartu*. *Nature* 557, 424-428
- 668 Loh P-R, Danecek P, Palamara PF, Fuchsberger C, Reshef YA, Finucane HK, *et al.* (2016) Reference-
669 based phasing using the Haplotype Reference Consortium panel. *Nature Genetics* 48, 1443-1448 DOI:
670 10.1038/ng.3679
- 671 Luo MC, Gu YQ, Puiu D, Wang H, Twardziok SO, Deal KR, *et al.* (2017) Genome sequence of the
672 progenitor of the wheat D genome *Aegilops tauschii*. *Nature* 551, 498-502
- 673 Maccaferri M, Harris NS, Twardziok SO, Pasam RK, Gundlach H, Spannagl M, *et al.* (2019) Durum wheat
674 genome highlights past domestication signatures and future improvement targets. *Nature Genetics*
675 51,885 DOI: 10.1038/s41588-019-0381-3
- 676 Mascher M, Gundlach H, Himmelbach A, Beier S, Twardziok SO, Wicker T, *et al.* (2017) A chromosome
677 conformation capture ordered sequence of the barley genome. *Nature* 544, 427-433

- 678 Mascher M, Schreiber M, Scholz U, Graner A, Reif JC, Stein N (2019) Genebank genomics bridges the
679 gap between the conservation of crop diversity and plant breeding. *Nature Genetics* 51, 1076-1091
680 DOI: 10.1038/s41588-019-0443-6
- 681 Milner SG, Jost M, Taketa S, *et al.* (2019) Genebank genomics highlights the diversity of a global barley
682 collection. *Nature Genetics* 51, 319-326 DOI: 10.1038/s41588-018-0266-x
- 683 Monat C, Padmarasu S, Lux T, Wicker T, Gundlach H, Himmelbach A, *et al.* (2019) TRITEX: chromosome-
684 scale sequence assembly of Triticeae genomes with open-source tools. *Genome Biology* 20, 284
- 685 Money D, Gardner K, Migicovsky Z, Schwaninger H, Zhong G-Y, Myles S (2015) LinkImpute: Fast and
686 accurate genotype imputation for non-model organisms. *Genes, Genomics, Genetics* 5, 2383-2390
687 DOI: 10.1534/g3.115.021667
- 688 Montenegro JD, Golicz AA, Bayer PE, Hurgobin B, Lee H, Chan C-KK, *et al.* (2017) The pangenome of
689 hexaploidy bread wheat. *Plant Journal* 90, 1007-1013 DOI: 10.1111/tpj.13515
- 690 Negro SS, Millet EJ, Madur D, Bauland C, Combes V, Welcker C, *et al.* (2019) Genotyping-by-sequencing
691 and SNP-arrays are complementary for detecting quantitative trait loci by tagging different haplotypes
692 in association studies. *BMC Plant Biology* 19, 318 DOI: 10.1186/s12870-019-1926-4
- 693 Nyine M, Wang S, Kiani, Jordan K, Liu S, Byrne P, *et al.* (2019) Genotype imputation in winter wheat
694 using first-generation haplotype map SNPs improves genome-wide association mapping and genomic
695 prediction of traits. *Genes, Genomes, Genetics* 9, 125-133 DOI: 10.1534/g3.118.200664
- 696 Ogonnaya FC, Ye G, Trethowan R, Dreccer F, Lush D, Shepperd J, *et al.* (2007) Yield of synthetic
697 backcross-derived lines in rainfed environments of Australia. *Euphytica* 157, 321-336 DOI:
698 10.1007/s10681-007-9381-y
- 699 Pasam RP, Bansal U, Daetwyler HD, Forrest KL, Wong D, Petkowski J, *et al.* (2016) Detection and
700 validation of genomic regions associated with three rust resistances to rust diseases in a worldwide
701 hexaploid wheat landrace collection using BayesR and Mixed Linear Model approaches. *Theoretical*
702 *and Applied Genetics* 130, 777-793 DOI: 10.1007/s00122-016-2851-7
- 703 Pont C, Leroy T, Seidel M, Tondelli A, Duchemin W, Armisen D, *et al.* (2019) Tracing the ancestry of
704 modern bread wheats. *Nature Genetics* 51, 905-911 DOI: 10.1038/s41588-019-0393-z
- 705 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, *et al.* (2007) PLINK: a tool set for
706 whole-genome association and population-based linkage analyses. *American Journal of Human*
707 *Genetics* 81,559-75 DOI: 10.1086/519795
- 708 Rasheed A, Hao Y, Xia X, Khan A, Xu Y, Varshney RK, *et al.* (2017) Crop breeding chips and genotyping
709 platforms: Progress, challenges, and perspectives. *Molecular Plant* 10, 1047-1064 DOI:
710 10.1016/j.molp.2017.06.008
- 711 Rimbart H, Darrier B, Navarro J, Kitt J, Choulet F, Leveugle M, *et al.* (2018) High throughput SNP
712 discovery and genotyping in hexaploid wheat. *PLOS One* 13, e0186329 DOI:
713 10.1371/journal.pone.0186329
- 714 Russell J, Mascher M, Dawson IK, Kyriakidis S, Calixto C, Freund F, *et al.* (2016) Exome sequencing of
715 geographically diverse barley landraces and wild relatives gives insights into environmental
716 adaptation. *Nature Genetics* 48, 1024-1030 DOI: 10.1038/ng.3612

- 717 Shi C, Zhao L, Zhang X, Lv G, Pan Y, Chen F (2019) Gene regulatory network and abundant genetic
718 variation play critical roles in heading stage of polyploidy wheat. *BMC Plant Biology* 19, 6 DOI:
719 10.1186/s12870-018-1591-z
- 720 Soleimani B, Lehnert H, Keilwagen J, Plieske J, Ordon F, Naseri Rad S, *et al.* (2020) Comparison between
721 core set selection methods using different Illumina marker platforms: A case study of assessment of
722 diversity in wheat. *Frontiers in Plant Science* 11,1040 DOI: 10.3389/fpls.2020.01040
- 723 Sun C, Dong Z, Zhao L, Ren Y, Zhang N, Chen F (2020) The Wheat 660K SNP array demonstrates great
724 potential for marker-assisted selection in polyploid wheat. *Plant Biotechnology Journal* DOI:
725 10.1111/pbi.13361
- 726 The International Wheat Genome Sequencing Consortium (IWGSC) (2018) Shifting the limits in wheat
727 research and breeding using a fully annotated reference genome. *Science* 361, eaar7191
- 728 Wang S, Wong D, Forrest K, Allen A, Chao S, Huang BE, *et al.* (2014) Characterization of polyploid wheat
729 genomic diversity using a high-density 90,000 single nucleotide polymorphism array. *Plant*
730 *Biotechnology Journal* 12, 787-96 DOI: 10.1111/pbi.12183
- 731 Walkowiak S, Gao L, Monat C, Haberer G, Kassa MT, Brinton J, *et al.* (2020) Multiple wheat genomes
732 reveal global variation in modern breeding. *Nature* DOI: 10.1038/s41586-020-2961-x
- 733 Weir BS, Cockerham CC (1984) Estimating F-statistics for the analysis of population structure.
734 *Evolution* 38, 1358-1370 DOI: doi.org/10.2307/2408641
- 735 Winfield MO, Allen AM, BurrIDGE AJ, Barker GL, Benbow HR, Wilkinson PA, *et al.* (2016) High-density
736 SNP genotyping array for hexaploid wheat and its secondary and tertiary gene pool. *Plant*
737 *Biotechnology Journal* 14, 1195-1206 DOI: 10.1111/pbi.12485
- 738 Zhu T, Wang L, Rimbert H, Rodriguez JC, Deal KR, De Oliveira R, Choulet F, Keeble-Gagnère G, Tibbits
739 J, Rogers J, Eversole K, Appels R, Gu YQ, Mascher M, Dvorak J, Luo MC. Optical maps refine the bread
740 wheat *Triticum aestivum* cv. Chinese Spring genome assembly. *Plant J.* 2021 Apr 24. doi:
741 10.1111/tpj.15289

742

743 **Supporting Information**

744 **Table S1.** Detailed description of Infinium Wheat Barley 40K SNP array content

745 **Figure S1.** Cumulative number of SNP tagged by tSNP at $r^2 \geq 0.90$ in each chromosome in wheat and
746 barley. Curves shown until the first singleton SNP is reached on each chromosome

747 **Figure S2.** PCA based on (a) 17,600 SNP described in Maccaferri *et al.* (2019) from the Infinium wheat
748 90K SNP array and (b) 2,609 SNP selected for inclusion on the Infinium Wheat Barley 40K SNP array
749 showing differentiation among 1,856 tetraploid wheat accessions representing wild emmer wheat
750 from North Eastern Fertile Crescent (WEW-NE), wild emmer wheat from Southern Levant Fertile
751 Crescent (WEW-SL), domesticated emmer wheat (DEW), domesticated emmer wheat from Ethiopia
752 (DEW-ETH), durum wheat landraces (DWL) and durum wheat cultivars (DWC)

753 **Figure S3.** PCA based on (a) 37,105 called SNP from the Infinium wheat 90K SNP array, and (b) 20,665
754 SNP on the Infinium Wheat Barley 40K SNP array showing differentiation among bread wheat (green),
755 synthetics derivatives (blue) and hexaploid wheat derived from crosses between bread and durum
756 accessions (red) (number of accessions=1219)

758 **Table 1.** Accuracy for imputing from the tSNP on the array to the sets of SNP tagged at $r^2 \geq 0.50$, 0.70
 759 and 0.90 respectively in wheat and barley. Correlation is the Pearson r^2 between SNP called in both
 760 genotypes being compared. Concordance is the fraction of SNP in agreement between those called in
 761 both genotypes being compared. Standard deviations are shown in brackets

	Set of SNP tagged at r^2	Wheat	Barley
Correlation (including heterozygous calls)	0.50	93.7 (4.0)	86.0 (3.1)
	0.70	95.3 (3.8)	92.4 (2.6)
	0.90	97.0 (3.4)	96.8 (1.6)
Correlation (excluding heterozygous calls)	0.50	97.6 (1.3)	91.5 (2.9)
	0.70	98.7 (1.0)	96.9 (2.3)
	0.90	99.3 (0.7)	98.7 (1.3)
Concordance (including heterozygous calls)	0.50	96.9 (2.2)	92.8 (1.4)
	0.70	97.4 (2.1)	95.2 (1.2)
	0.90	98.3 (2.0)	98.1 (0.8)
Concordance (excluding heterozygous calls)	0.50	99.6 (0.2)	99.7 (0.3)
	0.70	99.8 (0.2)	99.3 (0.5)
	0.90	99.9 (0.1)	99.7 (0.2)

762

763

764 **Table 2.** SNP content of the Infinium Wheat Barley 40K SNP bead chip array

	Wheat	Barley	Total
Tagging SNP for imputation	21,012	13,469	34,481
Trait associated SNP	427	178	605
SNP linking germplasm resources	3,924	614	4,538
Total number of SNP	25,363	14,261	39,624

765

766

767 **Figure Legends**

768 **Figure 1.** PCA plots showing genetic diversity of wheat and barley accessions used for SNP discovery.
769 (a) 6,087 wheat accessions genotyped with the iSelect wheat 90K SNP array (Wang *et al.* 2014) (black),
770 exome-sequenced accessions used for LD analysis (red) and synthetic derivative accessions capturing
771 D-genome diversity (blue); and (b) 19,778 barley accessions genotyped with GBS (black), with exome-
772 sequenced accessions used for LD analysis (red).

773 **Figure 2.** Cumulative number of SNP tagged by tSNP at $r^2 \geq 0.9$, 0.7 and 0.5 respectively in wheat and
774 barley. The curves are shown until the first singleton SNP (at $r^2 \geq 0.90$) is reached

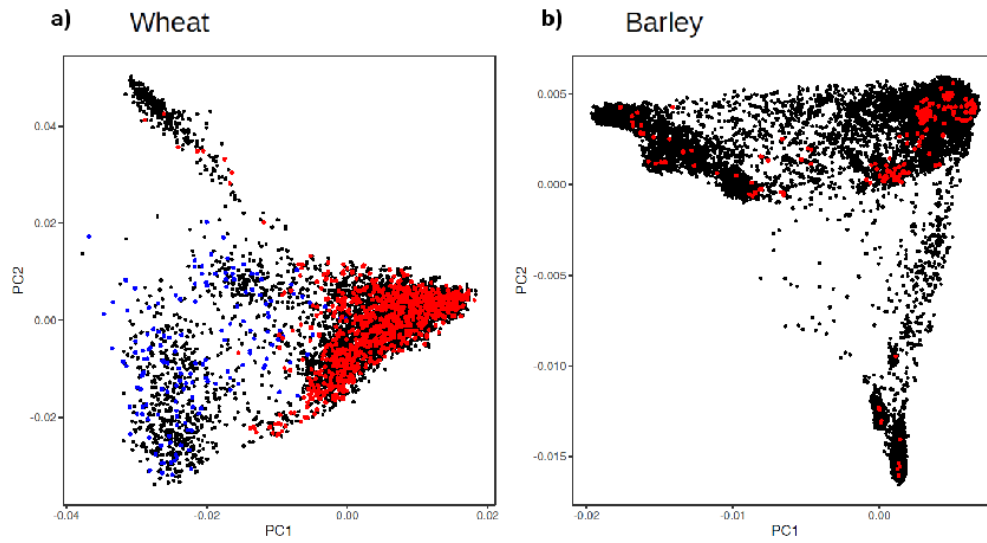
775 **Figure 3.** MAF distribution of all SNP used for LD analysis, selected tSNP and the set of SNP tagged by
776 the tSNP at $r^2 \geq 0.70$ in the globally diverse wheat (n=790) and barley (n=157) collections

777 **Figure 4.** Imputation accuracy from the tSNP on the array to the set of SNP tagged at $r^2 \geq 0.5$, 0.7 and
778 0.9 respectively in wheat and barley. Metrics plotted are correlation r^2 including heterozygous calls
779 (orange line), r^2 excluding heterozygous calls (cyan line), concordance including heterozygous calls
780 (green line) and concordance excluding heterozygous calls (orange line). The accessions are ranked
781 ordered based on the r^2 including heterozygous calls

782 **Figure 5.** Cluster positions and theta separation of SNP in single sample hybridisation assays. Scatter
783 plot of cluster positions (left) and density plot of difference in theta value between REF and ALT
784 clusters (right) for (a) 14,261 barley and (b) 24,598 wheat SNP revealing polymorphism in the globally
785 diverse wheat and barley populations

786 **Figure 6.** Cluster positions and theta separation of SNPs in dual hybridisation assays. Scatter plot of
787 cluster positions (left) and density plot of difference in theta value between REF and ALT clusters
788 (right) for (a) 9,826 barley and (b) 9,118 wheat SNP revealing polymorphism among 576 wheat and
789 barley breeding lines

790



791

792

793 **Figure 1.** PCA plots showing genetic diversity of wheat and barley accessions used for SNP discovery.

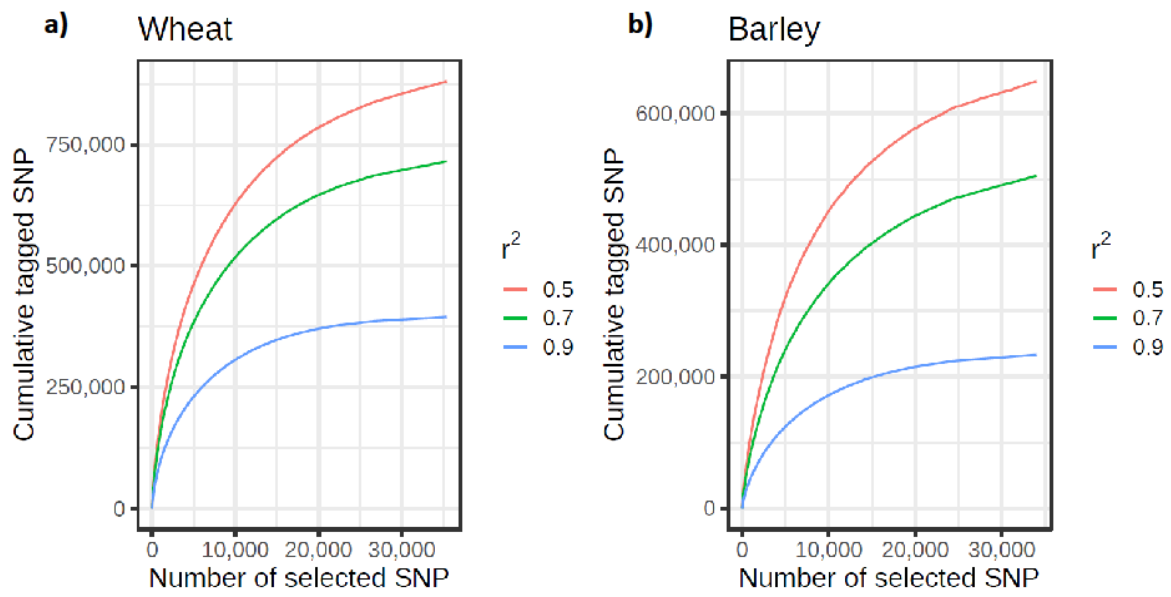
794 (a) 6,087 wheat accessions genotyped with the iSelect wheat 90K SNP array (Wang *et al.* 2014) (black),

795 exome-sequenced accessions used for LD analysis (red) and synthetic derivative accessions capturing

796 D-genome diversity (blue); and (b) 19,778 barley accessions genotyped with GBS (black), with exome-

797 sequenced accessions used for LD analysis (red).

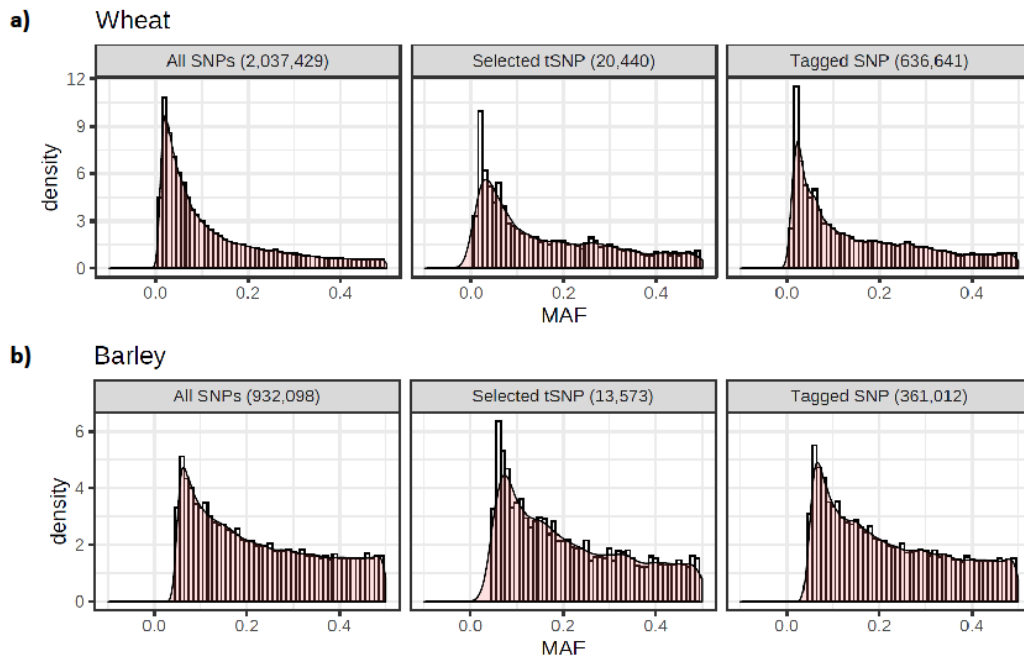
797



798

799 **Figure 2.** Cumulative number of SNP tagged by tSNP at $r^2 \geq 0.9$, 0.7 and 0.5 respectively in wheat and
800 barley. The curves are shown until the first singleton SNP (at $r^2 \geq 0.90$) is reached

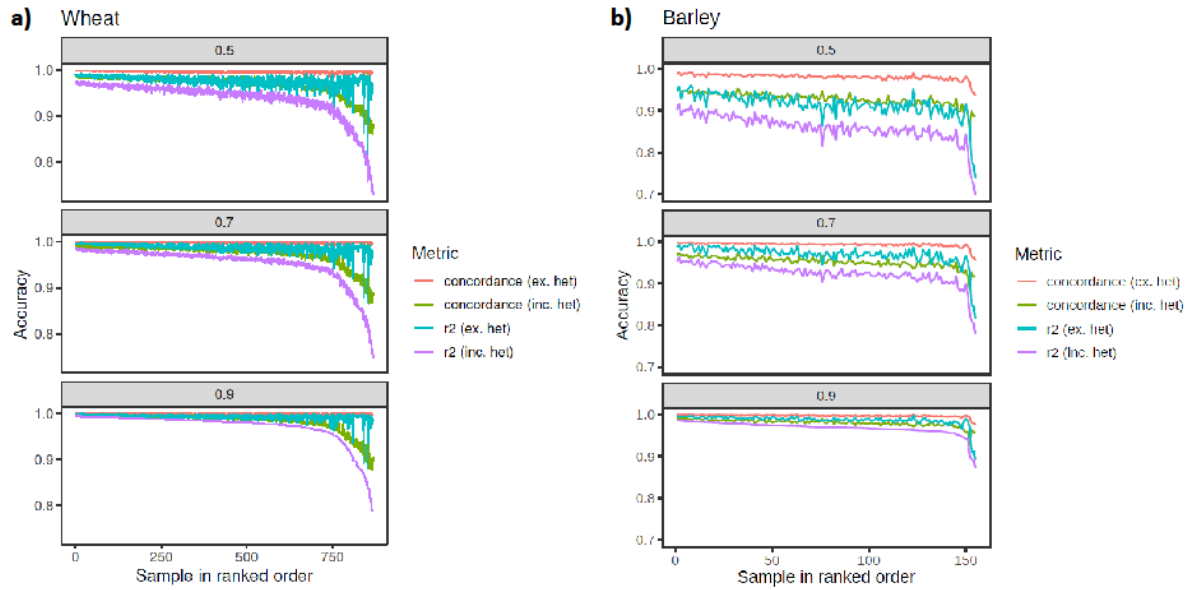
801



802

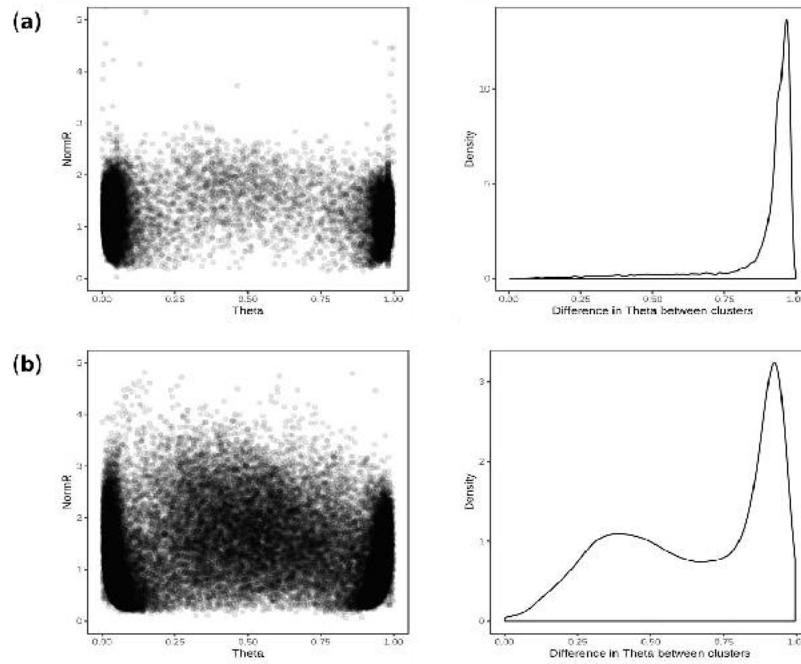
803 **Figure 3.** MAF distribution of all SNP used for LD analysis, selected tSNP and the set of SNP tagged by
804 the tSNP at $r^2 \geq 0.70$ in the globally diverse wheat ($n=790$) and barley ($n=157$) collections

805



806
807 **Figure 4.** Imputation accuracy from the tSNP on the array to the set of SNP tagged at $r^2 \geq 0.5$, 0.7 and
808 0.9 respectively in wheat and barley. Metrics plotted are correlation r^2 including heterozygous calls
809 (orange line), r^2 excluding heterozygous calls (cyan line), concordance including heterozygous calls
810 (green line) and concordance excluding heterozygous calls (orange line). The accessions are ranked
811 ordered based on the r^2 including heterozygous calls

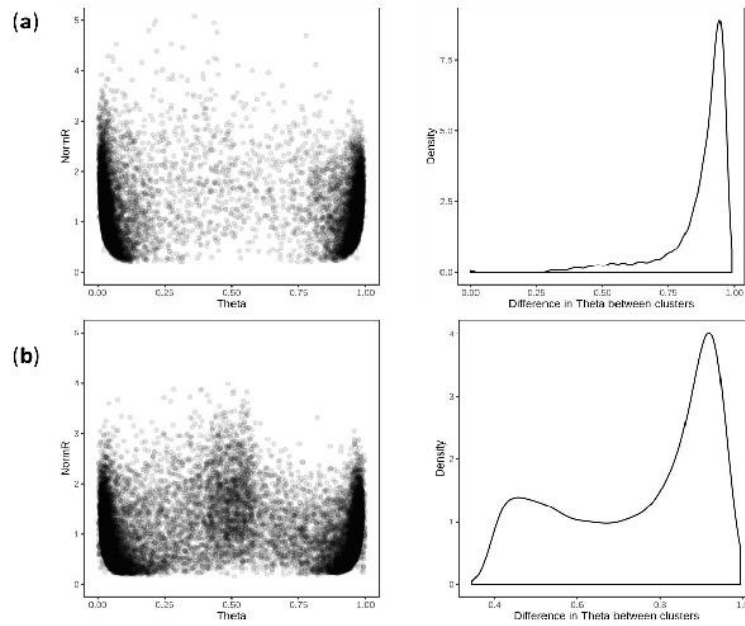
812



813

814 **Figure 5.** Cluster positions and theta separation of SNP in single sample hybridisation assays. Scatter
815 plot of cluster positions (**left**) and density plot of difference in theta value between REF and ALT
816 clusters (**right**) for **(a)** 14,261 barley and **(b)** 24,598 wheat SNP revealing polymorphism in the globally
817 diverse wheat and barley populations

818



819

820 **Figure 6.** Cluster positions and theta separation of SNPs in dual hybridisation assays. Scatter plot of
821 cluster positions (**left**) and density plot of difference in theta value between REF and ALT clusters
822 (**right**) for **(a)** 9,826 barley and **(b)** 9,118 wheat SNP revealing polymorphism among 576 wheat and
823 barley breeding lines

824