

Transcription factor retention through multiple polyploidisation steps in wheat

Catherine EB Evans^{1,2}, Ramesh Arunkumar¹, Philippa Borrill¹

¹ Department of Crop Genetics, John Innes Centre, Norwich Research Park, NR4 7UH, UK

² School of Biosciences, University of Birmingham, Birmingham, B15 2TT, UK

Author for correspondence:

Philippa Borrill, Tel: +44 (0)1603 450533, Email: philippa.borrill@jic.ac.uk

Total word count for the main body of the text (excluding summary, references and legends): 5,883 words

Summary: 199 words

Introduction: 874 words

Materials and Methods: 1,699 words

Results: 1,981 words

Discussion: 1,324 words

Colour figures: 4

Supplemental Figures: 8

Supplemental Tables: 4

1 **Summary**

- 2 • Whole genome duplication (WGD) is widespread in plant evolutionary history, but the
3 mechanisms of non-random gene loss after WGD are debated. The gene balance hypothesis
4 proposes that dosage-sensitive genes such as regulatory genes are retained in polyploids. To
5 test this hypothesis, we analysed the retention of transcription factors (TFs) in the recent
6 allohexaploid bread wheat (*Triticum aestivum*).
- 7 • We annotated TFs in hexaploid, tetraploid and diploid wheats; compared the co-expression of
8 homoeologous TF and non-TF triads; and analysed single nucleotide variation in TFs across
9 cultivars.
- 10 • We found that, following each of two hybridisation and WGD events, the proportion of TFs
11 in the genome increased. TFs were preferentially retained over other genes as homoeologous
12 groups in tetraploid and hexaploid wheat. Across cultivars, TF triads contain fewer
13 deleterious missense mutations than non-TFs.
- 14 • TFs are preferentially retained as three functional homoeologs in hexaploid wheat, in support
15 of the gene balance hypothesis. High co-expression between TF homoeologs suggests that
16 neo- and sub-functionalisation are not major drivers of TF retention in this young polyploid.
17 Knocking out one TF homoeolog to alter gene dosage, using TILLING or CRISPR, could be
18 a way to further test the gene balance hypothesis and generate new phenotypes for wheat
19 breeding.

21 **Keywords**

22 Transcription factors, evolution, polyploidy, *Triticum aestivum* L. (wheat), gene balance hypothesis

24 **Introduction**

25 Gene duplication plays a major role in the evolution of genetic and phenotypic diversity and in
26 speciation events in eukaryotes (Lynch & Conery, 2000; Van de Peer *et al.*, 2009b). Ancient whole
27 genome duplications (WGD) are observed throughout the angiosperm plant phylogeny and occurred
28 at the base of major clades such as the seed plants, core eudicots and monocots (Van de Peer *et al.*,
29 2009a; Tasdighian *et al.*, 2017). Following WGD, most gene duplicates are eventually lost from the
30 genome by pseudogenisation or deletion over the course of millions of years (Freeling, 2009),
31 although the retention of a significant portion of duplicates has been observed across multiple plant
32 lineages (Lloyd *et al.*, 2014; Li *et al.*, 2016).

33 Several explanations for the retention of gene duplicates have been proposed including selection for
34 genetic redundancy (or genetic buffering) where effects of a null mutation are compensated by the

35 presence of an intact duplicate copy (Nowak *et al.*, 1997) and gene dosage increase where increased
36 dosage is advantageous, for example by increasing flux through a pathway (Ohno, 1970). Other
37 mechanisms include sub-functionalisation of the duplicate copies which may be mediated by
38 complementary degenerative mutations in each copy (Force *et al.*, 1999) and paralog interference
39 where degenerative mutations in one duplicate copy interfere with the function of the other copy
40 thereby promoting the retention of both functional copies (Baker *et al.*, 2013). Duplicate gene copies
41 may also be retained if they evolve new and distinct functions through neo-functionalisation (Ohno,
42 1970).

43 An alternative model for duplicate retention is the gene balance hypothesis, which proposes that
44 dosage sensitive genes tend to be retained as duplicates (Birchler *et al.*, 2005; Birchler & Veitia,
45 2007). This hypothesis explains the observation that the loss of genes after WGD is non-random and
46 certain classes of gene are preferentially retained including genes involved in regulatory interactions
47 or in protein complexes which are dosage sensitive (Blanc & Wolfe, 2004). Conversely, these dosage
48 sensitive genes are less frequently found in segmental duplications in which they would upset the
49 dosage balance with interacting partners (Maere *et al.*, 2005), in contrast to WGD where their
50 interacting partners would also be duplicated.

51 Studies across multiple angiosperms have revealed that transcription factors (TFs), a major type of
52 dosage sensitive regulatory gene, tend to be retained as duplicates after WGD for millions of years
53 (Lloyd *et al.*, 2014; Li *et al.*, 2016). In a comparative study of 37 sequenced angiosperm genomes, Li
54 *et al.* (2016) found that duplicate genes that originated at the Cretaceous-Paleogene boundary ~50 to
55 70 million years ago (mya) when a large number of WGD events occurred (Van de Peer *et al.*, 2009a),
56 were enriched for TFs. However these angiosperm-wide studies focussed on relatively old WGD
57 events >5mya, whilst more recent WGD events which have occurred in individual lineages and are
58 found in several major crop species are less well studied, perhaps due to a lack of genome sequences.
59 Preferential retention of dosage sensitive genes such as TFs has been observed in the young polyploid
60 *Tragopogon miscellus* which underwent WGD only ~80 years before (Buggs *et al.*, 2012). This study
61 used a limited number of loci, therefore there remains a need to understand the effects of recent (<5
62 mya) WGD at a genome-wide scale. The recent publication of the wheat genome sequence (IWGSC
63 *et al.*, 2018) provides an opportunity to examine the retention of dosage-sensitive genes from more
64 recent WGDs.

65 Hexaploid bread wheat evolved from two hybridisation and WGD events: allotetraploid wild emmer
66 wheat (*Triticum turgidum* ssp. *dicoccoides*) was formed approximately 0.4 mya when the A genome
67 progenitor *Triticum urartu* hybridised with the B genome progenitor species (Feldman & Levy, 2012).
68 The allotetraploid emmer was domesticated and hybridised with the D genome progenitor *Aegilops*
69 *tauschii* approximately 10,000 years ago to form hexaploid bread wheat (*Triticum aestivum* L.)

70 (Dubcovsky & Dvorak, 2007). This two-step recent history of WGD events has resulted in >50% of
71 genes being present with three homoeologous copies in bread wheat (IWGSC *et al.*, 2018). Previous
72 studies in wheat have shown that 58% of NAC TFs and 63% of MIKC-type MADS-box TFs have
73 three homoeologs (Borrill *et al.*, 2017; Schilling *et al.*, 2020), but a systematic study has not been
74 carried out to establish whether the preferential retention of TFs is observed across all TF families in
75 this recent polyploidy.

76 In this study we investigated whether these two recent WGD events resulted in the preferential
77 retention of TFs in hexaploid wheat, as would be predicted by the gene balance hypothesis. Using the
78 extensive curated expression data available for wheat, we explored alternative hypotheses about TF
79 retention, such as sub- or neo-functionalisation, based on expression patterns. Moreover, since genetic
80 variation in several TFs has been instrumental in wheat adaptation during domestication including the
81 free-threshing gene *Q* (Simons *et al.*, 2006) and the vernalisation gene *VRN1* (Yan *et al.*, 2003), we
82 examined the natural variation in TF homoeologs observed in wheat. Specifically we examined the
83 propensity of wheat TFs to be retained as functional copies without deleterious mutations at a
84 population level. Hence our study addresses not only an evolutionary question about the retention of
85 TFs in young polyploids, but also provides insight into TF expression diversity and genetic variation
86 which lays a foundation for future research and breeding.

87

88 **Materials and Methods**

89 **Annotation of TFs in wheat and progenitor species**

90 Peptide sequences for genes in the RefSeqv1.1 gene annotation of *Triticum aestivum* cv. Chinese
91 Spring (IWGSC *et al.*, 2018) were downloaded from EnsemblPlants (Howe *et al.*, 2020). The file was
92 divided into three parts to contain <50,000 sequences per file and TFs were annotated in each file
93 using iTAK online v1.6 (Zheng *et al.*, 2016). In cases where different transcript isoforms were
94 assigned to different TF families (23 out of 6,128 genes) the family assigned to the longer transcript
95 isoform was retained (Table S1). Peptide sequences for genes in the *Aegilops tauschii* assembly (Luo
96 *et al.*, 2017) were downloaded from EnsemblPlants, divided into six smaller files and annotated using
97 iTAK online v1.6. Again, when different transcript isoforms were assigned to different TF families
98 (186 out of 2,120 genes) the family assigned to the longer transcript isoform was retained (Table S2).
99 In general, discrepancies between TF families were due to one isoform being truncated, with the
100 truncated isoform lacking a protein domain that allowed a more specific TF family to be assigned to
101 the longer isoform. Coding sequences for the longest isoforms of genes in the *Triticum urartu* genome
102 (Ling *et al.*, 2018) were downloaded from <http://www.mbkbase.org/Tu/> and annotated using iTAK
103 online v1.6 (Table S3). TF annotations for *Triticum turgidum* ssp. *dicoccoides* cv. Zavitan (Avni *et*
104 *al.*, 2017) were downloaded from the iTAK database (update 18.12) (Zheng *et al.*, 2016) (Table S4).

105 **Identification of 1:1:1 triads in hexaploid and 1:1 diads in tetraploid wheat**

106 Homoeologs were downloaded from EnsemblPlants Biomart for the RefSeqv1.1 gene annotation
107 using only high confidence gene models. Only one2one homoeologs (assigned by EnsemblPlants)
108 were retained. There were 20,393 triads corresponding to 61,179 genes (56.7% of genes) (Table S1).
109 Homoeologs in *T. turgidum* ssp. *dicoccoides* were obtained from Avni *et al.* (2017) and filtered to
110 only retain 1:1 homoeologs by removing “singleton” and “hit2homolog” (i.e. paralog) groups (Table
111 S4). Only high confidence genes from the RefSeqv1.1 annotation were used in all subsequent
112 analyses.

113 **Adjusting for the effect of gene loss in tetraploid wheat on hexaploid wheat triad** 114 **numbers per TF family**

115 In order to adjust for the differences in triad proportions between TF families observed in hexaploid
116 due to the varying proportions in diads in tetraploid wheat, we calculated the normalised percentage
117 of genes in triads:

$$118 \quad \text{Normalised percentage of genes in triads} = \frac{\% \text{ of genes in triads in hexaploid wheat}}{\% \text{ of genes in diads in that TF family}}$$

119 For example if 60% of genes were in triads in hexaploid, but only 80% genes were in diads in
120 tetraploid, the normalised value will be 75% - i.e. 75% of the potential triads were formed because we
121 have accounted for the 20% which were already missing in tetraploid.

122 **Correlation of expression levels per family to homoeolog retention in triads**

123 To measure the gene expression level of each TF family we used RNA-seq data from 15 different
124 tissues and developmental stages from Chinese Spring (Choulet *et al.*, 2014). These included tissues
125 from seedling roots and shoots through to grain 30 days after anthesis. We downloaded gene
126 expression data in transcripts per million (tpm) for this dataset from expVIP ([www.wheat-](http://www.wheat-expression.com)
127 [expression.com](http://www.wheat-expression.com)) (Borrill *et al.*, 2016; Ramírez-González *et al.*, 2018). We calculated the mean
128 expression level for each gene across the 15 tissues, and then calculated the median expression level
129 for each TF family. We fitted a linear regression model between log(median expression level per TF
130 family) and the percentage of TFs in triads in the family.

131 **Correlation of tandem duplication per family to homoeolog retention in triads**

132 For all TF genes we defined tandem duplicates as genes which were adjacent in the genome assembly
133 according to their gene IDs ± 3 genes in either direction (gene IDs increase by 100 for adjacent genes
134 in this genome assembly). We allowed 1 or 2 genes between tandem duplicates because a tandem
135 duplication event may have occurred capturing a TF and non-TF in the same duplication event. Each
136 nearby duplicate was counted as one tandem duplication event (i.e. a cluster of 3 TF genes would be
137 counted as 2 tandem duplication events), and the total number of tandem duplication events was

138 divided by the total number of genes in each TF family to calculate the percentage of tandem
139 duplicated genes per TF family. We fitted a linear regression model between the percentage of genes
140 which are tandem duplicates per TF family and the percentage of TFs in triads in the family. We
141 repeated our analysis only considering ± 2 genes (with one gene between them) or ± 1 gene (with no
142 gene between them) as tandem duplicates.

143 **Calculation of homoeolog similarity of expression per family**

144 Using the same data from 15 different tissues and developmental stages from Chinese Spring we
145 filtered to only keep triads where at least 1 homoeolog was expressed >0.5 tpm in one tissue
146 (calculated as the mean value of two biological replicates). To account for differences in expression
147 level between TFs and non-TFs, we normalised the expression level of each triad per tissue to sum to
148 1 as in Ramírez-González *et al.* (2018) before calculating the standard deviation of expression level
149 between homoeologs. For 58 out of 19,391 triads (0.3%) the TF family was inconsistent between
150 homoeologs (e.g. MYB and MYB-related) so the family assigned to two of the three homoeologs was
151 retained. A Mann-Whitney test was used to determine whether the standard deviation within TF triads
152 was different from non-TF triads for each tissue.

153 **Calculation of homoeolog co-expression per family**

154 To calculate the Pearson's correlation between the three homoeologs we used the same data from 15
155 different tissues and developmental stages from Chinese Spring. We filtered to only keep triads where
156 at least 1 homoeolog was expressed >0.5 tpm in one tissue (calculated as the mean value of two
157 biological replicates), and triads where all three homoeologs were expressed (tpm >0 in at least one
158 tissue). The Pearson's correlation was calculated between homoeologs within a triad in a pairwise
159 fashion (A vs B, B vs D, A vs D) and the three correlations were plotted for each triad. To calculate
160 the median Pearson's correlation for TF triads and non-TF triads, the Pearson's correlation values
161 were Z transformed using DescTools v0.99.44 (Signorell, 2021) before obtaining the median, then
162 back-transformed to reduce bias (Corey *et al.*, 1998).

163 As an alternative measure of co-expression we used information about module assignment from a
164 Weight Gene Co-expression Network Analysis (WGCNA) across 850 wheat RNA-samples
165 (Langfelder & Horvath, 2008; Ramírez-González *et al.*, 2018). The co-expression network was built
166 using RefSeq v1.0 annotation. To enable compatibility with our TF annotation which was carried out
167 using RefSeq v1.1 annotation, only genes which were 99% identical with $>90\%$ coverage from v1.0
168 to v1.1 were included in this analysis. To calculate the percentage of triads with homoeologs in the
169 same module only triads in which all three homoeologs had a module assigned, excluding module 0,
170 were considered. Module 0 largely contains genes with invariable expression patterns between
171 samples (Ramírez-González *et al.*, 2018).

172 **Analysis of SNP variation data**

173 To investigate the types of single nucleotide polymorphisms (SNPs) in wheat TFs, we used exome
174 capture data of 811 hexaploid wheat landraces and cultivars representing global genetic diversity (He
175 *et al.*, 2019). Filtered and imputed SNPs (~3 million) were downloaded May 2021 from
176 <http://wheatgenomics.plantpath.ksu.edu/1000EC/>.

177 We selected SNPs in genes in triads and used the Ensembl Variant Effect Predictor (VEP v99.2) to
178 predict the effect of SNPs on these genes (McLaren *et al.*, 2016). From an input of 529,066 SNPs in
179 triad genes, VEP output 1,146,195 SNP effects. We selected 216,285 SNPs predicted in the coding
180 sequence of the canonical transcript of a triad gene. Using R, we filtered to exclude: SNPs which were
181 also splice region variants; missense variants without SIFT scores; and SNPs with >25% missing
182 calls. 210,578 SNPs remained (97% of unfiltered SNPs in coding sequences of canonical transcripts).

183 To exclude potential bias from rare SNPs, we filtered to retain SNPs with a minor allele frequency
184 (MAF) of at least 0.01, resulting in a total of 74,442 SNPs. To focus on SNPs more likely to have a
185 functional effect *in planta*, we only retained SNPs in genes that were expressed at >0.5 tpm in at least
186 one tissue using data from (Choulet *et al.*, 2014). We excluded SNPs in regions that He *et al.* (2019)
187 identified as being under environmental adaptation, improvement selection or within a selective
188 sweep, as positive and purifying selection have similar impacts on nucleotide diversities in
189 populations (Cvijovic *et al.*, 2018). Introgressed sites were also excluded as they would have had a
190 different demographic history compared to the remainder of the genome. Synonymous sites that had
191 more than one annotation were excluded from analyses. This left 16,119 SNPs (1020 TF, 15,099 non-
192 TF).

193 We categorised the SNPs according to variant effect (stop gained, missense and synonymous).
194 Missense mutations were further categorised as deleterious or tolerated according to their Sorting
195 Intolerant from Tolerant (SIFT) prediction (Sim *et al.*, 2012). A SIFT score of ≤ 0.05 is predicted to be
196 deleterious, affecting the protein phenotype and a score > 0.05 is predicted to be tolerated, not
197 affecting phenotype.

198 Per site nucleotide diversity was estimated using VCFtools c0.1.16 (Danecek *et al.*, 2011). Mann-
199 Whitney tests were used to compare the TF and non-TF nucleotide site diversity distributions.
200 Mutation load was estimated by calculating the number of homozygous alternate alleles for each site
201 type, divided by the summed lengths of all the canonical transcripts for TFs and non-TFs separately.
202 A linear regression with mutation load as the response and the category of sites (stop gained,
203 deleterious missense, tolerated missense, synonymous) and the group of genes (TF and non-TF) was
204 fitted, and an ANOVA was performed to test for the significance of the fixed effects. Further, a
205 Tukey's test was used to compare TFs and non-TFs for each site category. Individuals with extreme
206 mutation loads were classed as those with loads in the 2.5% tails in any of the distributions. For the

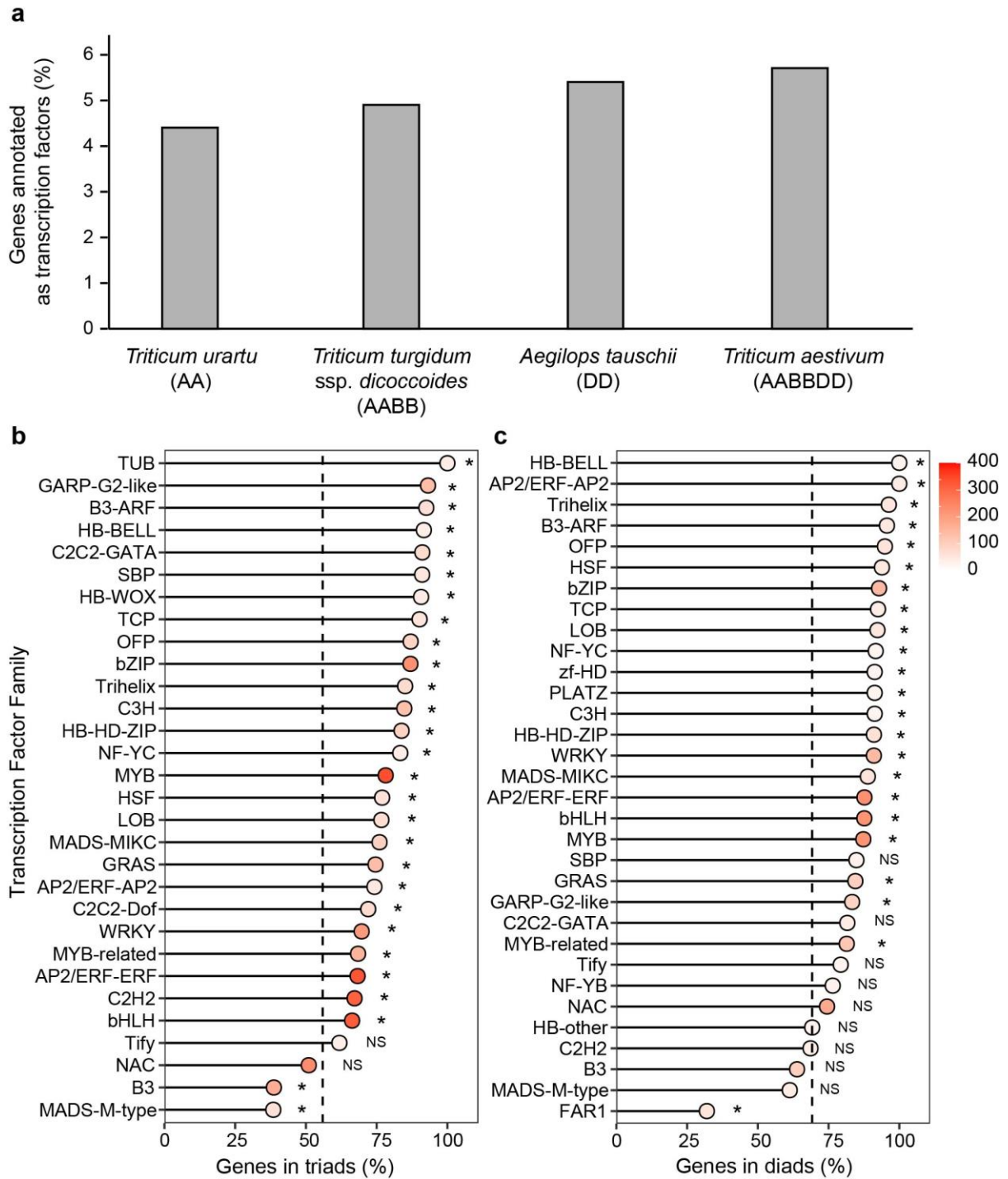
207 TF families plot, we excluded SNPs which are only represented in individuals with extreme mutation
208 loads. We plotted the proportion of SNPs by variant effect for TF families containing more than 10
209 triads and ≥ 5 SNPs and for non-TFs.

210

211 **Results**

212 **TFs homoeologs are retained across polyploidisation events more frequently than non-** 213 **TFs**

214 To explore TF evolution and conservation in polyploid wheat we annotated TFs in the hexaploid *T.*
215 *aestivum* (AABBDD), the tetraploid ancestor *T. turgidum* ssp. *dicoccoides* (AABB) and the diploid
216 ancestral species *T. urartu* (AA) and *A. tauschii* (DD) (Table S1-4). We found that overall the
217 percentage of genes in the genome which were annotated as TFs increased as the ploidy level
218 increased (Figure 1a), from 4.4% in diploid *T. urartu*, to 4.9% in tetraploid *T. turgidum* ssp
219 *dicoccoides*, to 5.7% in hexaploid *T. aestivum*. A higher percentage of genes were TFs in *A. tauschii*
220 (5.4%) than in the other diploid progenitor *T. urartu* (4.4%), although this was still lower than in the
221 hexaploid wheat (5.7%). This supports the hypothesis that TFs are preferentially retained, compared
222 to other types of genes, in polyploid wheat. The retained TFs were distributed similarly across the
223 genomes in tetraploid (50.4% on A genome, 49.6% on B genome) and hexaploid wheat (33.7% on A
224 genome, 33.1% on B genome and 33.3% on D genome), consistent with previous reports that wheat
225 does not show biased sub-genome fractionation associated with preferential loss of genes associated
226 with one subgenome (IWGSC, 2014).



227

228 **Figure 1.** Transcription factor (TF) genes in *T. aestivum* and ancestral species. a) Percentage of genes
 229 annotated as TFs in hexaploid *T. aestivum* and the tetraploid and diploid ancestral species. b)
 230 Percentage of genes in triads in *T. aestivum* TF families with >10 triads and c) Percentage of genes in
 231 diads in *T. turgidum* ssp. *dicoccoides* TF families with >10 diads. In b) and c) the dotted black line
 232 indicates the mean value for non-transcription factors and asterisks (*) denote families which are
 233 significantly different from non-TFs (Fisher's exact test, $p < 0.05$, FDR corrected for multiple testing).
 234 NS= non-significant. The fill colour of the dots indicates the number of genes in the TF family.
 235

236 We hypothesised that the higher proportion of TF genes in polyploid wheats compared to their wheat
237 progenitors were due to the preferential retention of TF homoeologs, whilst other types of genes were
238 less often retained with all homoeologs. Consistent with this hypothesis we found that in polyploid
239 wheat, TFs were more frequently present with all homoeologs than other types of genes. Across TF
240 and non-TF genes in hexaploid *T. aestivum*, 56.7% of genes are in triads with a single A homoeolog, a
241 single B homoeolog and a single D homoeolog. TF genes were more commonly found in triads with
242 70.5% of TFs in triads, compared to other types of genes (55.9% in triads; $p < 0.001$, Fisher's exact
243 test). This enrichment for triads was observed in nearly all TF families (Fig 1b, Fig S1). Similar trends
244 were observed in tetraploid *T. turgidum* ssp. *dicoccoides*. Across TF and non-TF genes 69.8% of
245 genes in the tetraploid were in diads with a single A homoeolog and a single B homoeolog, but this
246 figure rose to 82.5% of TFs, compared to 69.2% of other types of genes ($p < 0.001$, Fisher's exact test).
247 The enrichment for diads was common to most TF families (Figure 1c, Fig S2).

248 In general, TF families with a lower percentage of triads in hexaploid wheat already had a lower
249 proportion of diads in tetraploid. For example, the B3 and MADS-M-type families had fewer
250 triads/diads in both wheat species than non-TF genes, with tetraploid having 63.9% and 61.3% of
251 genes for the B3 and MADS-M-type family in diads respectively, and hexaploid having 38.7% and
252 38.5% of genes in triads respectively (Figure 1b and 1c). The NAC TF family, which is one of the
253 largest TF families in wheat, is one of the less well retained TF families in tetraploid (74.5% of genes
254 in diads), although this is still higher than for non-TFs. However, in hexaploid wheat only 51.0% of
255 NACs are in triads which is lower than for non-TFs. After accounting for gene loss in tetraploid
256 wheat, the B3, MADS-M-type and NAC families in hexaploid wheat still had significantly fewer
257 genes in triads (60.5%, 62.8% and 68.5% respectively) than non-TFs (80.8%; FDR adjusted $p < 0.001$
258 Fisher's exact test). This indicates that homoeolog loss in specific TF families occurred across both
259 polyploidisation steps and was not solely due to pre-existing gene loss in the tetraploid.

260 **Differential conservation of TF families as triads is correlated with expression level and** 261 **tandem duplications**

262 To understand why certain TF families are more prone to homoeolog loss we explored two previously
263 proposed hypotheses. The first is the “highly expressed gene retention idea” (Freeling, 2009), which
264 proposes that gene families which are highly expressed are more likely to be retained with
265 homoeologous copies. Secondly we investigated the “balanced gene drive hypothesis” (Freeling,
266 2009) which proposes that gene families which have more tandem duplications are less likely to be
267 retained with homoeologous copies.

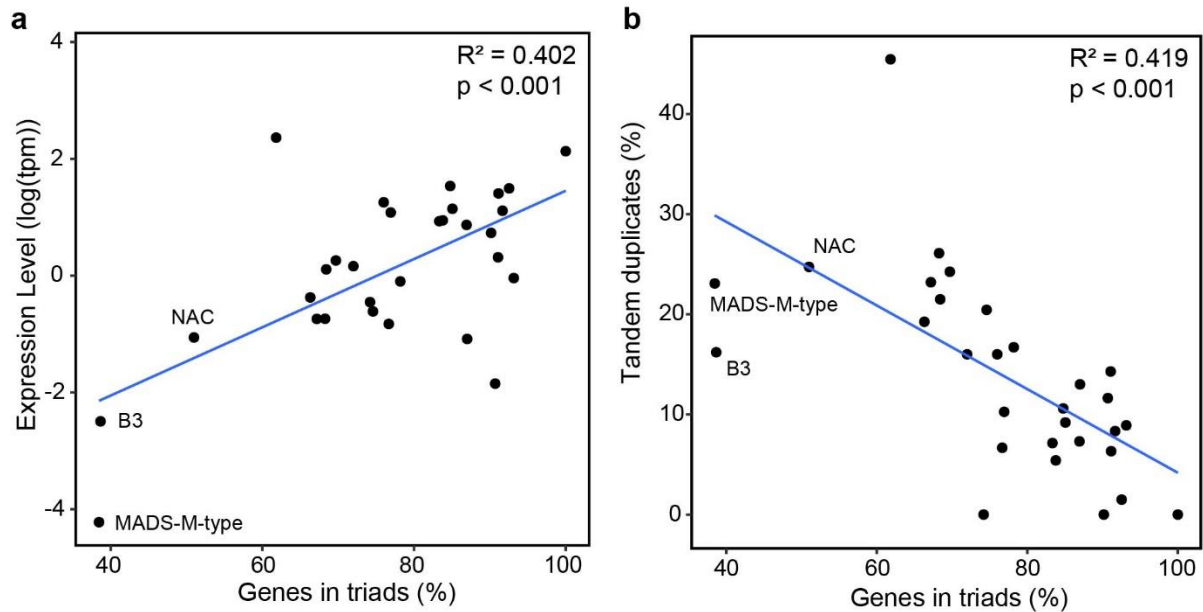
268 To test the correlation between gene expression level and gene retention in hexaploid wheat we used
269 RNA-seq data from 15 different tissues from Chinese Spring from a developmental timecourse
270 (Choulet *et al.*, 2014). We calculated the mean expression level for each gene across the 15 tissues,

271 and then calculated the median expression level for each TF family. Focussing on TF families with
272 >10 triads we found a significant positive correlation between the expression level of the TF family
273 and the percentage of genes in the TF family which are in triads ($R^2 = 0.40$, $p < 0.001$; Figure 2a). This
274 relationship also held across all TF families regardless of size, although the correlation was weaker
275 due to small families that were outliers ($R^2 = 0.21$, $p < 0.001$; Fig S3). Consistent with this relationship,
276 the three TF families with a lower retention of homoeologs in hexaploid wheat than non-TFs (NAC,
277 MADS-M-type and B3) all had low median expression levels (Figure 2a).

278 We also explored the relationship between tandem duplications and gene retention. Focussing on TF
279 families with >10 triads we found that the degree of tandem duplication in a TF family was negatively
280 correlated with the percent of triads within the TF family, consistent with the “balanced gene drive
281 hypothesis” ($R^2 = 0.42$, $p < 0.001$, permitting up to 2 genes between tandem duplicated TFs; Figure
282 2B). This correlation held with a more stringent criteria for tandem duplicates only permitting 1 gene
283 between tandem duplicates ($R^2 = 0.38$, $p < 0.001$; Fig S4a) or 0 genes between tandem duplicates ($R^2 =$
284 0.28 , $p = 0.003$; Fig S4B). These relationships also held when including all TF families regardless of
285 size, although the correlation was weaker ($R^2 = 0.14$ to 0.22 , $p < 0.003$) due to variability within small
286 families (Fig S4c-e). The NAC TF family that had low retention of homoeologs in hexaploid wheat
287 had quite high levels of tandem duplication (Figure 2b). However, the MADS-M-type and B3 TF
288 families had lower levels of tandem duplication than the trendline across all TF families (Figure 2b),
289 suggesting that low expression levels (Figure 2a) may be driving the lack of homoeolog retention in
290 these families. Together these results indicate that different retention levels in individual TF families
291 are associated with gene expression level and the degree of tandem duplication.

292

293



294

295 **Figure 2.** Factors explaining differential retention of homoeologs in different transcription factor (TF)
296 families. a) Median expression level per TF family plotted against the percentage of the TF family in
297 triads for TF families with >10 triads. The mean expression level of each gene in transcripts per
298 million (tpm) was calculated using 15 tissues of Chinese Spring RNA-seq data and these gene level
299 values were used to calculate median expression level within the TF family. b) The percentage of
300 tandem duplicated genes within each TF family plotted against the percentage of the TF family in
301 triads for TF families with >10 triads. TFs were considered to be tandem duplicates when they were
302 up to ± 3 genes away from each other (i.e. up to two genes in between duplicates).

303

304 **TF triads do not show increased sub- or neo-functionalisation of expression or co-** 305 **expression patterns**

306 There are several different mechanisms which can contribute to the retention of homoeologs
307 following polyploidisation. Conant *et al.* (2014) proposed a pluralist framework in which dosage
308 effects, sub-functionalisation and neo-functionalisation interplay to preserve duplicated genes, in a
309 time dependent manner. Although transcriptomics cannot provide a definitive answer about the
310 contributions of these different mechanisms (Conant *et al.*, 2014), it can provide a starting point to
311 understand potential mechanisms operating.

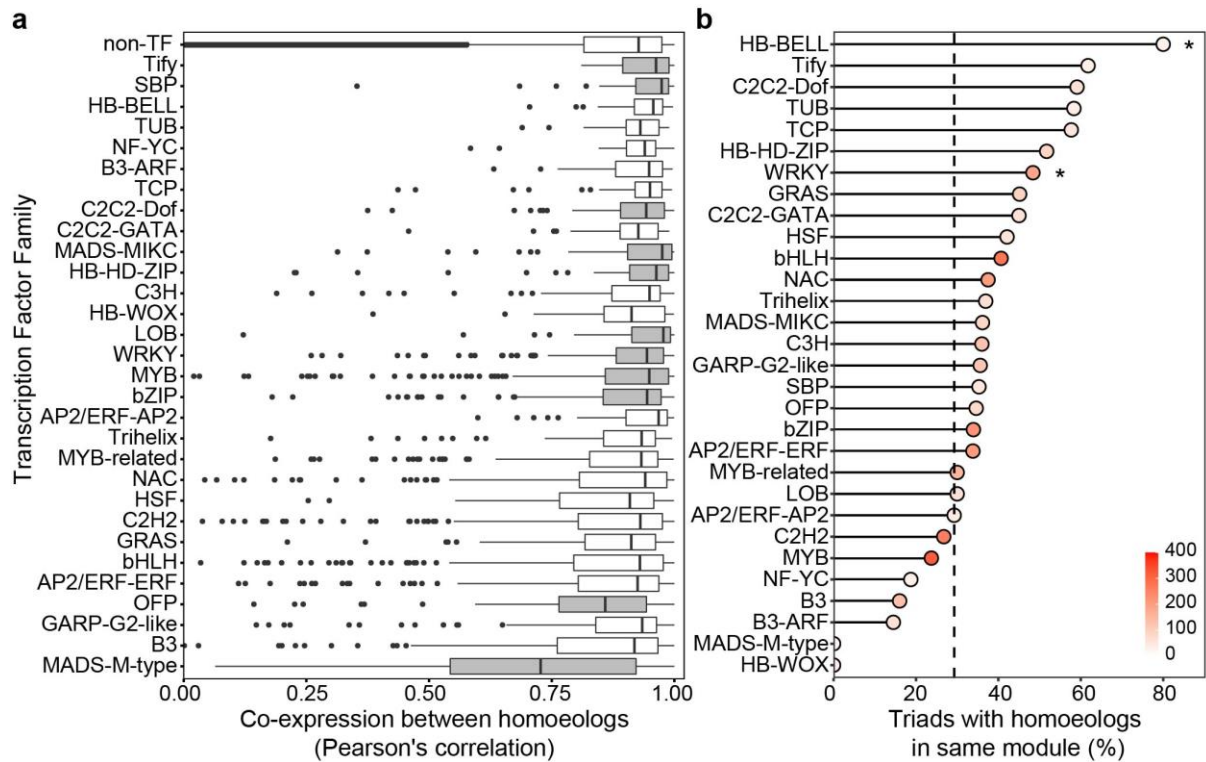
312 First, we used the same RNA-seq samples from 15 tissues from Chinese Spring to test whether TF
313 triads had more different expression levels between homoeologs than non-TF triads, which would
314 support sub- or neo-functionalisation of TF homoeologs at the gene expression level leading to TF
315 retention. We normalised the global expression of each triad so that total expression level of the triad
316 was 1 as described in (Ramírez-González *et al.*, 2018), to account for differences in expression level
317 between TFs and non TFs. We found that the standard deviation between the expression levels of

318 homoeologs within TF triads was not significantly different from non-TF triads in 14 out of 15 tissues
319 (Mann-Whitney test, $p > 0.05$). Only roots at Zadoks stage 39 (flag leaf ligule just visible) had a
320 significantly lower standard deviation between homoeolog expression levels in TF triads than in non-
321 TF triads (median 0.093 for TF triads, 0.099 for non-TF triads, $p = 0.036$, Mann-Whitney test).
322 Overall, the standard deviation between homoeolog expression levels was not higher in TFs than non-
323 TFs in any tissue suggesting that sub- or neo-functionalisation is not occurring at the gene expression
324 level between TF homoeologs globally.

325 Building upon this finding, we explored co-expression between homoeologs across different tissues.
326 We calculated the Pearson's correlation coefficient pairwise between homoeologs across the 15
327 Chinese Spring tissues. Co-expression was higher for TF triads than non-TF triads (Pearson's
328 correlation coefficient 0.938 vs 0.923, p -value < 0.001 , Mann-Whitney test). Amongst TF families
329 with over 10 triads, most TF families showed higher homoeolog co-expression than non-TFs, and the
330 differences were significant for nine TF families (Figure 3a). Two TF families has significantly lower
331 homoeolog co-expression than non-TFs (OFP and MADS-M-type, Figure 3a) The trend for higher co-
332 expression within TF families than non-TFs was also observed in TF families with fewer than 10
333 triads (Fig S5).

334 As an alternative measure of triad co-expression we explored a previously generated co-expression
335 network made using WGCNA across 850 wheat RNA-samples (Langfelder & Horvath, 2008;
336 Ramírez-González *et al.*, 2018). We found that TF homoeologs were more frequently assigned to the
337 same co-expression module than non-TF homoeologs (35.5% vs 29.3%; $p < 0.001$ Fisher's exact test),
338 consistent with our Pearson's correlation approach. A higher level of co-expression in TFs than non-
339 TFs was consistent across most TF families in this WGCNA based approach although the difference
340 was only statistically significant in a few families after adjustment for multiple testing (Figure 3b and
341 Fig S6). TF families which showed higher co-expression were quite consistent with both measures of
342 co-expression, e.g. Tify and WRKY, whilst some other families such as MADS-M-type TFs had
343 lower co-expression using both measures (Figure 3). Overall, we did not find support for higher levels
344 of sub- or neo-functionalisation at the expression or co-expression level in TF triads than in non-TFs,
345 suggesting that other mechanisms such as dosage may be important for TF retention.

346



347

348 **Figure 3.** Co-expression of homoeologs within triads in transcription factor (TF) families with >10
 349 triads. a) Pearson's correlation coefficient between homoeologs across 15 tissues per TF family. TF
 350 families which were significantly different to non-TFs are highlighted in grey (Mann-Whitney test,
 351 $p < 0.05$, FDR corrected for multiple testing). The correlation between non-TF homoeologs is shown in
 352 the top row. b) Homoeologs in same module in 850 sample WGCNA network per TF family. Black
 353 dotted line in b) represents mean value of non-TFs and asterisks (*) denote families which are
 354 statistically significant different from non-TFs (Fisher's exact test, $p < 0.05$, FDR corrected for
 355 multiple testing). The fill colour of the dots in b) indicates the number of genes in the TF family.
 356

357 **Reduced deleterious mutation load in TF triads compared to non-TFs**

358 To investigate how TFs evolve in wheat populations we explored single nucleotide polymorphisms
 359 (SNPs) in TFs and non-TFs using an exome capture dataset of 811 diverse hexaploid wheat cultivars
 360 and landraces (He *et al.*, 2019). We hypothesised that TF triads would accumulate fewer mutations
 361 deleterious to gene function than non-TF triads, which would be consistent with their preferential
 362 retention during polyploidisation. We did not observe significant differences in the distribution of
 363 deleterious or synonymous nucleotide site diversities, estimated using π , between TFs and non-TFs
 364 (Fig S7). π is low when allele frequency is low or high (Fig S8) and, therefore, it does not capture the
 365 deleterious load burden in TFs and non-TFs. To identify the mutational burden, we calculated the
 366 number of homozygous deleterious and synonymous mutations in TF and non-TF triads. Numbers of
 367 homozygous mutations per individual scaled by the total length of all canonical transcripts differ
 368 between TF and non-TF genes (ANOVA, $F=66.5$, $df=1$, $p < 0.001$). There were 32.0% fewer
 369 deleterious missense mutations per kilobase in TFs compared to non-TFs (Figure 4a; $p < 0.001$,
 370 Tukey's test). Frequencies of homozygous stop gained mutation were not significantly different

371 between TFs and non-TFs. However, only 7 stop gained mutations were detected in TFs making the
372 comparison underpowered. There were 5.7% more tolerated missense mutations and 17.6% fewer
373 synonymous mutations per kilobase in TFs compared to non-TF genes. As sites occurring in regions
374 associated with adaptation, introgression, or domestication were removed, the lower synonymous site
375 diversity and load in TFs likely reflects background selection.

376 To explore the distribution of SNP effects across TF families, we plotted the proportion of SNPs of
377 different effects in the coding sequence of TFs in families containing >10 triads and ≥ 5 SNPs (Figure
378 4b). 17 out of 26 TF families had fewer deleterious missense plus stop gained SNPs relative to non-
379 TFs, while 9 had more. The lowest proportion of deleterious plus stop gained SNPs were found in the
380 MADS-M-Type (0.0%), TCP (0.0%) and HSF (5.0%) families, and the highest proportion in the
381 AP2/ERF-AP2 (27.8%), HB-HD-ZIP (21.1%) and C2H2 (21.1%) families. Overall TF families vary
382 widely in the level of deleterious polymorphism in triads.

383

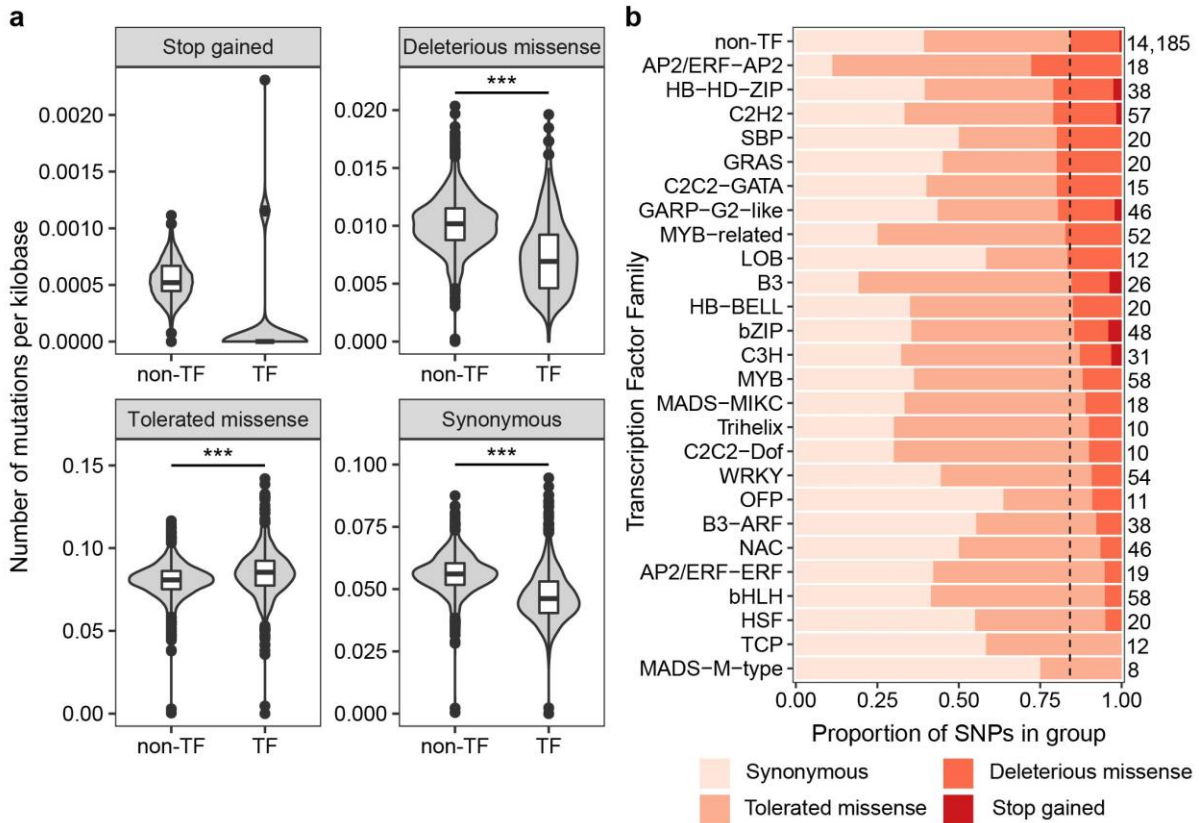
384

385

386

387

388



389

390 **Figure 4.** Transcription factors (TFs) accumulate fewer harmful mutations than non-TFs. a)
 391 Mutational burden in TFs compared to non-TFs for different categories of variants in 811 individuals.
 392 Mutational burden was calculated as the number of homozygous mutations of each type per individual
 393 scaled by the total length of all canonical transcripts for TFs and non-TFs. The total number of
 394 polymorphic sites analysed for TFs and non-TFs were: stop gained (TF = 7, non-TF = 107),
 395 deleterious missense (TF = 125, non-TF = 2,277), tolerated missense (TF = 474, non-TF = 6,723),
 396 synonymous (TF = 414, non-TF = 5,989). *** indicates $p < 0.001$ following a Tukey's test comparing
 397 the different classes of variants. b) The proportion of single nucleotide polymorphisms (SNPs) by
 398 variant effect in TF families containing > 10 triads and ≥ 5 SNPs. These SNPs are in genes expressed
 399 in at least one tissue and have minor allele frequency ≥ 0.01 . The number of SNPs in each group is
 400 shown to the right of the bars. TF families are sorted according to the proportion of deleterious
 401 missense plus stop gained SNPs, and non-TFs are shown in the top row. The black dotted line
 402 represents the split between (synonymous + tolerated missense) and (deleterious missense + stop
 403 gained) SNPs in non-TFs.

404

405

406 Discussion

407 TF retention is observed in both polyploidisation steps in wheat

408 In this study we found that across recurrent polyploidisation steps, wheat retains TF homoeologs more
 409 frequently than non-TF homoeologs. This is consistent with previous studies where TF retention was
 410 observed in both paleopolyploid events (>5 mya) (Li *et al.*, 2016) and neopolyploid events ($\sim 7,500$ to
 411 12,500 years ago) (Zhang *et al.*, 2021). Our findings agree with previous work in allotetraploid cotton
 412 *Gossypium hirsutum*, which formed through hybridisation 1-2 mya (Wendel, 1989), suggesting that

413 TF retention is observed regardless of the time since polyploidisation. Consistent with the gene
414 balance hypothesis, the degree of retention between different TF families was associated with both
415 expression level and the degree of tandem duplication, demonstrating that even within a functional
416 class, this hypothesis can make accurate predictions.

417 **Lack of gene expression support for homoeolog sub- or neo-functionalisation**

418 Using the gene expression data from 15 tissues we found that overall TFs in wheat do not show sub-
419 or neo-functionalisation at the expression or co-expression level which differs from results in other
420 species (Liang & Schnable, 2018). For example TF duplicates formed by paleopolyploidisation events
421 in Arabidopsis during the α , β and γ events (all > 15mya) and maize (5-12 mya) tended to have
422 divergent expression patterns with one copy retaining ancestral expression patterns, whilst the other
423 diverged in expression patterns (Pophaly & Tellier, 2015; Panchy *et al.*, 2019). In Arabidopsis the
424 copy with divergent expression tended to have more novel *cis*-regulatory sites, suggesting that neo-
425 functionalisation might be happening (Panchy *et al.*, 2019). One possible explanation for a lack of
426 divergence in hexaploid wheat homoeolog expression is that the polyploidisation event is much more
427 recent than previously studied paleopolyploidisation events. Alternatively this difference may be
428 because wheat does not show biased genome fractionation (Juery *et al.*, 2021) unlike many other
429 studied allopolyploid species.

430 Although our global analysis did not show divergent patterns of expression, we found that homoeolog
431 co-expression levels were variable between TF families. It was previously reported that a subset of
432 triads which are dynamic in their homoeolog expression between tissues have divergent *cis*-regulation
433 (Ramírez-González *et al.*, 2018) suggesting that a small number of these changes may already be
434 occurring in wheat. Given the highly similar expression and co-expression patterns observed in most
435 TF families, it seems more likely that maintenance of gene dosage underlies TF retention in wheat,
436 although sub- or neo-functionalisation of homoeolog expression may play a role in homoeolog
437 retention in TF families that show weaker co-expression. It would require further study to establish
438 whether *cis*-regulatory changes might explain differences in co-expression between TF families.

439 **Deleterious variation is reduced in TF triads indicating purifying selection**

440 We found that hexaploid wheat TF triads have fewer deleterious missense mutations than non-TF
441 genes. This could reflect selection against gene loss, selection against neo-functionalisation or both,
442 i.e. purifying selection for retaining each homoeolog in its original function. Our results are consistent
443 with Brassica allotetraploids in which TFs were enriched amongst genes without any missense
444 mutations compared to their diploid ancestors (Zhang *et al.*, 2021). However, this contrasts with
445 paleopolyploid TF homoeologs in Brassicas which have more frequent missense mutations than other
446 genes (Zhang *et al.*, 2021). This apparent contradiction could be explained by findings from 37
447 angiosperm species in which TFs were enriched amongst genes which were retained in duplicate for

448 millions of years after WGD but eventually returned to singleton status (Li *et al.*, 2016). Therefore,
449 Zhang *et al.* (2021) hypothesised that selection pressure on TFs is dynamic, with a strong purifying
450 selection for a short period after polyploidisation (hence reduced missense mutations observed in
451 hexaploid wheat), followed by a period with lower selection pressure once the target genes are lost
452 through the diploidisation process. Further studies will be needed on polyploids which formed 1-5
453 million years ago to test this hypothesis.

454 **Differences between TF families**

455 We found that TF families showed quantitative variation in their degree of diad and triad retention,
456 degree of tandem duplication, co-expression within triads and deleterious SNP variation. While most
457 TF families fell within a continuum of variation, the MADS-M-type family was an outlier in several
458 analyses with the lowest percentage of genes in triads (Figure 1b) and exceptionally low co-
459 expression levels (Figure 3) out of all 30 TF families with >10 triads. Selection to retain MADS-M-
460 type genes appears to have been weaker than that for other TF families at both polyploidisation steps,
461 with a gradual decrease from tetraploid to hexaploid wheat. This is consistent with previous reports
462 that genes in the MADS-M-type family experience a high rate of birth-and-death evolution, weaker
463 purifying selection and are less conserved between species than MADS-MIKC genes (Nam *et al.*,
464 2004). Counter-intuitively we found that MADS-M-type triads that are retained, are highly conserved
465 between wheat cultivars with no stop gained mutations or deleterious missense SNPs. One
466 explanation for the contradiction of low MADS-M-type retention during polyploidisation but high
467 conservation within hexaploid wheat cultivars could be due to their role in maintaining speciation
468 boundaries and importance in plant reproduction (Masiero *et al.*, 2011). Alternatively, the apparent
469 high level of conservation may be due to the low number of SNPs in the MADS-M-type family
470 included in our analysis, which is a consequence of the low level of expression of many of these
471 genes. The MADS-M-type family contrasts strongly with the related MADS-MIKC family which
472 behaves more similarly to other TF families and is frequently retained as triads, consistent with a
473 previous study on the MADS-MIKC family (Schilling *et al.*, 2020). While not the focus of this study,
474 there is also likely to be extensive variation within the non-TF genes which consist of a highly
475 heterogeneous set of genes for both function and propensity to be retained as triads.

476 **Implications for wheat breeding**

477 In general we found that TF triads are retained in hexaploid wheat and have relatively few deleterious
478 mutations, consistent with negative consequences to changing TF dosage. However, mutations in TFs
479 which affect dosage, such as dominant mutations, have been very important in wheat breeding for
480 their beneficial agronomic effects, for example to adapt flowering time (e.g. *VRN1* and *PPD1* (Yan *et al.*
481 *et al.*, 2003; Beales *et al.*, 2007)). Therefore, there is the potential to further alter gene dosage of TFs for
482 agronomic benefit. It has been proposed that TFs with lower co-expression across tissues, termed
483 dynamic genes (Ramírez-González *et al.*, 2018), have fewer common targets (Harrington *et al.*, 2020)

484 which might release selective pressure to retain all three copies to maintain genetic balance.
485 Therefore, one promising avenue to influence wheat phenotype by altering TF function would be to
486 focus on TF triads with high co-expression which are more likely to have stronger phenotypic
487 consequences if just one copy is removed. Conversely, one could focus on TF triads with low co-
488 expression because the three homoeologs may have diverged in function, and therefore mutating one
489 copy might lead to a phenotypic effect due to limited genetic redundancy. The recent developments in
490 wheat functional genomics such as TILLING and gene editing now make it possible to test the
491 effectiveness of these strategies (Krasileva *et al.*, 2017; Gao, 2021).

492 Although the possibility to alter the sequence of one homoeolog and induce a phenotypic change in
493 wheat is attractive, there is evidence that this will not be effective for all TFs. For example *VRN1* null
494 mutants in a tetraploid background flower much later than wild type plants and single mutants in the
495 A homoeolog have an intermediate flowering time, however single mutants in the B homoeolog of
496 *VRN1* do not differ in their flowering time to WT (Chen & Dubcovsky, 2012). A similar lack of
497 phenotype in a single mutant was observed for *NAM2* mutants which senesce at a similar time to wild
498 type, whereas null mutants had a significant delay in senescence (Borrill *et al.*, 2019). Therefore, there
499 will still be a need for detailed functional characterisation of individual TFs, although this could be
500 guided by predictions informed by the gene balance hypothesis.

501

502 **Acknowledgements**

503 This work was supported by the UK Biotechnology and Biological Science Research Council
504 (BBSRC) through grant BB/T013524/1 and the Designing Future Wheat Institute Strategic
505 Programme (BB/P016855/1). CEBE received a BBSRC CASE Doctor Training Partnership
506 studentship in collaboration with RAGT Seeds Ltd (BB/M01116X/1). PB also acknowledges funding
507 from the Rank Prize New Lecturer Award. This research was also supported by the NBI Research
508 Computing group through HPC resources and the University of Birmingham's BlueBEAR HPC
509 resources.

510

511 **Author contributions**

512 PB conceived and designed the study with contributions from CEBE and RA. PB wrote the
513 manuscript with contributions from CEBE and RA. PB carried out transcription factor annotation,
514 triad/diad identification, tandem duplication and expression analysis. CEBE carried out SNP variant
515 effect prediction, RA analysed nucleotide site diversity and mutation load, CEBE analysed SNP
516 distribution in TF families. PB, CEBE and RA prepared figures and supplemental files.

517

518 **Data availability**

519 The data that supports the findings of this study are available in the supplementary material of this
520 article and from public repositories mentioned in the methods section. Scripts and input files are
521 available at https://github.com/Borrill-Lab/TF_Triads.

522

523 **References**

- 524 **Avni R, Nave M, Barad O, Baruch K, Twardziok SO, Gundlach H, Hale I, Mascher M,**
525 **Spannagl M, Wiebe K, et al. 2017.** Wild emmer genome architecture and diversity
526 elucidate wheat evolution and domestication. *Science* **357**: 93.
- 527 **Baker CR, Hanson-Smith V, Johnson AD. 2013.** Following gene duplication, paralog
528 interference constrains transcriptional circuit evolution. *Science* **342**: 104-108.
- 529 **Beales J, Turner A, Griffiths S, Snape JW, Laurie DA. 2007.** A Pseudo-Response
530 Regulator is misexpressed in the photoperiod insensitive Ppd-D1a mutant of wheat
531 (*Triticum aestivum* L.). *Theoretical and Applied Genetics* **115**: 721-733.
- 532 **Birchler JA, Riddle NC, Auger DL, Veitia RA. 2005.** Dosage balance in gene regulation:
533 biological implications. *Trends in Genetics* **21**: 219-226.
- 534 **Birchler JA, Veitia RA. 2007.** The Gene Balance Hypothesis: from classical genetics to
535 modern genomics. *The Plant Cell* **19**: 395-402.
- 536 **Blanc G, Wolfe KH. 2004.** Functional divergence of duplicated genes formed by polyploidy
537 during *Arabidopsis* evolution. *The Plant Cell* **16**: 1679-1691.
- 538 **Borrill P, Harrington SA, Simmonds J, Uauy C. 2019.** Identification of transcription
539 factors regulating senescence in wheat through gene regulatory network modelling.
540 *Plant Physiology* **180**: 1740-1755.
- 541 **Borrill P, Harrington SA, Uauy C. 2017.** Genome-wide sequence and expression analysis
542 of the NAC transcription factor family in polyploid wheat. *G3:*
543 *Genes/Genomes/Genetics* **7**: 3019-3029.
- 544 **Borrill P, Ramirez-Gonzalez R, Uauy C. 2016.** expVIP: a customizable RNA-seq data
545 analysis and visualization platform. *Plant Physiology* **170**: 2172-2186.
- 546 **Buggs Richard JA, Chamala S, Wu W, Tate Jennifer A, Schnable Patrick S, Soltis**
547 **Douglas E, Soltis Pamela S, Barbazuk WB. 2012.** Rapid, repeated, and clustered
548 loss of duplicate genes in allopolyploid plant populations of independent origin.
549 *Current Biology* **22**: 248-252.
- 550 **Chen A, Dubcovsky J. 2012.** Wheat TILLING mutants show that the vernalization gene
551 *VRN1* down-regulates the flowering repressor *VRN2* in leaves but is not essential for
552 flowering. *PLOS Genetics* **8**: e1003134.
- 553 **Choulet F, Alberti A, Theil S, Glover N, Barbe V, Daron J, Pingault L, Sourdille P,**
554 **Couloux A, Paux E, et al. 2014.** Structural and functional partitioning of bread wheat
555 chromosome 3B. *Science* **345**: 1249721.
- 556 **Conant GC, Birchler JA, Pires JC. 2014.** Dosage, duplication, and diploidization:
557 clarifying the interplay of multiple models for duplicate gene evolution over time.
558 *Current Opinion in Plant Biology* **19**: 91-98.
- 559 **Corey DM, Dunlap WP, Burke MJ. 1998.** Averaging correlations: expected values and bias
560 in combined Pearson r_s and Fisher's z transformations. *The Journal of General*
561 *Psychology* **125**: 245-261.
- 562 **Cvijovic I, Good BH, Desai MM. 2018.** The effect of strong purifying selection on genetic
563 diversity. *Genetics* **209**: 1235-1278.

- 564 **Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE,**
565 **Lunter G, Marth GT, Sherry ST, et al. 2011.** The variant call format and VCFtools.
566 *Bioinformatics* **27**: 2156-2158.
- 567 **Dubcovsky J, Dvorak J. 2007.** Genome plasticity a key factor in the success of polyploid
568 wheat under domestication. *Science* **316**: 1862-1866.
- 569 **Feldman M, Levy AA. 2012.** Genome evolution due to allopolyploidization in wheat.
570 *Genetics* **192**: 763-774.
- 571 **Force A, Lynch M, Pickett FB, Amores A, Yan YL, Postlethwait J. 1999.** Preservation of
572 duplicate genes by complementary, degenerative mutations. *Genetics* **151**: 1531-1545.
- 573 **Freeling M. 2009.** Bias in plant gene content following different sorts of duplication:
574 tandem, whole-genome, segmental, or by transposition. *Annual Review of Plant*
575 *Biology* **60**: 433-453.
- 576 **Gao C. 2021.** Genome engineering for crop improvement and future agriculture. *Cell* **184**:
577 1621-1635.
- 578 **Harrington SA, Backhaus AE, Singh A, Hassani-Pak K, Uauy C. 2020.** The wheat
579 GENIE3 network provides biologically-relevant information in polyploid wheat. *G3:*
580 *Genes/Genomes/Genetics* **10**: 3675-3686.
- 581 **He F, Pasam R, Shi F, Kant S, Keeble-Gagnere G, Kay P, Forrest K, Fritz A, Hucl P,**
582 **Wiebe K, et al. 2019.** Exome sequencing highlights the role of wild-relative
583 introgression in shaping the adaptive landscape of the wheat genome. *Nature Genetics*
584 **51**: 896-904.
- 585 **Howe KL, Contreras-Moreira B, De Silva N, Maslen G, Akanni W, Allen J, Alvarez-**
586 **Jarreta J, Barba M, Bolser DM, Cambell L, et al. 2020.** Ensembl Genomes 2020—
587 enabling non-vertebrate genomic research. *Nucleic Acids Research* **48**: D689-D695.
- 588 **IWGSC. 2014.** A chromosome-based draft sequence of the hexaploid bread wheat (*Triticum*
589 *aestivum*) genome. *Science* **345**: 1251788.
- 590 **IWGSC, Appels R, Eversole K, Stein N, Feuillet C, Keller B, Rogers J, Pozniak CJ,**
591 **Choulet F, Distelfeld A, et al. 2018.** Shifting the limits in wheat research and
592 breeding using a fully annotated reference genome. *Science* **361**: eaar7191.
- 593 **Juery C, Concia L, De Oliveira R, Papon N, Ramírez-González R, Benhamed M, Uauy**
594 **C, Choulet F, Paux E. 2021.** New insights into homoeologous copy number
595 variations in the hexaploid wheat genome. *The Plant Genome* **14**: e20069.
- 596 **Krasileva KV, Vasquez-Gross HA, Howell T, Bailey P, Paraiso F, Clissold L, Simmonds**
597 **J, Ramirez-Gonzalez RH, Wang X, Borrill P, et al. 2017.** Uncovering hidden
598 variation in polyploid wheat. *Proceedings of the National Academy of Sciences* **114**:
599 E913-E921.
- 600 **Langfelder P, Horvath S. 2008.** WGCNA: an R package for weighted correlation network
601 analysis. *BMC Bioinformatics* **9**: 559.
- 602 **Li Z, Defoort J, Tasdighian S, Maere S, Van de Peer Y, De Smet R. 2016.** Gene
603 duplicability of core genes is highly consistent across all angiosperms. *The Plant Cell*
604 **28**: 326-344.
- 605 **Liang Z, Schnable JC. 2018.** Functional divergence between subgenomes and gene pairs
606 after whole genome duplications. *Molecular Plant* **11**: 388-397.
- 607 **Ling H-Q, Ma B, Shi X, Liu H, Dong L, Sun H, Cao Y, Gao Q, Zheng S, Li Y, et al.**
608 **2018.** Genome sequence of the progenitor of wheat A subgenome *Triticum urartu*.
609 *Nature* **557**: 424-428.
- 610 **Lloyd AH, Ranoux M, Vautrin S, Glover N, Fourment J, Charif D, Choulet F, Lassalle**
611 **G, Marande W, Tran J, et al. 2014.** Meiotic gene evolution: can you teach a new
612 dog new tricks? *Molecular Biology and Evolution* **31**: 1724-1727.

- 613 **Luo M-C, Gu YQ, Puiu D, Wang H, Twardziok SO, Deal KR, Huo N, Zhu T, Wang L,**
614 **Wang Y, et al. 2017.** Genome sequence of the progenitor of the wheat D genome
615 *Aegilops tauschii*. *Nature* **551**: 498-502.
- 616 **Lynch M, Conery JS. 2000.** The evolutionary fate and consequences of duplicate genes.
617 *Science* **290**: 1151-1155.
- 618 **Maere S, De Bodt S, Raes J, Casneuf T, Van Montagu M, Kuiper M, Van de Peer Y.**
619 **2005.** Modeling gene and genome duplications in eukaryotes. *Proceedings of the*
620 *National Academy of Sciences* **102**: 5454-5459.
- 621 **Masiero S, Colombo L, Grini PE, Schnittger A, Kater MM. 2011.** The emerging
622 importance of type I MADS box transcription factors for plant reproduction. *The*
623 *Plant Cell* **23**: 865-872.
- 624 **McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, Flicek P, Cunningham**
625 **F. 2016.** The Ensembl Variant Effect Predictor. *Genome Biology* **17**: 122.
- 626 **Nam J, Kim J, Lee S, An G, Ma H, Nei M. 2004.** Type I MADS-box genes have
627 experienced faster birth-and-death evolution than type II MADS-box genes in
628 angiosperms. *Proceedings of the National Academy of Sciences* **101**: 1910-1915.
- 629 **Nowak MA, Boerlijst MC, Cooke J, Smith JM. 1997.** Evolution of genetic redundancy.
630 *Nature* **388**: 167-171.
- 631 **Ohno S. 1970.** *Evolution by gene duplication*. Berlin Heidelberg: Springer.
- 632 **Panchy NL, Azodi CB, Winship EF, O'Malley RC, Shiu S-H. 2019.** Expression and
633 regulatory asymmetry of retained *Arabidopsis thaliana* transcription factor genes
634 derived from whole genome duplication. *BMC Evolutionary Biology* **19**: 77.
- 635 **Pophaly SD, Tellier A. 2015.** Population level purifying selection and gene expression shape
636 subgenome evolution in maize. *Molecular Biology and Evolution* **32**: 3226-3235.
- 637 **Ramírez-González RH, Borrill P, Lang D, Harrington SA, Brinton J, Venturini L,**
638 **Davey M, Jacobs J, van Ex F, Pasha A, et al. 2018.** The transcriptional landscape of
639 polyploid wheat. *Science* **361**: eaar6089.
- 640 **Schilling S, Kennedy A, Pan S, Jermiin LS, Melzer R. 2020.** Genome-wide analysis of
641 MIKC-type MADS-box genes in wheat: pervasive duplications, functional
642 conservation and putative neofunctionalization. *New Phytologist* **225**: 511-529.
- 643 **Signorell Aema. 2021.** DescTools: tools for descriptive statistics. *R package version 0.99.44*
644 <https://cran.r-project.org/package=DescTools>.
- 645 **Sim N-L, Kumar P, Hu J, Henikoff S, Schneider G, Ng PC. 2012.** SIFT web server:
646 predicting effects of amino acid substitutions on proteins. *Nucleic Acids Research* **40**:
647 W452-W457.
- 648 **Simons KJ, Fellers JP, Trick HN, Zhang Z, Tai Y-S, Gill BS, Faris JD. 2006.** Molecular
649 characterization of the major wheat domestication gene *Q*. *Genetics* **172**: 547-555.
- 650 **Tasdighian S, Van Bel M, Li Z, Van de Peer Y, Carretero-Paulet L, Maere S. 2017.**
651 Reciprocally retained genes in the angiosperm lineage show the hallmarks of dosage
652 balance sensitivity. *The Plant Cell* **29**: 2766-2785.
- 653 **Van de Peer Y, Fawcett JA, Proost S, Sterck L, Vandepoele K. 2009a.** The flowering
654 world: a tale of duplications. *Trends in Plant Science* **14**: 680-688.
- 655 **Van de Peer Y, Maere S, Meyer A. 2009b.** The evolutionary significance of ancient
656 genome duplications. *Nature Reviews Genetics* **10**: 725-732.
- 657 **Wendel JF. 1989.** New World tetraploid cottons contain Old World cytoplasm. *Proceedings*
658 *of the National Academy of Sciences* **86**: 4132-4136.
- 659 **Yan L, Loukoianov A, Tranquilli G, Helguera M, Fahima T, Dubcovsky J. 2003.**
660 Positional cloning of the wheat vernalization gene *VRN1*. *Proceedings of the National*
661 *Academy of Sciences* **100**: 6263-6268.

662 **Zhang H, Xie J, Wang W, Wang J. 2021.** Comparison of *Brassica* genomes reveals
663 asymmetrical gene retention between functional groups of genes in recurrent
664 polyploidizations. *Plant Molecular Biology* **106**: 193-206.
665 **Zheng Y, Jiao C, Sun H, Rosli HG, Pombo MA, Zhang P, Banf M, Dai X, Martin GB,**
666 **Giovannoni JJ, et al. 2016.** iTAK: A program for genome-wide prediction and
667 classification of plant transcription factors, transcriptional regulators, and protein
668 kinases. *Mol Plant* **9**: 1667-1670.

669

670 **Figure legends**

671 **Figure 1.** Transcription factor (TF) genes in *T. aestivum* and ancestral species. a) Percentage of genes
672 annotated as TFs in hexaploid *T. aestivum* and the tetraploid and diploid ancestral species. b)
673 Percentage of genes in triads in *T. aestivum* TF families with >10 triads and c) Percentage of genes in
674 diads in *T. turgidum* ssp. *dicoccoides* TF families with >10 diads. In b) and c) the dotted black line
675 indicates the mean value for non-transcription factors and asterisks (*) denote families which are
676 significantly different from non-TFs (Fisher's exact test, $p < 0.05$, FDR corrected for multiple testing).
677 NS= non-significant. The fill colour of the dots indicates the number of genes in the TF family.

678

679 **Figure 2.** Factors explaining differential retention of homoeologs in different transcription factor (TF)
680 families. a) Median expression level per TF family plotted against the percentage of the TF family in
681 triads for TF families with >10 triads. The mean expression level of each gene in transcripts per
682 million (tpm) was calculated using 15 tissues of Chinese Spring RNA-seq data and these gene level
683 values were used to calculate median expression level within the TF family. b) The percentage of
684 tandem duplicated genes within each TF family plotted against the percentage of the TF family in
685 triads for TF families with >10 triads. TFs were considered to be tandem duplicates when they were
686 up to ± 3 genes away from each other (i.e. up to two genes in between duplicates).

687

688 **Figure 3.** Co-expression of homoeologs within triads in transcription factor (TF) families with >10
689 triads. a) Pearson's correlation coefficient between homoeologs across 15 tissues per TF family. TF
690 families which were significantly different to non-TFs are highlighted in grey (Mann-Whitney test,
691 $p < 0.05$, FDR corrected for multiple testing). The correlation between non-TF homoeologs is shown in
692 the top row. b) Homoeologs in same module in 850 sample WGCNA network per TF family. Black
693 dotted line in b) represents mean value of non-TFs and asterisks (*) denote families which are
694 statistically significant different from non-TFs (Fisher's exact test, $p < 0.05$, FDR corrected for
695 multiple testing). The fill colour of the dots in b) indicates the number of genes in the TF family.

696

697 **Figure 4.** Transcription factors (TFs) accumulate fewer harmful mutations than non-TFs. a)
698 Mutational burden in TFs compared to non-TFs for different categories of variants in 811 individuals.
699 Mutational burden was calculated as the number of homozygous mutations of each type per individual
700 scaled by the total length of all canonical transcripts for TFs and non-TFs. The total number of
701 polymorphic sites analysed for TFs and non-TFs were: stop gained (TF = 7, non-TF = 107),
702 deleterious missense (TF = 125, non-TF = 2,277), tolerated missense (TF = 474, non-TF = 6,723),
703 synonymous (TF = 414, non-TF = 5,989). *** indicates $p < 0.001$ following a Tukey's test comparing
704 the different classes of variants. b) The proportion of single nucleotide polymorphisms (SNPs) by
705 variant effect in TF families containing > 10 triads and ≥ 5 SNPs. These SNPs are in genes expressed
706 in at least one tissue and have minor allele frequency ≥ 0.01 . The number of SNPs in each group is
707 shown to the right of the bars. TF families are sorted according to the proportion of deleterious
708 missense plus stop gained SNPs, and non-TFs are shown in the top row. The black dotted line
709 represents the split between (synonymous + tolerated missense) and (deleterious missense + stop
710 gained) SNPs in non-TFs.

711

712 **Supporting Information**

713 Table S1. Genes in *Triticum aestivum* (hexaploid wheat) assigned into transcription factor families
714 and homoeologous groups.

715 Table S2. *Aegilops tauschii* genes in transcription factor families.

716 Table S3. *Triticum urartu* genes in transcription factor families.

717 Table S4. Genes in *Triticum turgidum* ssp. *diccoides* (tetraploid wheat) assigned into transcription
718 factor families with homoeolog information.

719 Figure S1. Percentage of genes in triads in *Triticum aestivum* transcription factor (TF) families.

720 Figure S2. Percentage of genes in diads in *Triticum turgidum* ssp. *diccoides* transcription factor
721 (TF) families.

722 Figure S3. Median expression level per TF family plotted against the percentage of the transcription
723 factor (TF) family in triads.

724 Figure S4. Relationship between tandem duplication within each TF family and percentage of the
725 transcription factor (TF) family in triads.

726 Figure S5. Pearson's correlation coefficient between homoeologs across 15 tissues per transcription
727 factor (TF) family.

728 Figure S6. Homoeologs in same module in 850 sample WGCNA network per transcription factor (TF)
729 family.

730 Figure S7. Distribution of per-site nucleotide diversity (π) for transcription factors (TF) and
731 background genes (non-TF).

732 Figure S8. Association between per-site nucleotide diversity (π) and allele frequency for transcription
733 factors (TF) and background genes (non-TF).