

1 Article type: Research article

2 ***TaAPO-A1*, an ortholog of rice *ABERRANT PANICLE ORGANIZATION 1*, is**  
3 **associated with total spikelet number per spike in elite hexaploid winter wheat**  
4 **varieties (*Triticum aestivum* L.)**

5 Quddoos H. Muqaddasi<sup>1\*</sup>, Jonathan Brassac<sup>1</sup>, Ravi Koppolu<sup>1</sup>, Jörg Plieske<sup>2</sup>,  
6 Martin W. Ganal<sup>2</sup>, and Marion S. Röder<sup>1</sup>

7 <sup>1</sup>Leibniz Institute of Plant Genetics and Crop Plant Research (IPK), Corrensstraße 3,  
8 D-06466 Stadt Seeland, OT Gatersleben, Germany.

9 <sup>2</sup>SGS TraitGenetics, Am Schwabeplan 1b, D-06466 Stadt Seeland, OT Gatersleben,  
10 Germany.

11 \*Corresponding author: [muqaddasi@ipk-gatersleben.de](mailto:muqaddasi@ipk-gatersleben.de)

12 Date of submission: May 30, 2019

13 Number of tables: 1

14 Number of figures: 6

15 Number of supplementary tables: 4

16 Number of supplementary figures: 7

17 **Abstract**

18 We dissected the genetic basis of total spikelet number (TSN) along with other traits,  
19 namely spike length (SL) and flowering time (FT) in a panel of 518 elite European  
20 winter wheat varieties. Genome-wide association studies based on 39,908 SNP  
21 markers revealed highly significant quantitative trait loci (QTL) for TSN on  
22 chromosomes 2D, 7A, and 7B, for SL on 5A, and FT on 2D, with 2D-QTL being the  
23 functional marker for the gene *Ppd-D1*. The physical region of the 7A-QTL for TSN  
24 revealed the presence of an ortholog to *APO1* – a rice gene that positively controls  
25 spikelet number on panicles. Interspecific analyses of *TaAPO-A1* orthologs showed  
26 that it is a highly conserved gene important for floral development, and present in a  
27 wide range of terrestrial plants. Intraspecific studies of the wheat ortholog *TaAPO-A1*  
28 across wheat genotypes revealed a polymorphism in the highly conserved F-box  
29 domain, defining two haplotypes. A KASP marker developed on the polymorphic site  
30 showed a highly significant association of *TaAPO-A1* with TSN, explaining 23.2% of  
31 the genotypic variance. Also, the *TaAPO-A1* alleles showed weak but significant  
32 differences for SL and grain yield. Our results demonstrate the importance of wheat  
33 sequence resources to identify candidate genes for important traits based on genetic  
34 analyses.

35 **Keywords:** Wheat; total spikelet number; GWAS; QTL; physical mapping; *TaAPO-*  
36 *A1*

## 37 Introduction

38 The wheat spike and its architecture are key components for improving grain yield. In  
39 the recent past, several genes controlling spike morphology have been investigated  
40 and described in temperate cereals (Gauley and Boden 2019; Koppolu and  
41 Schnurbusch 2019). From a plant breeder's viewpoint, most spike morphological  
42 traits in wheat such as spike length and spikelet number behave as quantitative traits,  
43 and various QTL and association studies have recently been published (Deng et al.  
44 2017; Guo et al. 2017; Liu et al. 2018; Sakuma et al. 2019; Würschum et al. 2018;  
45 Zhai et al. 2016). High associations and prediction abilities for total and fertile spikelet  
46 number as well as spike length and grain yield were also reported (Guo et al. 2018).

47 Nevertheless, only a few cloned genes for the trait number of spikelet pairs in  
48 wheat are available, among them is the *Q*-gene which played a major role in wheat  
49 domestication and encodes an *AP2* transcription factor (Faris et al. 2003). The  
50 domesticated allele *Q* confers a free-threshing character, a sub-compact spike  
51 (Greenwood et al. 2017), and is regulated by microRNA172 (Debernardi et al. 2017).  
52 Also, genes related to heading date are involved in spikelet meristem identity  
53 determination. For example, the photoperiodism gene *Ppd* was reported to influence  
54 spikelet primordia initiation (Ochagavía et al. 2018). Mutants of the *FLOWERING*  
55 *LOCUS T2 (FT2)* in wheat showed a significant increase in the number of spikelets  
56 per spike with an extended spike development period accompanied by delayed  
57 heading time (Shaw et al. 2018). Moreover, *Ppd-1* and *FT* were reported as  
58 regulators of paired spikelet formation resulting in increased number of grain  
59 producing spikelets (Boden et al. 2015). Mutants of the MADS-box genes, e.g., *VRN1*  
60 or *FUL2* showed increased number of spikelets per spike, likely due to a delayed  
61 formation of the terminal spikelet (Li et al. 2019) and a putative ortholog to rice *MOC1*  
62 regulating axillary meristem initiation and outgrowth was associated with spikelet  
63 number per spike in wheat (Zhang et al. 2015).

64 The *ABERRANT PANICLE ORGANIZATION 1 (APO1)* gene in rice was  
65 reported to affect the inflorescence structure severely (Ikeda et al. 2005). It encodes  
66 an F-box protein which is an ortholog of *UNUSUAL FLORAL ORGAN (UFO)*,  
67 regulating floral identity in *Arabidopsis* (Samach et al. 1999; Wilkinson and Haughn  
68 1995). Characterization of rice *apo1* mutants revealed that *APO1* positively controls  
69 spikelet number by suppressing the precocious conversion of inflorescence  
70 meristems to spikelet meristems. Besides this, *APO1* was associated with the  
71 regulation of the plastochron, floral organ identity, and floral determinacy (Ikeda et al.  
72 2007). Four dominant mutants of *APO1* with elevated expression levels of *APO1*  
73 produced increased number of spikelets by a delay in the programmed shift to  
74 spikelet formation. Ectopic overexpression of *APO1* resulted in increased meristem  
75 size caused by different rates of cell proliferation. It was concluded that the level of  
76 *APO1* activity regulates the inflorescence form through the control of meristematic  
77 cell proliferation (Ikeda-Kawakatsu et al. 2009).

78 In the present study, we investigated the inheritance and genetic basis of total  
79 spikelet number (TSN) per spike, spike length and flowering time as component traits  
80 of grain yield in an elite European winter wheat panel. Our findings show the complex  
81 genetic architecture of the investigated traits, and that *TaAPO-A1* – an ortholog of  
82 rice *APO1*, which is vital for inflorescence development, is associated with TSN  
83 determination in wheat. Intraspecific sequence analyses of *TaAPO-A1* revealed that  
84 polymorphisms were forming distinct haplotypes while intraspecific studies showed  
85 the conserved nature of this gene across terrestrial plant species.

## 86 **Materials and methods**

### 87 *Phenotypic data analyses*

88 The data for total spikelet number (TSN), spike length (SL), and flowering time (FT)  
89 were collected on an elite European winter wheat panel comprising of 518 varieties.  
90 The whole panel was grown in the experimental fields of Leibniz Institute of Plant  
91 Genetics and Crop Plant Research (IPK) Gatersleben, Germany in plots of 2 m<sup>2</sup> as  
92 single replication in three cropping seasons (2015/16; 2016/17; and 2017/18),  
93 henceforth called environments. The traits TSN and SL were recorded in two  
94 environments (2016/17 and 2017/18) from ten spikes per plot as the total number of  
95 spikelets and spike length in centimeters (cm) from basal spikelet to the top of a  
96 spike by excluding the awns. The arithmetic mean of TSN and SL from ten spikes  
97 were calculated to represent the genetic value of traits in the individual environments.  
98 Flowering time was recorded in all three environments by counting the number of  
99 days from the first of January to when approximately half of the spikes in a plot  
100 flowered. The phenotypic data for grain yield estimated in eight environments were  
101 taken from the previous study for comparison purposes (Schulthess et al. 2017). A  
102 linear mixed-effect model was used for across environment phenotypic data analysis  
103 as:

$$104 \quad y_{ij} = \mu + G_i + E_j + e_{ij}$$

105 where,  $y_{ik}$  is the phenotypic record of the  $i^{th}$  genotype in the  $j^{th}$  environment,  $\mu$  is  
106 the common intercept term,  $G_i$  is the effect of the  $i^{th}$  genotype,  $E_j$  is the effect of the  
107  $j^{th}$  environment, and  $e_{ij}$  denotes the corresponding error term. All effects, except the  
108 intercept, were assumed random to calculate the individual variance components.  
109 The broad-sense heritability ( $H^2$ ) was calculated as:

$$110 \quad H^2 = \frac{\sigma_G^2}{\sigma_G^2 + \left(\frac{\sigma_e^2}{nE}\right)}$$

111 where,  $\sigma_G^2$  and  $\sigma_e^2$  denote the variance components of the genotype and the error,  
112 respectively; and  $nE$  denotes the number of environments. To calculate the best  
113 linear unbiased estimations (BLUEs), the intercept and the genotypic effects were  
114 assumed fixed in the above model.

## 115 *Genotypic data analyses, population structure, and linkage disequilibrium*

116 All 518 varieties were extensively genotyped with the 35k Affymetrix and 90k  
117 iSELECT single nucleotide polymorphism (SNP) arrays (Allen et al. 2017; Wang et al.  
118 2014) which generated in total 116,730 SNP markers (35k = 35,143; 90k = 81,587).  
119 We also genotyped the whole panel with functional markers for the candidate genes  
120 such as photoperiodism (*Ppd-D1*), reduced height (*Rht*), and vernalization (*Vrn1*).  
121 The quality of the marker data was improved by removing the markers harboring  
122 >10% heterozygous or missing calls and markers with a minor allele frequency of  
123 <0.05. The mean of both alleles imputed the remaining missing data. The quality  
124 control resulted in a total of 39,908 markers, which were used in subsequent  
125 analyses.

126 Population structure based on marker genotypes was examined by principal  
127 component (PC) analysis. The first two PCs were drawn to see the clustering among  
128 varieties. Moreover, the genetic relatedness among varieties was evaluated by an  
129 additive variance-covariance genomic relationship matrix. To infer the hidden  
130 population sub-structuring, an inference algorithm LEA (Landscape and Ecological  
131 Association Studies) was used by assuming ten ancestral populations ( $K = 1-10$ ).  
132 The function *snmf*, which provides the least squares estimates of ancestry  
133 proportions and estimates an entropy criterion to evaluate the quality fit of the model  
134 by cross-validation, was used. The number of ancestral populations best explaining  
135 the data can be chosen by using the entropy criterion. We performed ten repetitions  
136 for each  $K$ , and the optimal repetition demonstrating the minimum cross-entropy  
137 value was used to visualize clustering among varieties via bar plots (Frichot and  
138 François 2015).

139 Linkage disequilibrium (LD), the non-random association of alleles at different  
140 loci, was measured as the squared correlation ( $r^2$ ) among markers. The genetic  
141 mapping positions of the markers for both arrays were adopted from the data  
142 generated for the International Triticeae Mapping Initiative (ITMI) DH population, as  
143 described in Sorrells et al. (2011). Although inter and intra-chromosomal LD among  
144 the loci varies, genome-wide calculation of LD gives a global estimate about the  
145 genetic map distance over which LD decays in the given population. The genome-  
146 wide (global) LD was calculated only from the mapped markers.

## 147 *Genome-wide association studies*

148 Genome-wide association studies (GWAS) were performed on data taken from the  
149 individual environment and SNPs passing the quality criteria *plus* the functional gene  
150 markers. Let  $n$  be the number of varieties and  $p$  be the predictor marker genotypes. A  
151 standard linear mixed-effect model following Yu et al. (2006) was used to perform  
152 GWAS as:

$$153 \quad y = \mu + E\tau + X\beta + Pv + Zu + e$$

154 where,  $y$  is the  $n \times 1$  vector of phenotypic record of each genotype in each  
155 environment,  $\mu$  is the common intercept,  $\tau, \beta, v, u$ , and  $e$  are the vectors of the  
156 environment, marker, population (principal components), polygenic background, and  
157 the error effects, respectively;  $E, X, P$ , and  $Z$  are the corresponding design matrices.  
158 In the model,  $\mu, \tau, \beta$ , and  $v$  were assumed fixed while  $u$  and  $e$  as random with  
159  $u \sim N(0, G\sigma_a^2)$  and  $e \sim N(0, I\sigma_e^2)$ . The  $n \times n$  variance-covariance additive relationship  
160 matrix ( $G$ ) was calculated from  $n \times p$  matrix  $W = (w_{ik})$  of marker genotypes (being 0,  
161 1 or 2) as  $G = \frac{\sum_{k=1}^p (w_{ik} - 2p_k)(w_{jk} - 2p_k)}{2 \sum_{k=1}^p p_k(1-p_k)}$  where,  $w_{ik}$  and  $w_{jk}$  are the profiles of the  $k^{th}$   
162 marker for the  $i^{th}$  and  $j^{th}$  variety, respectively;  $p_k$  is the estimated frequency of one  
163 allele in  $k^{th}$  marker, as described by VanRaden (2008).

164 As population stratification and familial relatedness can severely impact the  
165 power to detect true marker-trait association (MTA) in GWAS, different statistical  
166 models were used to avoid spurious MTA viz., (1) general linear model (*naive*), (2)  
167 population structure correction via principal components (*PCs*), (3) correction of  
168 familial relatedness via genomic relationship matrix ( $G$ ), and (4) correction of  
169 population structure and relatedness via *PCs* and  $G$ . It is expected that using both  
170 *PCs* and  $G$ , in the model can enhance the accuracy of GWAS. Along with this,  
171 environmental fixed effects were assigned in all model scenarios. The models  
172 described above were compared by plotting expected versus the observed  
173  $-\log_{10}(P - \text{value})$  in a quantile-quantile plot and the best model was determined by  
174 checking how well the observed  $-\log_{10}(P - \text{value})$  aligned with the expected.

175 To declare the presence of MTA, a false discovery rate (FDR)  $< 0.05$  to  
176 account for multiple testing was applied (Benjamini and Hochberg 1995). Following  
177 Utz et al. (2000), the percentage of total genotypic variance ( $p_G$ ) explained by all the  
178 QTL passing the FDR threshold was determined as  $p_G = [R_{adj}^2 / H^2] \times 100$  where,  
179  $R_{adj}^2$  was calculated by fitting all the MTA in a multiple linear regression model in the  
180 order of ascending  $P$ -values and  $H^2$  is the broad-sense heritability. The  $p_G$  values of  
181 individual QTL were accordingly derived from the sum of squares of the QTL ( $SS_{QTL}$ )  
182 in the linear model.

183 *Candidate gene identification, haplotype analysis by exploiting resources from The*  
184 *10+ Wheat Genome Project, and KASP marker development*

185 We narrowed-down the QTL region, and BLASTed sequences of all the significant  
186 markers present within the genetically defined region onto the physical map of the  
187 corresponding chromosome of the reference sequence of the wheat genome which  
188 yielded significant physical region (Altschul et al. 1990; Consortium 2018). Afterward,  
189 the gene identifiers (gene-IDs) present within the physical region and their annotated  
190 functional descriptions were retrieved. Among them was a most likely candidate gene  
191 *TaAPO-A1* for TSN.

192 *The 10+ Wheat Genome Project* is an international collaborative effort that  
193 aims to assemble the genomes of more than ten wheat varieties bred in different



194 countries to characterize the wheat pan-genome  
195 (<http://www.10wheatgenomes.com/>). We retrieved the genomic sequence of *TaAPO-*  
196 *A1* for ten wheat varieties from *The 10+ Wheat Genome Project* and aligned the  
197 sequences to observe the haplotype structures. The SNP that revealed a clear  
198 haplotype structure was used to design a Kompetitive Allele Specific PCR (KASP)  
199 marker in the candidate gene. The allele-wise phenotypic distribution of the  
200 investigated traits with the gene-specific KASP marker was analyzed by plotting the  
201 boxplots. The significance (*P*-values) between the mean values of genotypes  
202 harboring different KASP marker alleles was determined by two-sided t-test.  
203 Moreover, we performed a second round of GWAS by incorporating the gene-specific  
204 KASP marker in the original SNP matrix to determine whether it associates with the  
205 phenotypes. The GWAS parameters were kept the same as described above.

### 206 *Multiple sequence alignment and phylogenetic analyses*

207 The *TaAPO-A1* protein sequence (corresponding to *TraesCS7A01G481600*) was  
208 used as a BLAST query to retrieve the monocot, dicot and Bryophyte orthologs from  
209 EnsemblPlants (<http://plants.ensembl.org/index.html>) and Phytozome v12.1  
210 (<https://phytozome.jgi.doe.gov/pz/portal.html>) databases. The orthologous protein  
211 sequences were aligned using ClustalW in Geneious v.11.0.5 (Kearse et al. 2012).  
212 The protein alignment was used to infer a maximum likelihood (ML) phylogeny. The  
213 JTT matrix (Jones et al. 1992) was identified as the best-fitting model of protein  
214 evolution with ProtTest 3 (Darriba et al. 2011; Guindon and Gascuel 2003) and the  
215 Akaike Information Criterion (AIC). The evolutionary history among *TaAPO-A1*  
216 orthologs across various plant species was inferred using RAXML v8.2.12  
217 (Stamatakis 2014) with PROTGAMMAJTT model, rapid bootstrapping of 100  
218 replicates, and search for best-scoring ML tree (options “-f a -x 1 -# 100”). The  
219 consensus tree was further processed to collapse branches with bootstrap support  
220 lower than 50%, and the tree was rooted with the Bryophytes *Physcomitrella patens*  
221 and *Selaginella moellendorffii* as outgroup.

## 222 **Results**

### 223 *Total spikelet number per spike is significantly correlated with spike length, flowering* 224 *time, and grain yield*

225 The assessment of total spikelet number (TSN) per spike, spike length (SL), and  
226 flowering time (FT) were performed in the field trials on 518 elite European winter  
227 wheat varieties (including 15 spring type wheat varieties as an outgroup). The trait  
228 grain yield (GY) was assessed in multiple environment field trials on a subset (in  
229 total, 372) of varieties in a previous study (Schulthess et al. 2017). The best linear  
230 unbiased estimations (BLUEs) of all traits approximated normal distribution and  
231 showed wide variation (Fig. 1a-d; Table S1). The ANOVA showed that genotypic ( $\sigma_G^2$ )  
232 and environmental ( $\sigma_E^2$ ) variation was significantly ( $P < 0.001$ ) larger than zero  
233 (Table 1). The broad-sense heritability ranged from 0.68 to 0.89 which indicates the  
234 good quality of the phenotypic data and its potential for use in genome-wide

235 association (GWAS) studies to map the quantitative trait loci (QTL) underlying the  
236 traits (Table 1). We analyzed the Pearson product moment correlation ( $r$ ) among the  
237 BLUEs of investigated traits, which revealed that TSN was positively and significantly  
238 correlated with SL, FT, and GY (Fig. 1e). The TSN and SL showed the highest  
239 correlation among the analyzed traits ( $r = 0.46; P < 0.001$ ) whereas SL showed  
240 almost a null correlation with FT and GY suggesting that FT augments GY mainly by  
241 influencing TSN in wheat.

#### 242 *High-density marker arrays reveal the absence of distinct sub-populations and sharp* 243 *LD decay*

244 The whole wheat panel was extensively genotyped with high-density SNP arrays and  
245 functional markers for the genes *Ppd-D1*, *Rht-B1*, *Rht-D1*, *Vrn-A1*, *Vrn-B1*, and *Vrn-*  
246 *D1*, which resulted in 39,908 high-quality markers. The population structure analyzed  
247 with marker genotypes by PC analysis resulted in the absence of distinct sub-  
248 populations with the first two PCs representing only 11.3% of the variation (Fig. 2).  
249 The high familial relatedness and non-existence of distinct sub-populations were  
250 further supported by plotting a heat map of the genomic relationships among the  
251 wheat varieties (Fig. S1) and by the structure-like inference algorithm LEA, which  
252 resulted in the sub-populations being distinguished but with a slight entropy shift. The  
253 bar plots indicated admixed and weak sub-populations (Fig. S2).

254 Linkage disequilibrium (LD) between the marker genotypes determines the  
255 number of markers needed to perform GWAS. Genome-wide LD was performed with  
256 the mapped marker genotypes which resulted in rapid LD decay with increasing the  
257 genetic map (cM) distances, with first and third quantile dropping to 0.002 and 0.028,  
258 respectively; and the mean and median values equaling 0.051 and 0.008,  
259 respectively (Fig. 3a). The sub-genome-wise distribution of the markers varied, with  
260 the highest markers mapping on B-genome, followed by A- and D-genomes (Fig. 3b).  
261 Although the whole panel was genotyped with state-of-the-art genotyping arrays, the  
262 sub-genome-wise distribution of marker genotypes suggests that marker density  
263 could be improved especially for D-genome.

#### 264 *GWAS identifies large-effect QTL for TSN on chromosome 7A in wheat varieties*

265 Among the different GWAS models used in our study, we observed that the  $PC_{[1-3]}+G$   
266 model could best control the spurious MTA. Our GWAS analyses identified QTL on  
267 chromosomes 2D, 7A, and 7B for TSN (Fig. 4a-b; Table S1a), for SL on chromosome  
268 5A (Fig. S3, Table S1b), and for FT on chromosome 2D (Fig. S4; Table S1c). The  
269 QTL on chromosome 2D identified for TSN and FT was very likely the gene *Ppd-D1*.  
270 Of particular interest is the photoperiod insensitive allele *Ppd-D1a* that significantly  
271 reduced the TSN (Fig. 4a, f). The phenotypic data for GY were analyzed to  
272 investigate if there exists any significant correlation between the identified marker  
273 alleles and GY. The total proportion of genotypic variance ( $p_G$ ) imparted by the  
274 identified QTL amounted to 65.44% for TSN, 15.15% for SL, and 31.58% for FT. A

275 relatively low  $p_G$  for SL and FT is the result of the identification of only one mapped  
276 MTA for each trait.

277 Nevertheless, of interest is the large-effect QTL identified for TSN on  
278 chromosome 7A – for which the most significant marker *AX-95173991* is located at  
279 112.10 cM and explained 25.70% of the genotypic variance. This warrants, on the  
280 one hand, that the use of 7A-QTL would be beneficial for efficient marker-assisted  
281 selection. On the other hand, it made possible the further investigation of 7A-QTL at  
282 the physical sequence level to search for candidate genes.

283 *Significant physical region of chromosome 7A-QTL harbors TaAPO-A1 – a putative*  
284 *candidate gene for TSN in wheat varieties*

285 The significant 7A-QTL region for TSN spanned initially from 110.6 to 124.1 cM  
286 (Table S1a). We narrowed down the genetic region with the highly significant MTA  
287 with  $-\log_{10}(P - \text{value}) > 10$  within 2.3 cM starting from 111.3 to 113.6 cM (Fig. 4c).  
288 The alignment of marker sequences present within this most significant genetic  
289 region onto chromosome 7A revealed a physical region starting from 673.75 to  
290 674.30 Mb (Fig. 4d) that harbored only ten genes. The functional annotations of  
291 these ten genes revealed an interesting candidate gene *TraesCS7A01G481600*;  
292 (physical map position: 674,081,462 – 674,082,919 bp) with functional annotation as  
293 *Aberrant panicle organization 1 (APO1) protein*. The *APO1* in rice regulates  
294 inflorescence architecture and positively controls the total spikelet number by  
295 suppressing the precocious conversion of inflorescence meristems to spikelet  
296 meristems (Ikeda et al. 2007; Ikeda et al. 2005).

297 *A KASP marker developed for TaAPO-A1 shows significant association with TSN in*  
298 *wheat varieties*

299 *TaAPO-A1* is a 1,457 bp long gene, and like *APO1* in rice, it has two exons  
300 separated by one intron (Fig. 4e). We investigated the variation of *TaAPO-A1* in ten  
301 wheat varieties; the sequences were taken from *The 10+ Wheat Genome Project*,  
302 which revealed two haplotypes (Fig. 4e; Fig. S6). The first exon harbors a highly  
303 conserved F-box domain of 46 amino acid residues across the wheat varieties and  
304 other species (Figs. 4e, S6, and S7). Intraspecific sequence analysis of *TaAPO-A1*  
305 revealed a non-synonymous mutation in the F-box domain; and out of ten wheat  
306 varieties, four (including Chinese Spring) harbored T, while six had G allele. We  
307 developed a KASP marker for *TaAPO-A1* harboring this non-synonymous mutation in  
308 the F-box domain (Figs. 4e and S6; Table S1a, b). The alleles of the KASP marker  
309 were evenly distributed in the variety panel (Figure 2b) and were highly significantly  
310 associated with TSN (Fig. 4f; Table S2a). The second round of GWAS was  
311 performed by the *TaAPO-A1* KASP marker integrated into the original SNP matrix  
312 which further confirmed the significant association of *TaAPO-A1* with TSN, explaining  
313 23.21% of the genotypic variance (Fig. S5; Table S2). The reference allele in the  
314 population (represented by *TaAPO-A1a*, with nucleotide G translating to cysteine)  
315 was present in 50.62% of the investigated varieties and resulted in an average TSN



316 of 18.83. Whereas the variant allele (represented by *TaAPO-A1b*, with nucleotide T  
317 translating to phenylalanine) was present in 49.38% of the varieties and revealed an  
318 average TSN of 20.13 (Fig. 4f, Table S1a). The analysis of local linkage  
319 disequilibrium performed with the markers present in the 7A-QTL genetic region and  
320 the KASP marker for *TaAPO-A1* showed that *TaAPO-A1* was in tight linkage with  
321 other markers (Fig. 5). Furthermore, we also observed a rather weak but significant  
322 association of the *TaAPO-A1* KASP marker alleles with SL, FT, and GY (Fig. 4g-i).

323 The single nucleotide substitution G (low TSN allele) to T (high TSN allele) in  
324 the conserved functional domain of *TaAPO-A1* resulted in a non-synonymous amino  
325 acid substitution from cysteine (C) to phenylalanine (F). The amino acid cysteine  
326 appears to be well conserved across various grass species at this position potentially  
327 indicating the conservation of C residue across grasses. However, the SIFT (Sorting  
328 Intolerant from Tolerant) score (Sim et al. 2012) analysis showed no potential  
329 deleterious effect from C to F substitution at this position (Table S3). We then looked  
330 at the promoter region of *TaAPO-A1* in ten genotypes from *The 10+ Wheat Genome*  
331 *Project* and identified a 115 bp INDEL (insertion-deletion) polymorphism at -484 bp  
332 upstream of the transcription start site of *TaAPO-A1*. Interestingly, the low TSN  
333 haplotype “G” (coding for cysteine) always had a deletion of 115 bp in the promoter,  
334 whereas the high TSN haplotype “T” (coding for phenylalanine) had 115 bp insertion.  
335 It, nevertheless, remains to be established via functional studies if this INDEL affects  
336 the transcription rate of *TaAPO-A1* contributing to the observed phenotypic  
337 differences for TSN in two haplogroups.

338 *Phylogenetic analyses show that TaAPO-A1, an ortholog of UFO in Arabidopsis, is*  
339 *conserved across terrestrial plant species*

340 The BLAST search of *TaAPO-A1* orthologs across diverse plant species from the  
341 EnsemblPlants and the protein databases Phytozome v12.1 retrieved 64 protein  
342 sequences from 37 genera (52 species, Table S4) including Bryophytes, eudicots,  
343 and monocots. The final alignment consisted of 670 positions. The obtained ML  
344 topology reflects the evolution of terrestrial plants with *Amborella trichopoda* at the  
345 basis of the two main clades, monocotyledons and eudicotyledons (Fig. 6). The  
346 protein is relatively well conserved as seen from the very small branches especially  
347 within the grass tribe Triticeae, including *Triticum aestivum* and *Hordeum vulgare*,  
348 which diverged about ten million years ago (Ma) (Bernhardt et al. 2017) or even the  
349 Poaceae, whose most recent common ancestor probably occurred 50-75 Ma  
350 (Bouchenak-Khelladi et al. 2010).

## 351 Discussion

352 *Exploiting significant, heritable genetic variation of TSN as well as a positive*  
353 *correlation with other traits can help to improve the grain yield in wheat*

354 Grain yield (GY) improvement is considered as the top focus of virtually every wheat  
355 breeding program. However, an extremely complex genetic nature of GY often

356 hampers its genetic improvement as it is the product of several yield components,  
357 e.g., the number of spikes per plant, grains per spike, thousand-grain weight. The  
358 number of grains per spike is a product of TSN and fertility. Therefore, an essential  
359 consideration in wheat breeding has been to employ a reductionist approach, i.e., to  
360 exploit the information about the individual component traits; most of which are  
361 negatively associated with each other. In this study, we analyzed a winter wheat  
362 panel comprising of 518 varieties for grain yield component traits such as TSN and  
363 SL along with the flowering time (FT). The grain yield data based on previous studies  
364 were taken for comparison purposes (Schulthess et al. 2017). In all observed traits,  
365 besides significant genetic variation, we observed a significant genotype-by-  
366 environment (year) interaction. Nevertheless, the broad-sense heritability estimates  
367 ranging from 0.68 to 0.89 suggested that genetic variation is heritable – an essential  
368 indicator of high selection response (Table 1). Similar heritability values for the  
369 studied traits have been reported recently in other diverse mapping populations (Guo  
370 et al. 2017; Würschum et al. 2018).

371 In addition to significant genetic variation, TSN showed a positive and  
372 significant correlation with SL, FT, and GY (Fig. 1e). This showed that albeit being  
373 weak (which is by virtue of the extreme quantitative genetic nature of GY), the  
374 correlation with GY could help improve the genetic gain. Moreover, it should be noted  
375 that the genetic architecture of yield component traits *per se* is also important which  
376 means that if the component traits possess complex genetic architecture, the  
377 problem of grain yield improvement would be further compounded. Nevertheless, a  
378 reasonably high heritability value suggests that TSN is strongly genetically inherited  
379 and that mapping of the underlying quantitative trait loci (QTL) would be efficient.

### 380 *High marker density governs the efficacy of genetic and physical mapping*

381 The efficiency of GWAS depends on the size of the population and genetic diversity.  
382 Genome-wide marker density with many polymorphic sites is therefore vital and  
383 coupled with a sharp decline in linkage disequilibrium (LD) between marker loci; it  
384 increases GWAS resolution. In our study, the size of the population, high-density  
385 genotyping, and the use of stringent linear mixed-effect models warranted the genetic  
386 mapping of true marker-trait-associations (MTA). As noted in another study based on  
387 a subset of varieties, the absence of distinct sub-populations in this panel suggests  
388 that the European winter wheat varieties have been bred, by and large, from a  
389 narrow genetic base and with similar goals (Muqaddasi et al. 2019) which is in line  
390 with other reports based on studies using similar genetic material but different marker  
391 platforms (Kollers et al. 2013; Würschum et al. 2013).

392 To identify the candidate genes, high marker density in a given QTL genetic  
393 region is necessary since it helps to narrow-down to the physical region harboring the  
394 gene underlying the trait. Moreover, since GWAS hinges on the principle that  
395 markers work as proxies to the genes/QTL underlying the traits, a high density of  
396 markers in the QTL genetic region becomes vital for the success of fine mapping. In  
397 this study, we exploited this premise to identify a candidate gene physically.

398 *Physical mapping shows that TaAPO-A1 is a likely candidate gene for TSN in wheat*

399 Our GWAS analysis revealed a significant QTL for total spikelet number on  
400 chromosome 7A, which explained ~25% of the total genotypic variance. Also,  
401 Würschum et al. (2018) recently reported a QTL for TSN on chromosome 7A in a  
402 similar type of elite winter wheat germplasm. Zhang et al. (2015) reported a putative  
403 *MOC1* ortholog to be associated with spikelet number, which is also located on  
404 chromosome 7A.

405 The strategy to investigate orthologous genes of rice with the known function  
406 was already successfully applied for various genes associated with grain size, grain  
407 weight as well as yield in wheat (Ma et al. 2016; Su et al. 2011; Wang et al. 2015;  
408 Zhang et al. 2012; Zhang et al. 2014; Zheng et al. 2014). The highly significant region  
409 of the detected TSN-QTL in our study corresponded to a physical interval of <1 Mb,  
410 containing a block of only ten genes, all in high LD (Figure 5). Based on the  
411 functional annotations, the rice gene *ABERRANT PANICLE ORGANIZATION 1*  
412 (*APO1*), an ortholog of *Arabidopsis UFO* (Ikeda et al. 2007; Ikeda et al. 2005;  
413 Samach et al. 1999; Wilkinson and Haughn 1995) was considered as the most likely  
414 candidate gene and was named as *TaAPO-A1* in wheat. The functional analyses in  
415 both rice and *Arabidopsis* revealed that the F-box containing protein is involved in the  
416 regulation and development of floral organs; more specifically *APO1* in rice that  
417 controls the number of spikelets per panicle by regulating the cell proliferation in  
418 meristems (Ikeda-Kawakatsu et al. 2009).

419 *Functional diversity among orthologs of TaAPO-A1 reveals the conserved F-box*  
420 *domain*

421 The availability of genomic data for several wheat varieties from *The 10+ Wheat*  
422 *Genome Project* allowed the investigation of the intraspecific diversity of *TaAPO-A1*  
423 gene among wheat varieties. The *TaAPO-A1* contains two exons, each containing a  
424 SNP which causes an amino acid substitution. In the first exon, a T/G polymorphism  
425 at base 140 was related to the exchange of phenylalanine to cysteine, and in the  
426 second exon, at base 1284, a G/A polymorphism mutated aspartic acid to asparagine  
427 (Figure S6). It was possible to develop a functional KASP marker for the SNP in the  
428 first exon and to screen the germplasm panel. Both alleles were present in almost  
429 identical frequencies with 49.38% of the varieties carrying the allele of Chinese  
430 Spring with nucleotide T (referred to as *TaAPO-A1b*) and 50.62% of the varieties  
431 carrying the G nucleotide (referred to as *TaAPO-A1a*). The Chinese Spring allele was  
432 strongly associated ( $P < 2.2e-16$ ) with an increase in TSN and moderately associated  
433 with an increase in spike length ( $P = 3.4e-04$ ) and yield ( $P = 9.0e-04$ ) (Figure 4). For  
434 the B- and D-genomes, the orthologs of *TaAPO-A1* were related to the genes  
435 *TraesCS7B01G384000* and *TraesCS7D01G468700*. However, no MTAs were  
436 discovered on these genomes. The identified *TaAPO-A1* variants reflect natural  
437 allelic diversity with mild phenotypic effects, which is beneficial for practical breeding.

438 The presence of *TaAPO-A1* orthologs in a wide range of plants including  
439 Bryophytes, monocotyledons, and eudicotyledons suggests a central role of this  
440 gene class in the evolution and development of terrestrial plants (Figures 6, S7). The  
441 *Arabidopsis* gene *UFO* and rice *APO1* (orthologs of *TaAPO-A1*) encode for an F-box  
442 containing protein. It has been shown that the rice *APO1* and *Arabidopsis UFO* are  
443 important for floral development in respective species (Ikeda et al. 2007; Samach et  
444 al. 1999). Molecularly, the proteins SKP1, cullin like and F-box containing  
445 polypeptides form SCF protein complexes to function as E3-ubiquitin ligases that  
446 target specific proteins for degradation (Kaiser et al. 1998; Patton et al. 1998). It has  
447 been shown that *Arabidopsis UFO* indirectly regulates the expression of class B floral  
448 homeotic gene *APETALA 3* by targeting the degradation of proteins which negatively  
449 regulate its transcription (Samach et al. 1999). The rice *apo1* mutants show a  
450 reduction in the number of primary branches and thereby the number of spikelets due  
451 to the precocious conversion of inflorescence meristem (IM) to spikelet meristem  
452 (SM). Such a mutant phenotype offers an indication that *APO1* might target proteins  
453 that promote the precocious conversion of IM to SM for degradation in a functional  
454 state. In line with this idea, the dominant gain of function *APO1* alleles with an  
455 elevated expression as well as overexpression transgenic lines of *APO1* showed  
456 prolonged inflorescence development resulting in more branch iterations and  
457 consequently more spikelets (Ikeda-Kawakatsu et al. 2009).

458 From our promoter analysis, we found an INDEL where the 115 bp insertion  
459 was always associated with high TSN haplotype, whereas the deletion with low TSN  
460 haplotype. From this finding, it may be inferred that winter wheat genotypes in the  
461 haplogroup with insertion polymorphism have slightly elevated expression of *TaAPO-*  
462 *A1* leading to prolonged maturation of inflorescence meristem eventually producing  
463 more spikelets per spike. Conversely, the deletion haplotype has a comparatively  
464 reduced expression level of *TaAPO-A1*, leading to less number of spikelets.  
465 Nevertheless, validation of the INDEL haplotype across the whole winter wheat panel  
466 as well as expression analysis of *TaAPO-A1* in the two haplogroups with high and  
467 low TSN may offer further insights into the regulation of TSN in wheat.

## 468 **Conclusions**

469 Our results demonstrate that with the availability of modern genomic tools such as  
470 the wheat reference sequence and the access to *The 10+ Wheat Genome Project*,  
471 the way from phenotype to a candidate gene is shortened considerably.  
472 Nevertheless, a robust genetic analysis including appropriate mapping populations,  
473 accurate and high-density genotyping, and proper phenotypic analyses are  
474 prerequisites to detecting significant QTL regions from which the causative genes  
475 could be deduced.

476

## 477 **References**

- 478 Allen AM, Winfield MO, BurrIDGE AJ, Downie RC, Benbow HR, Barker GL, Wilkinson  
479 PA, Coghill J, Waterfall C, Davassi A, Scopes G, Pirani A, Webster T, Brew F, Bloor  
480 C, Griffiths S, Bentley AR, Alda M, Jack P, Phillips AL, Edwards KJ (2017)  
481 Characterization of a Wheat Breeders' Array suitable for high-throughput SNP  
482 genotyping of global accessions of hexaploid bread wheat (*Triticum aestivum*). Plant  
483 Biotechnol J 15:390-401
- 484 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment  
485 search tool. Journal of Molecular Biology 215:403-410
- 486 Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and  
487 powerful approach to multiple testing. Journal of the Royal Statistical Society Series  
488 B (Methodological):289-300
- 489 Bernhardt N, Brassac J, Kilian B, Blattner FR (2017) Dated tribe-wide whole  
490 chloroplast genome phylogeny indicates recurrent hybridizations within Triticeae.  
491 BMC Evolutionary Biology 17:141
- 492 Boden SA, Cavanagh C, Cullis BR, Ramm K, Greenwood J, Finnegan EJ, Trevaskis  
493 B, Swain SM (2015) *Ppd-1* is a key regulator of inflorescence architecture and paired  
494 spikelet development in wheat. Nature Plants 1:14016
- 495 Bouchenak-Khelladi Y, Verboom GA, Savolainen V, Hodkinson TR (2010)  
496 Biogeography of the grasses (Poaceae): a phylogenetic approach to reveal  
497 evolutionary history in geographical space and geological time. Botanical Journal of  
498 the Linnean Society 162:543-557
- 499 Consortium IWGS (2018) Shifting the limits in wheat research and breeding using a  
500 fully annotated reference genome. Science 361:eaar7191
- 501 Darriba D, Taboada GL, Doallo R, Posada D (2011) ProtTest 3: fast selection of best-  
502 fit models of protein evolution. Bioinformatics 27:1164-1165
- 503 Debernardi JM, Lin H, Chuck G, Faris JD, Dubcovsky J (2017) microRNA172 plays a  
504 crucial role in wheat spike morphogenesis and grain threshability. Development  
505 144:1966-1975
- 506 Deng Z, Cui Y, Han Q, Fang W, Li J, Tian J (2017) Discovery of consistent QTLs of  
507 wheat spike-related traits under nitrogen treatment at different development stages.  
508 Front Plant Sci 8:2120
- 509 Faris JD, Fellers JP, Brooks SA, Gill BS (2003) A bacterial artificial chromosome  
510 contig spanning the major domestication locus Q in wheat and identification of a  
511 candidate gene. Genetics 164:311-321
- 512 Frichot E, François O (2015) LEA: An R package for landscape and ecological  
513 association studies. Methods in Ecology and Evolution 6:925-929
- 514 Gauley A, Boden SA (2019) Genetic pathways controlling inflorescence architecture  
515 and development in wheat and barley. Journal of Integrative Plant Biology 61:296-  
516 309



- 517 Greenwood JR, Finnegan EJ, Watanabe N, Trevaskis B, Swain SM (2017) New  
518 alleles of the wheat domestication gene *Q* reveal multiple roles in growth and  
519 reproductive development. *Development* 144:1959-1965
- 520 Guindon S, Gascuel O (2003) A simple, fast, and accurate algorithm to estimate  
521 large phylogenies by maximum likelihood. *Systematic Biology* 52:696-704
- 522 Guo Z, Chen D, Alqudah AM, Röder MS, Ganai MW, Schnurbusch T (2017)  
523 Genome-wide association analyses of 54 traits identified multiple loci for the  
524 determination of floret fertility in wheat. *New Phytologist* 214:257-270
- 525 Guo Z, Zhao Y, Röder MS, Reif JC, Ganai MW, Chen D, Schnurbusch T (2018)  
526 Manipulation and prediction of spike morphology traits for the improvement of grain  
527 yield in wheat. *Sci Rep-Uk* 8:14435
- 528 Ikeda-Kawakatsu K, Yasuno N, Oikawa T, Iida S, Nagato Y, Maekawa M, Kyojuka J  
529 (2009) Expression level of *ABERRANT PANICLE ORGANIZATION1* determines rice  
530 inflorescence form through control of cell proliferation in the meristem. *Plant Physiol*  
531 150:736-747
- 532 Ikeda K, Ito M, Nagasawa N, Kyojuka J, Nagato Y (2007) Rice *ABERRANT*  
533 *PANICLE ORGANIZATION 1*, encoding an F-box protein, regulates meristem fate.  
534 *The Plant Journal* 51:1030-1040
- 535 Ikeda K, Nagasawa N, Nagato Y (2005) *ABERRANT PANICLE ORGANIZATION 1*  
536 temporally regulates meristem identity in rice. *Developmental Biology* 282:349-360
- 537 Jones DT, Taylor WR, Thornton JM (1992) The rapid generation of mutation data  
538 matrices from protein sequences. *Bioinformatics* 8:275-282
- 539 Kaiser P, Sia RA, Bardes EG, Lew DJ, Reed SI (1998) Cdc34 and the F-box protein  
540 Met30 are required for degradation of the Cdk-inhibitory kinase Swe1. *Genes &*  
541 *Development* 12:2587-2597
- 542 Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S,  
543 Cooper A, Markowitz S, Duran C (2012) Geneious Basic: an integrated and  
544 extendable desktop software platform for the organization and analysis of sequence  
545 data. *Bioinformatics* 28:1647-1649
- 546 Kollers S, Rodemann B, Ling J, Korzun V, Ebmeyer E, Argillier O, Hinze M, Plieske J,  
547 Kulosa D, Ganai MW, Röder MS (2013) Genetic architecture of resistance to  
548 *Septoria tritici blotch* (*Mycosphaerella graminicola*) in European winter wheat.  
549 *Molecular Breeding* 32:411-423
- 550 Koppolu R, Schnurbusch T (2019) Developmental pathways for shaping spike  
551 inflorescence architecture in barley and wheat. *Journal of Integrative Plant Biology*  
552 61:278-295
- 553 Li C, Lin H, Chen A, Lau M, Jernstedt J, Dubcovsky J (2019) Wheat *VRN1* and *FUL2*  
554 play critical and redundant roles in spikelet meristem identity and spike determinacy.  
555 bioRxiv:510388

- 556 Liu J, Xu Z, Fan X, Zhou Q, Cao J, Wang F, Ji G, Yang L, Feng B, Wang T (2018) A  
557 Genome-Wide Association Study of Wheat Spike Related Traits in China. *Front Plant*  
558 *Sci* 9
- 559 Ma L, Li T, Hao C, Wang Y, Chen X, Zhang X (2016) *TaGS5-3A*, a grain size gene  
560 selected during wheat improvement for larger kernel and yield. *Plant Biotechnol J*  
561 14:1269-1280
- 562 Muqaddasi QH, Zhao Y, Rodemann B, Plieske J, Ganal MW, Röder MS (2019)  
563 Genome-wide Association Mapping and Prediction of Adult Stage *Septoria tritici*  
564 Blotch Infection in European Winter Wheat via High-Density Marker Arrays. *The Plant*  
565 *Genome* 12:180029
- 566 Ochagavía H, Prieto P, Savin R, Griffiths S, Slafer G (2018) Dynamics of leaf and  
567 spikelet primordia initiation in wheat as affected by *Ppd-1a* alleles under field  
568 conditions. *J Exp Bot* 69:2621-2631
- 569 Patton EE, Willems AR, Tyers M (1998) Combinatorial control in ubiquitin-dependent  
570 proteolysis: don't Skp the F-box hypothesis. *Trends in Genetics* 14:236-243
- 571 Sakuma S, Golan G, Guo Z, Ogawa T, Tagiri A, Sugimoto K, Bernhardt N, Brassac J,  
572 Mascher M, Hensel G, Ohnishi S, Jinno H, Yamashita Y, Ayalon I, Peleg Z,  
573 Schnurbusch T, Komatsuda T (2019) Unleashing floret fertility in wheat through the  
574 mutation of a homeobox gene. *Proceedings of the National Academy of Sciences*  
575 116:5182-5187
- 576 Samach A, Klenz JE, Kohalmi SE, Risseuw E, Haughn GW, Crosby WL (1999) The  
577 *UNUSUAL FLORAL ORGANS* gene of *Arabidopsis thaliana* is an F-box protein  
578 required for normal patterning and growth in the floral meristem. *The Plant Journal*  
579 20:433-445
- 580 Schulthess AW, Reif JC, Ling J, Plieske J, Kollers S, Ebmeyer E, Korzun V, Argillier  
581 O, Stiewe G, Ganal MW, Röder MS, Jiang Y (2017) The roles of pleiotropy and close  
582 linkage as revealed by association mapping of yield and correlated traits of wheat  
583 (*Triticum aestivum* L.). *J Exp Bot* 68:4089-4101
- 584 Shaw LM, Lyu B, Turner R, Li C, Chen F, Han X, Fu D, Dubcovsky J (2018)  
585 *FLOWERING LOCUS T2* regulates spike development and fertility in temperate  
586 cereals. *J Exp Bot* 70:193-204
- 587 Sim N-L, Kumar P, Hu J, Henikoff S, Schneider G, Ng PC (2012) SIFT web server:  
588 predicting effects of amino acid substitutions on proteins. *Nucleic acids research*  
589 40:W452-W457
- 590 Sorrells ME, Gustafson JP, Somers D, Chao S, Benscher D, Guedira-Brown G,  
591 Huttner E, Kilian A, McGuire PE, Ross K (2011) Reconstruction of the Synthetic  
592 W7984 x Opata M85 wheat reference population. *Genome* 54:875-882
- 593 Stamatakis A (2014) RAxML version 8: a tool for phylogenetic analysis and post-  
594 analysis of large phylogenies. *Bioinformatics* 30:1312-1313

- 595 Su Z, Hao C, Wang L, Dong Y, Zhang X (2011) Identification and development of a  
596 functional marker of *TaGW2* associated with grain weight in bread wheat (*Triticum*  
597 *aestivum* L.). *Theor Appl Genet* 122:211-223
- 598 Utz HF, Melchinger AE, Schön CC (2000) Bias and sampling error of the estimated  
599 proportion of genotypic variance explained by quantitative trait loci determined from  
600 experimental data in maize using cross validation and validation with independent  
601 samples. *Genetics* 154:1839-1849
- 602 VanRaden PM (2008) Efficient methods to compute genomic predictions. *Journal of*  
603 *Dairy Science* 91:4414-4423
- 604 Wang S, Zhang X, Chen F, Cui D (2015) A single-nucleotide polymorphism of *TaGS5*  
605 gene revealed its association with kernel weight in Chinese bread wheat. *Front Plant*  
606 *Sci* 6:1166
- 607 Wang SC, Wong DB, Forrest K, Allen A, Chao SM, Huang BE, Maccaferri M, Salvi S,  
608 Milner SG, Cattivelli L, Mastrangelo AM, Whan A, Stephen S, Barker G, Wieseke R,  
609 Plieske J, International Wheat Genome Sequencing Consortium, Lillemo M, Mather  
610 D, Appels R, Dolferus R, Brown-Guedira G, Korol A, Akhunova AR, Feuillet C, Salse  
611 J, Morgante M, Pozniak C, Luo MC, Dvorak J, Morell M, Dubcovsky J, Ganal M,  
612 Tuberosa R, Lawley C, Mikoulitch I, Cavanagh C, Edwards KJ, Hayden M, Akhunov  
613 E (2014) Characterization of polyploid wheat genomic diversity using a high-density  
614 90 000 single nucleotide polymorphism array. *Plant Biotechnol J* 12:787-796
- 615 Wilkinson MD, Haughn GW (1995) *UNUSUAL FLORAL ORGANS* controls meristem  
616 identity and organ primordia fate in Arabidopsis. *The Plant Cell* 7:1485-1499
- 617 Würschum T, Langer SM, Longin CFH, Korzun V, Akhunov E, Ebmeyer E,  
618 Schachschneider R, Schacht J, Kazman E, Reif JC (2013) Population structure,  
619 genetic diversity and linkage disequilibrium in elite winter wheat assessed with SNP  
620 and SSR markers. *Theor Appl Genet* 126:1477-1486
- 621 Würschum T, Leiser WL, Langer SM, Tucker MR, Longin CFH (2018) Phenotypic  
622 and genetic analysis of spike and kernel characteristics in wheat reveals long-term  
623 genetic trends of grain yield components. *Theor Appl Genet* 131:2071-2084
- 624 Yu JM, Pressoir G, Briggs WH, Bi IV, Yamasaki M, Doebley JF, McMullen MD, Gaut  
625 BS, Nielsen DM, Holland JB, Kresovich S, Buckler ES (2006) A unified mixed-model  
626 method for association mapping that accounts for multiple levels of relatedness. *Nat*  
627 *Genet* 38:203-208
- 628 Zhai H, Feng Z, Li J, Liu X, Xiao S, Ni Z, Sun Q (2016) QTL analysis of spike  
629 morphological traits and plant height in winter wheat (*Triticum aestivum* L.) using a  
630 high-density SNP and SSR-based linkage map. *Front Plant Sci* 7:1617
- 631 Zhang B, Liu X, Xu W, Chang J, Li A, Mao X, Zhang X, Jing R (2015) Novel function  
632 of a putative *MOC1* ortholog associated with spikelet number per spike in common  
633 wheat. *Sci Rep-Uk* 5:12211

634 Zhang L, Zhao YL, Gao LF, Zhao GY, Zhou RH, Zhang BS, Jia JZ (2012) *TaCKX6-*  
635 *D1*, the ortholog of rice *OsCKX2*, is associated with grain weight in hexaploid wheat.  
636 *New Phytologist* 195:574-584

637 Zhang Y, Liu J, Xia X, He Z (2014) *TaGS-D1*, an ortholog of rice *OsGS3*, is  
638 associated with grain weight and grain length in common wheat. *Molecular breeding*  
639 34:1097-1107

640 Zheng J, Liu H, Wang Y, Wang L, Chang X, Jing R, Hao C, Zhang X (2014) *TEF-7A*,  
641 a transcript elongation factor gene, influences yield-related traits in bread wheat  
642 (*Triticum aestivum* L.). *J Exp Bot* 65:5351-5365

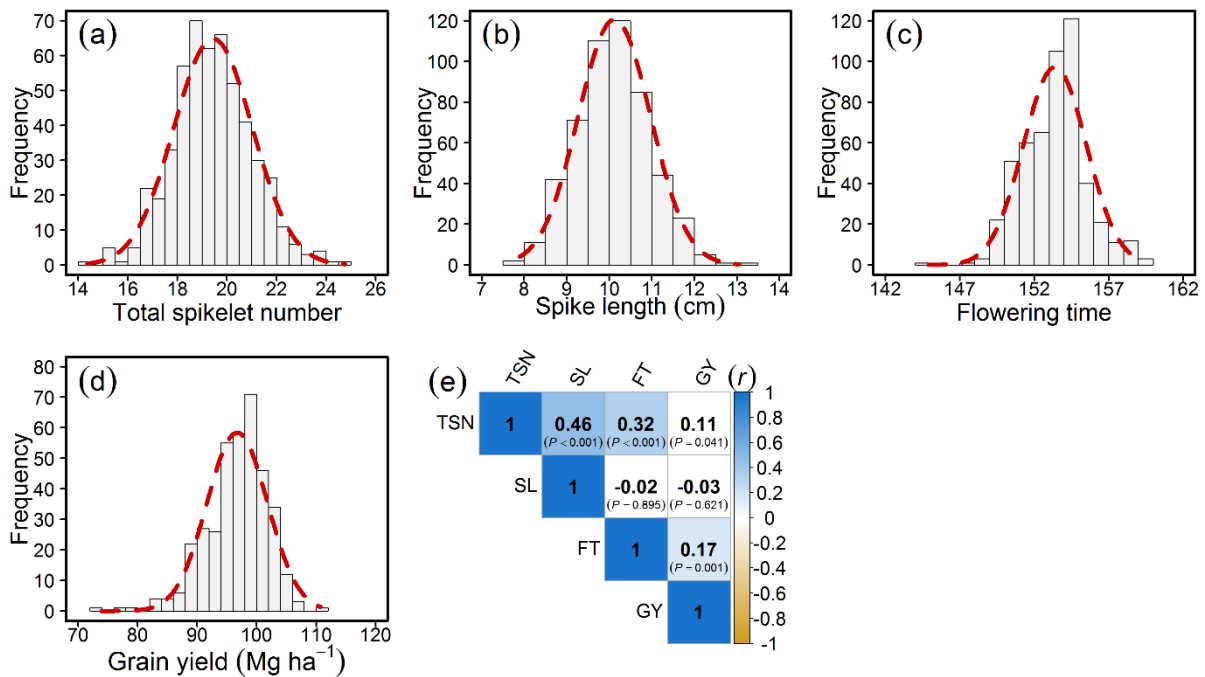
643 **Acknowledgments:** The genotyping data were produced in the project VALID  
644 funded by the German Federal Ministry of Education and Research (BMBF; project  
645 number 0315947). We are grateful to Ellen Weiß, Anette Heber, Ute Ostermann, and  
646 Sonja Allner for help in phenotypic data collection. We are thankful to *The 10+ Wheat*  
647 *Genome Project* for making the resources available before publication.

648 **Author contribution statement:** QHM and MSR conceived the idea. QHM analyzed  
649 the data, interpreted the results, and wrote the manuscript. JB and RK contributed to  
650 sequence and phylogenetic analyses. JP and MWG contributed the genotypic data.  
651 RK and MSR contributed to the interpretation of results and writing of the manuscript.

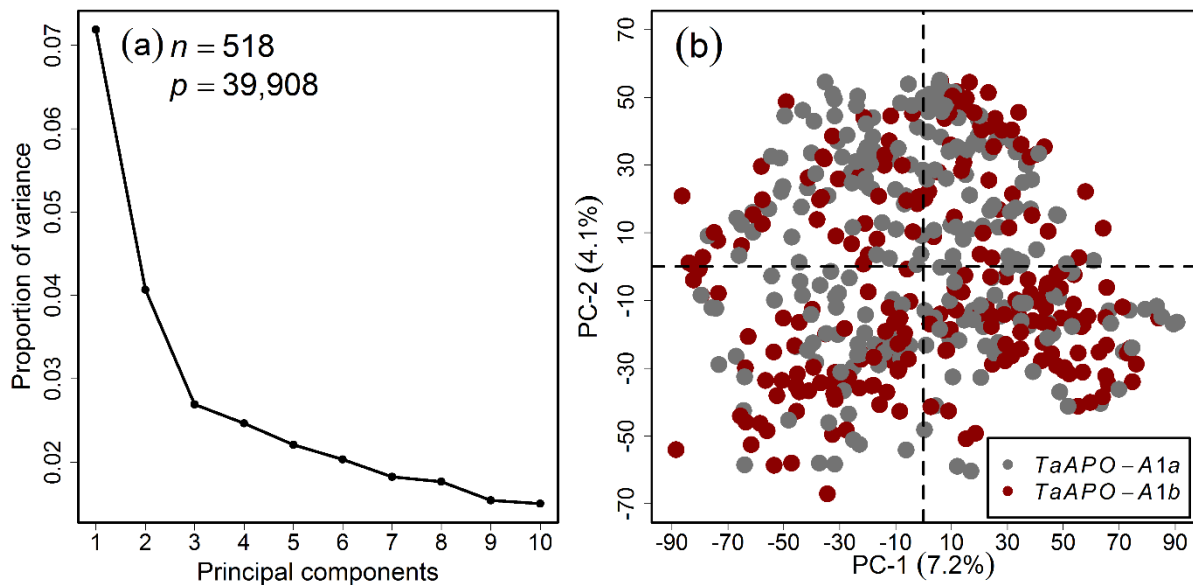
652 **Conflict of interest:** On behalf of all authors, the corresponding author states that  
653 there is no conflict of interest. JP and MWG are members of the company  
654 TraitGenetics. This does, however, in no way limit the availability or sharing of data  
655 and materials.

656

657 **Figures**

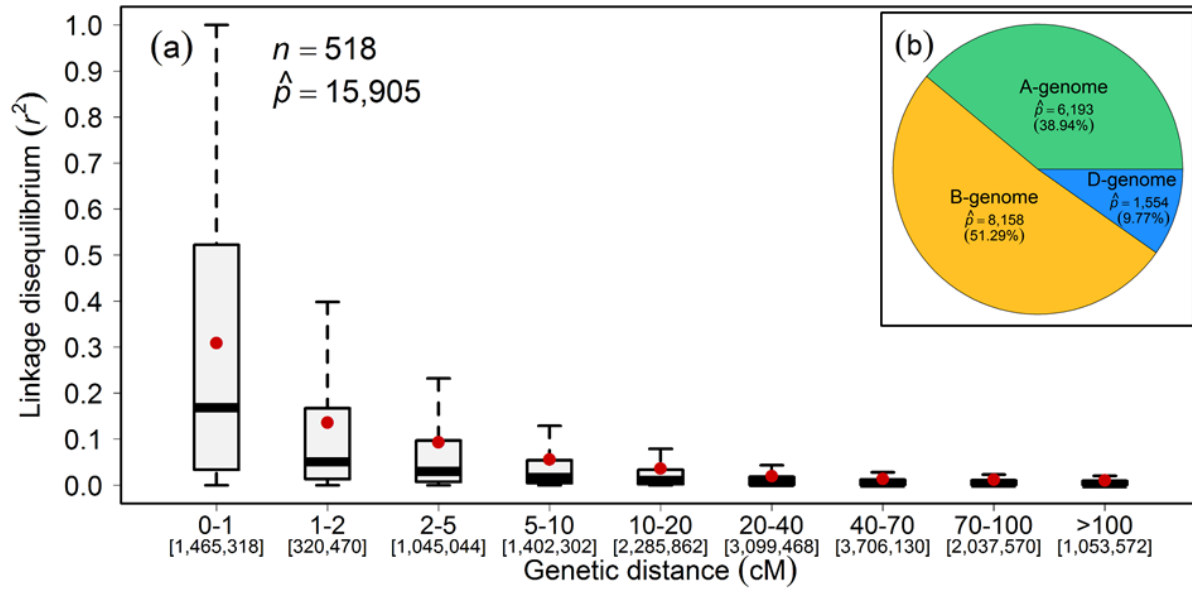


658 **Figure 1: Distribution and correlation of the investigated traits in a panel of 518**  
 659 **elite European winter wheat varieties. Distribution of (a) Total spikelet number**  
 660 **(TSN), (b) Spike length (SL), (c) Flowering time (FT), and (d) Grain yield (GY); (e)**  
 661 **Pearson product moment correlation (*r*) among the investigated traits. *P*-value**  
 662 **denotes the significance of the respective correlation.**  
 663

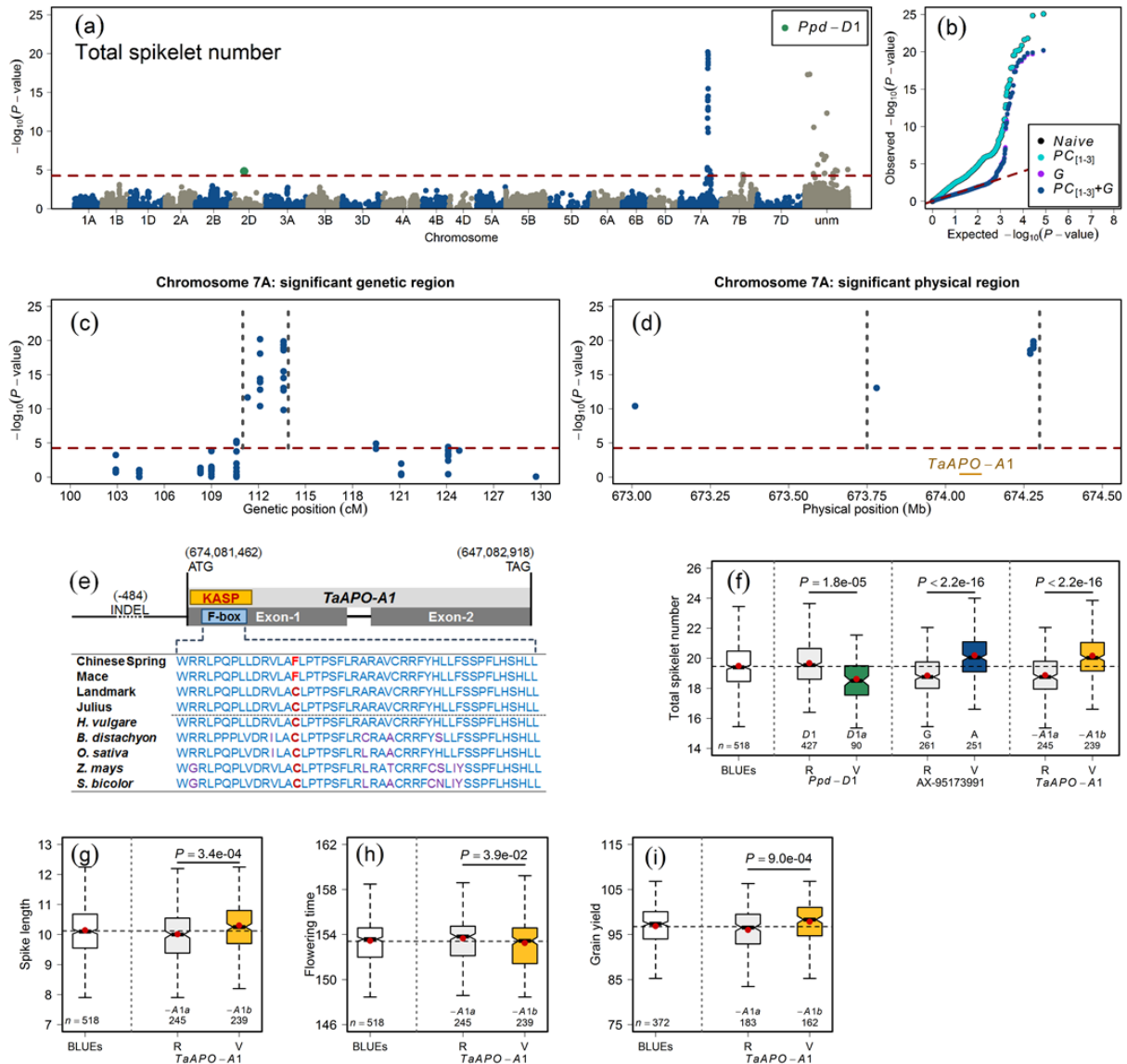


664 **Figure 2: Principal component (PC) analysis on the wheat marker loci**  
 665 **combined from the 35k and 90k single nucleotide polymorphism arrays. (a)**  
 666 **Scree plot showing the first ten PCs and their corresponding proportion of**  
 667 **variance, (b) Scatterplot showing the absence of pronounced clustering among**  
 668 **the varieties. Different colors represent the *TaAPO-A1* alleles. *n* and *p* denote**  
 669 **the number of varieties and the marker genotypes used in the analysis, respectively.**  
 670





671  
672 **Figure 3: Genome-wide decay of linkage disequilibrium (LD;  $r^2$ ) as a function of**  
673 **genetic map distance (cM) between the marker loci in the population of**  
674 **European winter wheat varieties. (a) Boxplots represent the LD-decay, (b) Sub-**  
675 **genome-wise distribution of mapped marker loci. Red dots within the boxplots**  
676 **represent the mean.  $n$  and  $\hat{p}$  denote the number of varieties and mapped marker loci,**  
677 **respectively.**

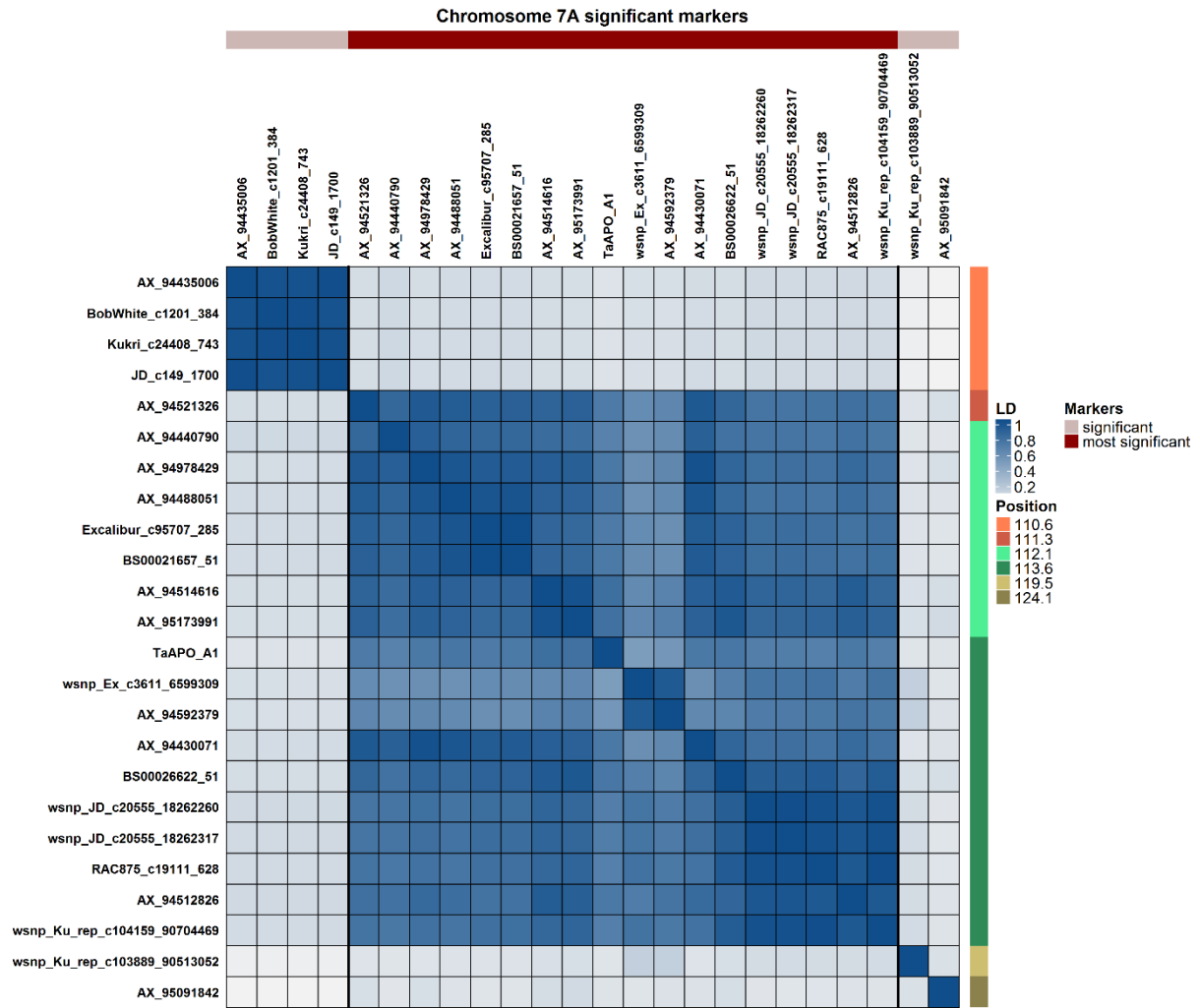


678

679 **Figure 4: Summary of genome-wide association studies of total spikelet**  
 680 **number per spike in the population of 518 European winter wheat varieties. (a)**  
 681 **Manhattan plot shows the distribution of marker significance  $-\log_{10}(P\text{-value})$**   
 682 **along the chromosomes. The correction for population stratification and familial**  
 683 **relatedness was performed by using the first three principal components ( $PC_{[1-3]}$ ) and**  
 684 **an additive genomic relationship matrix (G) in a linear mixed-effect model. The red**  
 685 **dashed line marks the multiple testing criteria of false discovery rate (FDR) <0.05, (b)**  
 686 **Quantile-quantile plot showing the distribution of observed versus expected**  
 687 **(red dashed line)  $-\log_{10}(P\text{-value})$ . The general linear model (naive) without**  
 688 **correction for population structure, the  $PC_{[1-3]}$  model (population structure corrected**  
 689 **with the first three PCs), the G model (familial relatedness corrected with a genomic**  
 690 **relationship matrix), and the  $PC_{[1-3]}+G$  model (population structure and familial**  
 691 **relatedness corrected with PCs and the G matrix). The color code for different**  
 692 **models is given in the figure legend, (c) Significant genetic region on**  
 693 **chromosome 7A for TSN in wheat. The gray vertical dashed lines mark the highly**  
 694 **significant physical region on chromosome 7A for**

695 **TSN in wheat.** The gray vertical dashed lines mark the highly significant physical  
696 region, **(e) Gene structure of *TaAPO-A1*.** The orange box represents the location of  
697 the KASP marker developed to exploit the variation in the F-box domain (highlighted  
698 in blue color). The horizontal line before the first exon depicts promotor region  
699 harboring INDEL and corresponding position. The first four rows represent the F-box  
700 sequences of wheat varieties (courtesy: *The 10+ Wheat Genomes Project*) and the  
701 second four rows represent the F-box domain of closely related species viz.,  
702 *Hordeum vulgare*, *Brachypodium distacyon*, *Oryza sativa*, *Zea mays*, and *Sorghum*  
703 *bicolor*. The non-synonymous mutation is highlighted in red color. The location of  
704 start and stop codons on chromosome 7A are given in the figure, **(f) Allele-wise**  
705 **phenotypic distribution of the most significant markers and KASP marker for**  
706 ***TaAPO-A1* associated with (f) TSN, (g) Spike length, (h) Flowering time, and (i)**  
707 **Grain yield.** *P* denotes the significance value of the two-sided t-test used to compare  
708 the mean value of marker alleles. In sub-figures (f) to (i), the first boxplots represent  
709 the phenotypic distribution of the best linear unbiased estimations (BLUEs) for the  
710 respective trait, whereas R and V denote the reference (major) and variant (minor)  
711 allele in the investigated population, respectively.

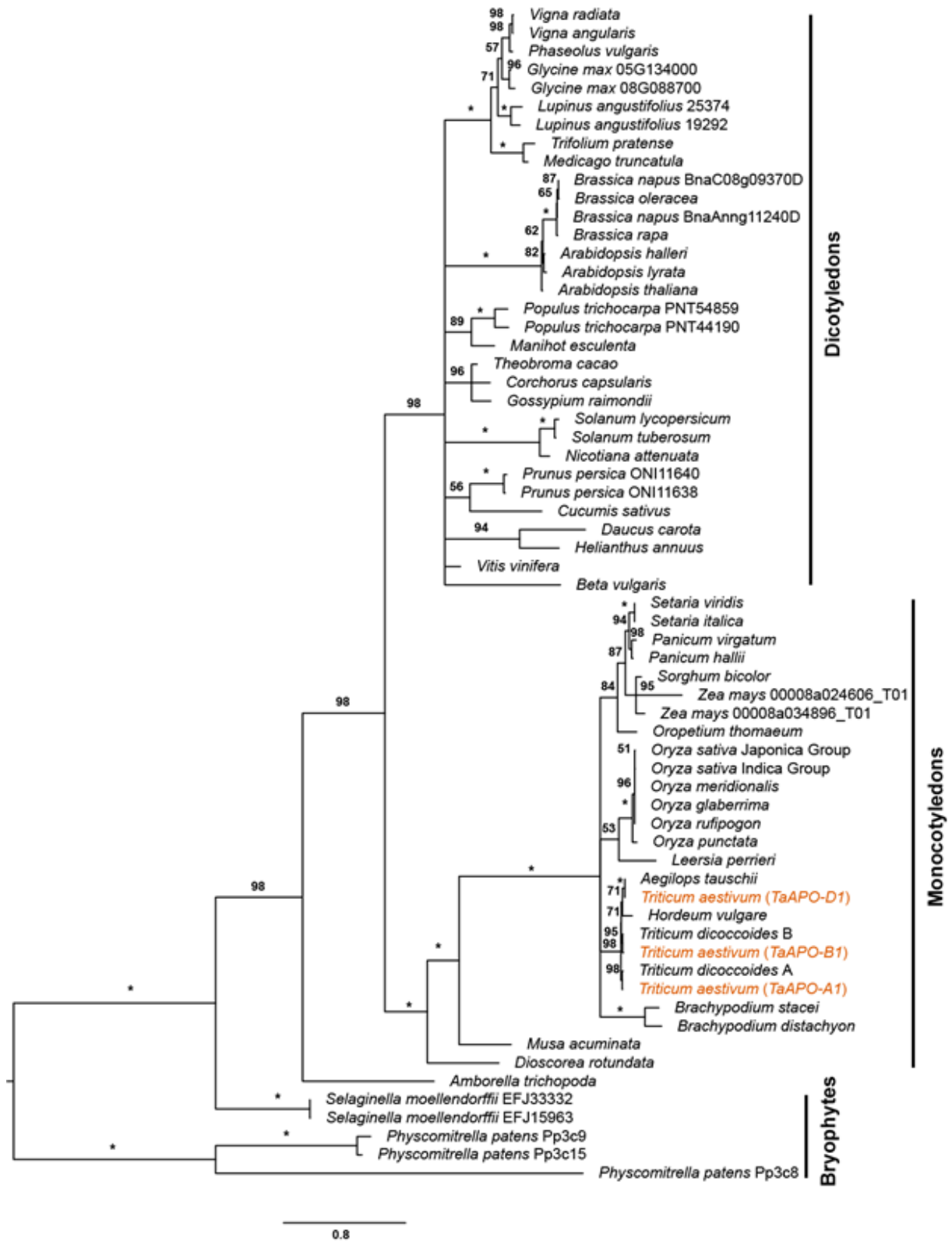
712



713

714 **Figure 5: Pairwise linkage disequilibrium ( $r^2$ ) among the marker loci (including**  
 715 **the KASP marker for *TaAPO-A1*) present in significant genetic region of TSN on**  
 716 **chromosome 7A in wheat. Based on the linkage blocks, markers are divided into**  
 717 **two categories viz., significant, and most significant. The color key is given in the**  
 718 **figure.**

719



720

721 **Figure 6: Maximum likelihood phylogenetic tree of *TaAPO-A1* orthologous**  
 722 **proteins across terrestrial plant species.** Bootstrap values are indicated along the  
 723 branches. Asterisks indicate >99% bootstrap values. The *TaAPO* homoeologs are  
 724 highlighted in orange color. The bars on the right side indicate the major clades. The  
 725 amino acid substitution scale is indicated at the bottom of the figure.



726 **Table 1. Summary statistics of the investigated traits, namely total spikelet**  
727 **number (TSN), spike length (SL), flowering time (FT), and grain yield (GY).**

Parameter	TSN	SL	FT	GY
Minimum	14.38	7.90	144.96	73.94
Mean	19.45	10.12	153.38	96.74
Maximum	24.75	13.05	159.61	110.71
$\sigma_G^2$	1.71 <sup>a</sup>	0.50 <sup>a</sup>	3.42 <sup>a</sup>	22.89 <sup>a</sup>
$\sigma_E^2$	1.75 <sup>a</sup>	1.63 <sup>a</sup>	6.30 <sup>a</sup>	94.51 <sup>a</sup>
$\sigma_e^2$	1.60	0.44	1.90	23.74
$H^2$	0.68	0.70	0.84	0.89
$nE$	2	2	3	8

728  $\sigma_G^2$  = genotypic variance;  $\sigma_E^2$  = environmental variance;  $\sigma_e^2$  = residual variance;  $H^2$  =  
729 broad-sense heritability;  $nE$  = number of environments; a = significant at <0.001  
730 probability level.