

1 Separating an allele associated with 2 late flowering and slow maturation 3 of *Arabidopsis thaliana* from 4 population structure

5 Yanjun Zan^{1†}, Xiao Feng^{2,3†}, Zheng Ning³, Weilin Xu^{3,4}, Qianhui Wan^{3,5},
6 Dongyu Zeng⁶, Ziyi Zeng⁷, Yang Liu^{6*}, Xia Shen^{3,8†*}

***For correspondence:**

xia.shen@ed.ac.uk (XS);
liuy353@mail.sysu.edu.cn (YL)

†These authors contributed
equally to this work

7 ¹Department of Medical Biochemistry and Microbiology, Uppsala
8 University, Uppsala, Sweden; ²State Key Laboratory of Biocontrol,
9 Guangdong Provincial Key Laboratory of Plant Resources, Key Laboratory
10 of Biodiversity Dynamics and Conservation of Guangdong Higher
11 Education Institutes, School of Life Sciences, Sun Yat-Sen University,
12 Guangzhou, China; ³Department of Medical Epidemiology and
13 Biostatistics, Karolinska Institutet, Stockholm, Sweden; ⁴Department of
14 Molecular Biology and Genetics, Cornell University, Ithaca, NY, USA;
15 ⁵Department of Mathematics, University of Wisconsin-Madison, Madison,
16 Wisconsin, USA; ⁶School of Life Sciences, Sun Yat-Sen University,
17 Guangzhou, China; ⁷School of Engineering, Sun Yat-Sen University,
18 Guangzhou, China; ⁸Center for Global Health Research, Usher Institute of
19 Population Health Sciences and Informatics, University of Edinburgh,
20 Edinburgh, Scotland, UK

22 **Abstract** Genome-wide association (GWA) analysis is a powerful tool to identify
23 individual loci underlying the complex traits. However, application of GWAS in natural
24 population comes with challenges, especially power loss due to population stratification.
25 Here, we introduce a bivariate analysis approach to a public GWAS dataset of *Arabidopsis*
26 *thaliana*. Using this powerful approach, a common allele, strongly confounded with
27 population structure, is discovered to be associated with late flowering and slow

28 maturation of the plant. The discovered genetic effect on flowering time is further
29 replicated in independent datasets. Using Mendelian randomization analysis based on
30 summary associated statistics from our GWAS and expression QTL (eQTL) scans, we
31 predicted and replicated a candidate gene *AT1G11560* that potentially causes this
32 association. Further analysis with flowering-time-related genes indicates that this locus is
33 also co-selected with many flowering-time-related genes. Our study demonstrates the
34 efficiency of multi-phenotype analysis to uncover hidden genetic loci masked by
35 population structure. The discovered pleiotropic genotype-phenotype map provides new
36 insights into understanding the genetic correlation of complex traits.

37

38 Introduction

39 Evolution has resulted in the speciation and adaptation of various organisms. Although
40 natural selection applies to all kinds of species, the resulted natural population structures
41 have dramatic difference. Especially, due to their lack of mobility, plants, comparing to
42 humans and most animals, have established much stronger population structure adaptive
43 to specific climate conditions (Ch. 11 in *Crawley, 2009*). This makes it difficult, for instance
44 in modern genomic studies, to distinguish genotypic effects on plants' phenotypes from
45 geographical stratification (*Atwell et al., 2010*).

46 Fast-developing genotyping techniques have made genome-wide association study
47 (GWAS) one of the most useful approaches for discovering genomic loci that regulate
48 phenotypes in various organisms (*Hirschhorn and Daly, 2005; Atwell et al., 2010; Huang
49 et al., 2010*). In human GWAS, we learnt that most of the discovered loci associated with
50 complex traits or disease have very small effects (*Yang et al., 2010*). The detected single
51 nucleotide polymorphisms (SNPs) need to have sufficiently high minor allele frequencies
52 (MAFs) for the statistical tests to gain enough power, while high-MAF variants tend to have
53 small effects on the studied phenotypes as these variants were under weak selection
54 pressure. Alleles that have high penetrance on a phenotype are normally under strong
55 selection, resulting in low MAFs of the corresponding SNPs. Thus, a major challenge in
56 human GWAS appears to be the trade-off between statistical power and the effect size
57 of the variant to detect (*Korte and Farlow, 2013; Yang et al., 2014; Wellenreuther and
58 Hansson, 2016*).

59 Although similar trade-off also applies to GWAS in plant populations, e.g. in the natural
60 population of *Arabidopsis thaliana*, in terms of discovery power, the major challenge is dif-
61 ferent. As each individual plant accession is sampled from a specific geographical location
62 in the world, accessions with different genotypes normally have much greater phenotypic

63 differences compared to those in humans. It appears that the genome can explain a
64 large proportion of variation in the plant phenotype, however, the population structure in
65 nature makes such a genomic effect heavily confounded with the environmental effect
66 due to geographical stratification. Therefore, there can be a number of alleles, who have
67 large genetic effects on a certain phenotype, but masked by the population structure.

68 As a community based effort, over 1000 natural *A. thaliana* accessions have been col-
69 lected from worldwide geographical locations (**1001 Genomes Consortium, 2016; Kawakatsu**
70 **et al., 2016**). Most of those plants have been sequenced for genome, transcriptome, and
71 even methylome, and these datasets have been made publicly available for worldwide
72 researchers. Many accessions in this collection have been phenotyped for developmen-
73 tal, metabolic, inomics, stress resistance traits (**Atwell et al., 2010**), and more and more
74 phenotypes are gradually releasing. Previous analysis in those datasets have revealed
75 substantial connections between genotypic and phenotypic variations in this species. The
76 application of association mapping have provided insights to the genetic basis of complex
77 traits (**Atwell et al., 2010; Shen et al., 2012; Wang et al., 2017**), adaptation (**Shen et al.,**
78 **2014**) and evolutionary process. Nevertheless, many essential genotype-phenotype links
79 are still difficult to establish based on the current GWAS data, due to the substantial pop-
80 ulation stratification highly correlated with the sampling origins of the plants. Therefore,
81 novel powerful analyses are required to further uncover hidden genetic regulation.

82 Based on publicly available *A. thaliana* datasets (**Atwell et al., 2010; Schmitz et al., 2013;**
83 **1001 Genomes Consortium, 2016; Kawakatsu et al., 2016**), here, we aim to use a bivariate
84 analysis method that combines the discovery power of two correlated phenotypes (**Shen**
85 **et al., 2017**), in order to map novel pleiotropic loci that simultaneously regulate both traits.
86 We interpret the statistical significance with a double-trait genotype-phenotype map. We
87 try to replicate and *in silico* functionally investigate the candidate genes that may drive
88 such associations.

89 RESULTS

90 **Bivariate genomic scan identifies a hidden locus simultaneously as-** 91 **sociated with flowering and maturation periods**

92 We re-analyzed a public dataset of a natural *A. thaliana* collection, where 43 developmental
93 phenotypes and 23 flowering-time-related phenotypes were previously published (**Atwell**
94 **et al., 2010**). The number of accessions with measured phenotypes varies from 93 to
95 193 with a median of 147 (Supplementary Table 1). We first excluded all variants with
96 minor allele frequencies (MAF) less than 0.1 and performed single-trait GWA analysis for
97 all these traits based on a linear mixed model, so that the confounded genetic effects due

98 to population stratification is adjusted. We then applied our recently developed multi-trait
99 GWAS method (*Shen et al., 2017*) to all pairwise combination of the phenotypes (Materials
100 & Methods). One novel locus, in one of the pairwise test, reached the most stringent 5%
101 Bonferroni-corrected genome-wide significance threshold for the 2,145 pairs of traits and
102 173,220 variants, i.e. $p < 1.35 \times 10^{-10}$ (Table 1, Fig. 1a). This signal also reaches single-trait
103 genome-wide significance in other six pairs of traits highly correlated with the top pair
104 (Supplementary Fig. 1), without Bonferroni-correction for the number of tested trait pairs
105 (Table 1, Supplementary Fig. 3-8).

106 For the most significant trait combination, 2W (days to flowering time under long day
107 with vernalized for 2 weeks) and MT GH (maturation period), the linkage disequilibrium
108 (LD) block of this locus (LD $r > 0.7$) covers about a 260 kb interval on chromosome 1, with
109 a top variant at 3,906,923 bp (double-trait $p = 9.9 \times 10^{-12}$, Fig. 1b, Table 1). The detected
110 locus shows joint effects on flowering and maturation, where the effect on flowering
111 time (2W) is notably large (15.3 days), and that on maturation period (MT GH) is 2.5 days
112 (Table 1). These correspond to narrow-sense heritability values of 24% and 10% of the
113 two phenotypes, respectively.

114 [TABLE 1 ABOUT HERE]

115 [FIGURE 1 ABOUT HERE]

116 **Double-trait analysis is sufficiently powerful to overcome the con-** 117 **founding population structure**

118 The detected joint-effect locus was missed in the corresponding single-trait GWA analysis
119 of 2W (effect = 15.3, $p = 2.26 \times 10^{-5}$ after correcting for population stratification) and that
120 of MT GH (effect = 2.5, $p = 3.70 \times 10^{-5}$). Notably, this locus was not even detectable at
121 the genome-wide significance level in a much larger population of more than 1,000 *A.*
122 *thaliana* accessions (*Kawakatsu et al., 2016; 1001 Genomes Consortium, 2016*) due to its
123 severe confounding with the natural population structure. The statistical significance
124 can only be identified when considering the joint distribution of the bivariate statistic.
125 According to the genome-wide Z-scores (student t-statistics), these two phenotypes are
126 negatively correlated, as the plant's lifespan is relatively stable (estimated and observed
127 phenotypic correlation = -0.55 and -0.68, respectively). However, the observed effects on
128 the two traits are both substantially positive, showing sufficient deviation from the joint
129 distribution that led to bivariate statistical significance (Fig. 2).

130 [FIGURE 2 ABOUT HERE]

131 The strong confounding with the population structure can also be visualized by the
132 allele frequency distribution of the top associated SNP across different *A. thaliana* sub-

133 populations based on the genome re-sequencing data from the *A. thaliana* 1001-genomes
134 project (**1001 Genomes Consortium, 2016**, Fig. 3). The sub-populations were divided by
135 admixture analysis using ADMIXTURE (**1001 Genomes Consortium, 2016; Alexander et al.,**
136 **2009**). The plus allele increasing flowering time was predominantly found in Sweden
137 and almost fixed in the Northern Sweden population (Fig. 3b; allele frequency = 0.97 in
138 Northern Sweden and 0.51 in Southern Sweden). Overall, the phenotype, e.g. flowering
139 time at 10 °C, highly correlates with the frequency of the plus allele (Fig. 3). The genotype
140 at this locus follows a latitude decline, where the northern accessions are enriched with
141 the plus allele and the southern accessions are enriched with the minus allele (Fig. 3). This
142 spatially imbalanced enrichment shows strong confounding with the population structure,
143 which is why standard single-trait GWAS loses power substantially.

144 [FIGURE 3 ABOUT HERE]

145 **Replication of the detected genetic effect on flowering time**

146 Although we are lack of an independent dataset of *A. thaliana* maturation duration
147 to replicate the bivariate statistical test, datasets containing additional independent *A.*
148 *thaliana* flowering time measurements are available. We downloaded a flowering time
149 GWAS dataset measured in 1,135 natural accessions from the 1001-genomes project
150 collection (**1001 Genomes Consortium, 2016**) and performed a single-trait association
151 analysis of our discovered top SNP with linear mixed model correction for the population
152 structure. The genetic effect was significantly replicated for flowering time at 10 °C (effect
153 = 1.7 days, $p = 0.037$) and flowering time at 16 °C (effect = 3.6 days, $p = 0.003$). The effects
154 on flowering time in the replication sample appear to be smaller than in the discovery
155 population, possibly due to Winner's curse in the discovery phase.

156 We also screened literature for conventional quantitative trait loci (QTL) studies in
157 intercrosses using natural *A. thaliana* accessions. Our detected signal is underneath a
158 reported QTL peak for flowering time from an intercross between a Swedish and an
159 Italian accession (**Dittmar et al., 2014**, Supplementary Fig. 2). This, together with the
160 replication above, justifies the detected association. Although the discovered genetic
161 effect on maturation period is not directly replicated, the effect does exists when the effect
162 on flowering is justified, as the pleiotropic signal must be driven by both phenotypes.

163 **Prediction and replication of candidate genes using summary-level** 164 **Mendelian randomization**

165 As a community-based effort, all the natural *A. thaliana* accessions from the 1001-genomes
166 project were measured for their transcriptome (**Kawakatsu et al., 2016; 1001 Genomes**
167 **Consortium, 2016**). Such a public gene expression dataset allows us to predict candidate

168 genes underlying the association signal. We extracted the expression levels of 19 genes
169 within a \pm 20kb window around the top associated SNP using RNA-seq gene expression
170 measurements from 140 accessions (*Schmitz et al., 2013*). Among these, the distributions
171 of 14 gene expression phenotypes significantly deviate from normality (Kolmogorov-
172 Smirnov test statistic > 0.8), and these genes were filtered out due to potential unreliable
173 measurements (*Zan et al., 2016*). The remaining 5 genes were passed onto eQTL mapping
174 at the discovered locus (Materials & Methods).

175 Based on the locus-specific eQTL mapping summary statistics, we applied the recently
176 developed Summary-level Mendelian randomization (SMR) method (*Zhu et al., 2016*)
177 to predict potential candidate genes among these five genes. The analysis integrates
178 summary association statistics from GWAS and eQTL scan to predict functional candidate
179 genes using multiple-instrument Mendelian randomization (*Burgess et al., 2015*), where
180 the complementary Heterogeneity In Dependent Instruments (HEIDI) test checks that
181 the gene expression and flowering time share the same underlying causal variant. One
182 significant candidate *AT1G11560* was detected after Bonferroni correction for five tests (Fig.
183 4, Table 2). This candidate gene prediction result was also replicated using an independent
184 eQTL mapping dataset (*Kawakatsu et al., 2016*).

185 [TABLE 2 ABOUT HERE]

186 [FIGURE 4 ABOUT HERE]

187 **Indication of co-selection with genes in flowering-related pathways**

188 As flowering time is a well-known polygenic trait, we expect multiple loci to be involved
189 and possibly co-selected as a result of parallel evolution. Therefore, we explored the
190 evidence of co-selection by associating the expression values of 288 known genes in
191 flowering-time-related pathways and 1 gene in the maturation related pathway with
192 our top SNP using transcriptome data from 648 *A. thaliana* accessions (*1001 Genomes*
193 *Consortium, 2016*, Materials & Methods). In total, six genes (*NF-YA8, AT5G53360, SPL15,*
194 *AGL42, FLC, AGL20*) were associated with our top SNP (false discovery rate < 0.05), where,
195 conservatively, four genes (*AT5G53360, AGL42, FLC, AGL20*) were replicated after Bonferroni
196 correction for six tests using data from an independent collection of 140 *A. thaliana*
197 (*Schmitz et al., 2013*, Table 3). This indicates that co-selected genes in multiple pathways
198 determine the flowering time variation in nature, and our detected locus contributes to a
199 part of that.

200 [TABLE 3 ABOUT HERE]

DISCUSSION

A serious issue of GWAS in natural population is the confounding between true underlying genetic effects and the population structure, which can lead to spurious associations between genotypes and phenotypes if population stratification is not properly adjusted (*Korte and Farlow, 2013; Yang et al., 2014; Wellenreuther and Hansson, 2016*). Incorporation of the random polygenic effect using linear mixed models can effectively control the population structure, but such correction often compromises the true signals. Here, we applied a bivariate analysis to a classic dataset and successfully separated a locus from strong population structure. The detected allele is associated with late flowering and slow maturation of *A. thaliana*, which was corrected away by the linear mixed model in standard single-trait analysis. The replication of the genetic effect on flowering time in an old intercross linkage analysis and another independent dataset improves the confidence of this association. The discovered association is a typical example that jointly modeling phenotypes that share genetic basis can boost discovery power and reveal pleiotropic genotype-phenotype map at the same time.

Together with our recent application of multivariate analysis in human isolated populations (*Shen et al., 2017*), the results further indicate that multi-phenotype analysis is an effective approach to detect hidden loci that are lack of discovery power in single-phenotype analysis thus is worth testing in broader applications. Multivariate analysis appears to have the greatest power when the locus-specific genetic correlation does not agree with the natural phenotypic correlation. For instance, like the discovery here, for two traits that are negatively correlated, loci that generate positive genetic correlation between the traits tend to have good chance to be detected in a joint analysis.

In GWAS, phenotypes are usually chosen based on morphological, physiological or economical features. Those features are usually feasible and simple to quantify; however, they might not be directly representative for the underlying genetic or biological factor that we try to detect. Fortunately, a certain degree of biological pathway sharing among complex traits is common, i.e. pleiotropy (*Visscher and Yang, 2016*). Nowadays, it is very common that multiple phenotypes are measured for same individuals in many GWAS datasets, especially in omics study where thousands of phenotypes are measured. Instead of focusing on one phenotype at a time, it is of essential value to jointly model multiple phenotypes, attempting to detect pleiotropic loci that affect multiple traits with biological relevance.

In this study, all the pairs of traits that are associated with the detected locus contain at least one flowering-time trait, and nearly all of them have maturation duration involved. Detection of the novel locus in a bivariate analysis indicates shared genetic basis for the two types of developmental traits, which measure the lengths of two important

238 period during the plant's life time. By integrating the expression level information and
239 GWAS result using SMR/HEIDI test, we were able to predict candidate genes in this region.
240 However, further work beyond the scope of this paper is still required to establish the
241 molecular biological basis underlying the replicate association.

242 Many genetic variants affecting flowering time have been mapped and many genes
243 promoting flowering times have been well characterized using standard lab accession,
244 Col-0 (*Brachi et al., 2010*). Unlike simple traits, where only one or a few alleles are
245 driving the trait's variation, there are many more variants throughout the genome that
246 contribute to the variation of flowering time. The associations between our top SNP and
247 the expression of many flowering-time-related genes serve as evidence of co-selection or
248 parallel adaptation.

249 In conclusion, our study demonstrates the efficiency of joint modeling multiple-
250 phenotypes which overcomes severe power loss due to population stratification in associ-
251 ation studies. We discover and replicate a pleiotropic allele that regulate flowering and
252 maturation periods simultaneously, providing novel insights in understanding the plant's
253 development over life time. By integrating gene expression information with the GWAS
254 results, we predict a functional candidate underneath the associated genomic region.
255 Analysis of gene expression with other flowering-time-related genes show evidence of
256 co-selection of the predicted candidate with many genes in flowering-time pathways.
257 We encourage wider applications of such a multivariate framework in future analyses of
258 genomic data.

259 **Acknowledgements**

260 X.S. was in receipt of a Swedish Research Council (VR) grant (No. 537-2014-371). Interna-
261 tional collaboration within this work was partly supported by the Swedish Foundation for
262 International Cooperation in Research and Higher Education (STINT) initiation grant to X.S.
263 (No. IB2015-6000) and Karolinska Institutet travel grant (No. 2017-00534).

264 **Author contributions**

265 X.S. initiated and coordinated the study. Y.Z. and X.F. performed the main data analysis.
266 Z.N. and X.S. contributed to statistical modeling and interpretation. W.X., Q.W., D.Z. and
267 Z.Z. contributed to data processing. Y.Z., X.F. and X.S. wrote the manuscript. Y.L. and X.S.
268 supervised the study.

269 **Competing interests statement**

270 The authors declare no competing financial interests.

271 Figure legends

272 **Figure 1: Bivariate genome-wide association analysis of two developmental trait.**

273 2W: Days to flowering time (FT) under long day (LD) with vernalized for 2 wks at 5°C,
274 8hrs daylight, MT GH: Maturation period. (a) Manhattan plots comparison of bivariate
275 and univariate analysis results, where the novel variants only discoverable when com-
276 bining two phenotypes are shown in green. The horizontal dashed line represents a 5%
277 Bonferroni-corrected genome-wide significant threshold for the number of variants and
278 also the number of tested trait pairs, respectively. (b) Zooming in the novel locus detected
279 using bivariate analysis. r : linkage disequilibrium measured as correlation coefficient
280 between the top variant and each variant in the region. .

281 **Figure 2: Hexbin scatter plot comparing all Z-scores of the two traits across the**

282 **genome, showing the bivariate statistical significance of the detected locus.** The
283 top variants of the locus is marked on the edge of the empirical bivariate normal distribu-
284 tion with a red circle. The black line with a slope of -1 is provided as a visual guide.

285 **Figure 3: a) Flowering time variation (10°C) among different sub-populations of *Ara-***

286 *bidopsis thaliana*. These populations are divided by admixture analysis (**1001 Genomes**
287 **Consortium, 2016**); b) Frequency of the top associated SNP at chromosome 1, 3,906,923
288 bp in different sub-populations. The association between the structure of the phenotype
289 and that of the allele frequency shows the population confounding at this locus.

290 **Figure 4: Prioritized candidate genes at the detected locus for flowering time using**

291 **SMR analysis.** a) Manhattan plot of association between flowering time at 10°C and SNPs
292 around 40kb of top associated SNP in bivariate analysis. The diamonds highlight top eQTL
293 for individual genes; b) Manhattan plot of association between expression of *AT1G11560*
294 and SNPs around 40kb of top associated SNP in bivariate analysis. Genes tested in SMR
295 analysis are highlighted using arrows; c) Similar linkage-disequilibrium structure at the
296 locus for the corresponding populations of GWA and eQTL analyses.

Tables

Table 1: Discovery and replication analyses results for the novel pleiotropic locus.

Reported association statistics are for the top variant at the locus for each pair of traits.

¹LD: Days to flowering time under Long Day. ²OW: Days to flowering time under long day without vernalization. ³2W: Days to flowering time under long day with vernalized for 2 weeks at 5°C, 8hrs daylight. ⁴4W: Days to flowering time under long day with vernalized for 4 weeks at 5°C, 8 hrs daylight. ⁵OW GH FT: Days to flowering time (greenhouse). ⁶FT GH: Days to flowering (greenhouse). ⁷MT GH: Maturation period (greenhouse), 20°C, 16 hrs daylight. ⁸RP GH: Reproduction period (greenhouse), 20°C, 16 hrs daylight. ⁹RA: Reference allele. ¹⁰EA: Effect allele. ¹¹MAF: Minor allele frequency. ¹²Correlation refers to observed phenotypic correlation. ¹³FT: Flowering time.

<i>Double-trait Analysis</i>									
Trait 1	Trait 2	Chr	Position	RA ⁹	EA ¹⁰	MAF ¹¹	<i>P</i>	Correlation ¹²	
LD ¹	MT GH ⁷	1	3904658	T	A	0.20	6.3×10 ⁻⁹	-0.39	
OW ²	MT GH ⁷	1	3896072	G	T	0.20	8.4×10 ⁻⁹	-0.58	
2W ³	MT GH ⁷	1	3906923	T	C	0.22	9.9×10 ⁻¹²	-0.68	
2W ³	RP GH ⁸	1	3978064	A	C	0.27	1.3×10 ⁻⁸	-0.17	
4W ⁴	MT GH ⁷	1	3906923	T	C	0.22	3.1×10 ⁻⁹	-0.64	
OW GH FT ⁵	MT GH ⁷	1	3906923	T	C	0.22	1.8×10 ⁻⁸	-0.36	
FT GH ⁶	MT GH ⁷	1	3896072	G	T	0.20	1.5×10 ⁻⁸	-0.60	

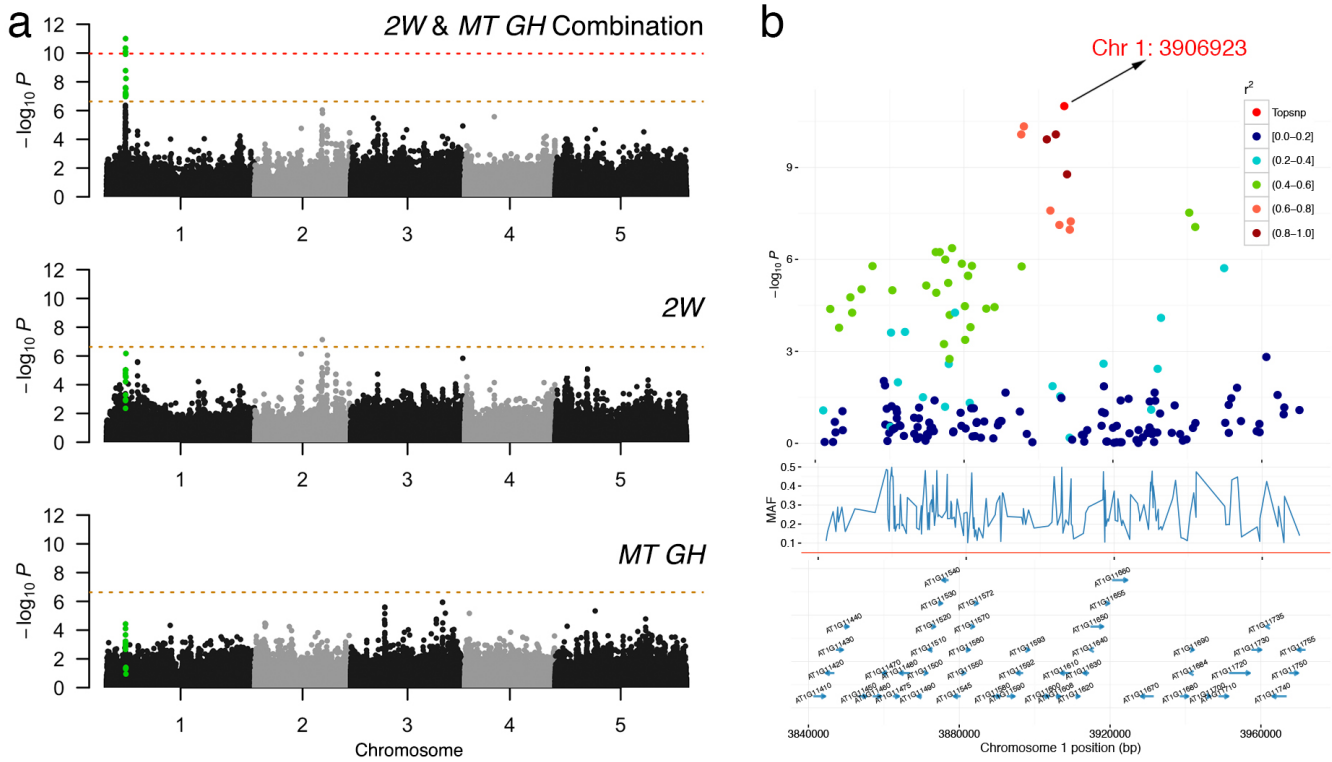
<i>Single-trait Analysis</i>						<i>Replication</i>			
Trait 1			Trait 2			FT ¹³ 10°C		FT ¹³ 16°C	
Effect	<i>P</i>	<i>h</i> ²	Effect	<i>P</i>	<i>h</i> ²	Effect	<i>P</i>	Effect	<i>P</i>
33.5	5.6×10 ⁻⁶	0.22	2.42	6.0×10 ⁻⁴	0.07	2.26	1.0×10 ⁻²	4.45	4.9×10 ⁻⁴
17.3	1.6×10 ⁻⁴	0.17	2.59	2.1×10 ⁻⁴	0.09	1.95	2.3×10 ⁻²	3.96	1.5×10 ⁻³
15.3	2.3×10 ⁻⁵	0.24	2.47	3.7×10 ⁻⁵	0.10	1.72	3.7×10 ⁻²	3.56	3.0×10 ⁻³
19.7	6.8×10 ⁻⁷	0.26	2.65	1.6×10 ⁻³	0.06	1.57	5.6×10 ⁻²	2.57	3.4×10 ⁻²
11.6	1.7×10 ⁻³	0.16	2.47	3.7×10 ⁻⁵	0.10	1.72	3.7×10 ⁻²	3.56	3.0×10 ⁻³
25.8	3.8×10 ⁻⁵	0.21	2.47	3.7×10 ⁻⁵	0.10	1.72	3.7×10 ⁻²	3.56	3.0×10 ⁻³
14.9	1.8×10 ⁻³	0.11	2.59	2.1×10 ⁻⁴	0.09	1.95	2.3×10 ⁻²	3.96	1.5×10 ⁻³

316 **Table 2: Summary of the SMR/HEIDI analysis results.** ¹Top SNP: The top SNP in ex-
 317 pression QTL analysis. ²MAF: Minor allele frequency of the top associated SNP. ³ P_{SMR} :
 318 p-value from SMR using a collection of 140 *A. thaliana* accessions. ⁴ P_{HEIDI} : p-value from
 319 HEIDI test using a collection of 140 *A. thaliana*. ⁵ P_{SMR} : p-value from SMR using a second
 320 collection of 648 accessions. ⁶ P_{HEIDI} : p-value from HEIDI test using a second collection of
 321 648 accessions.

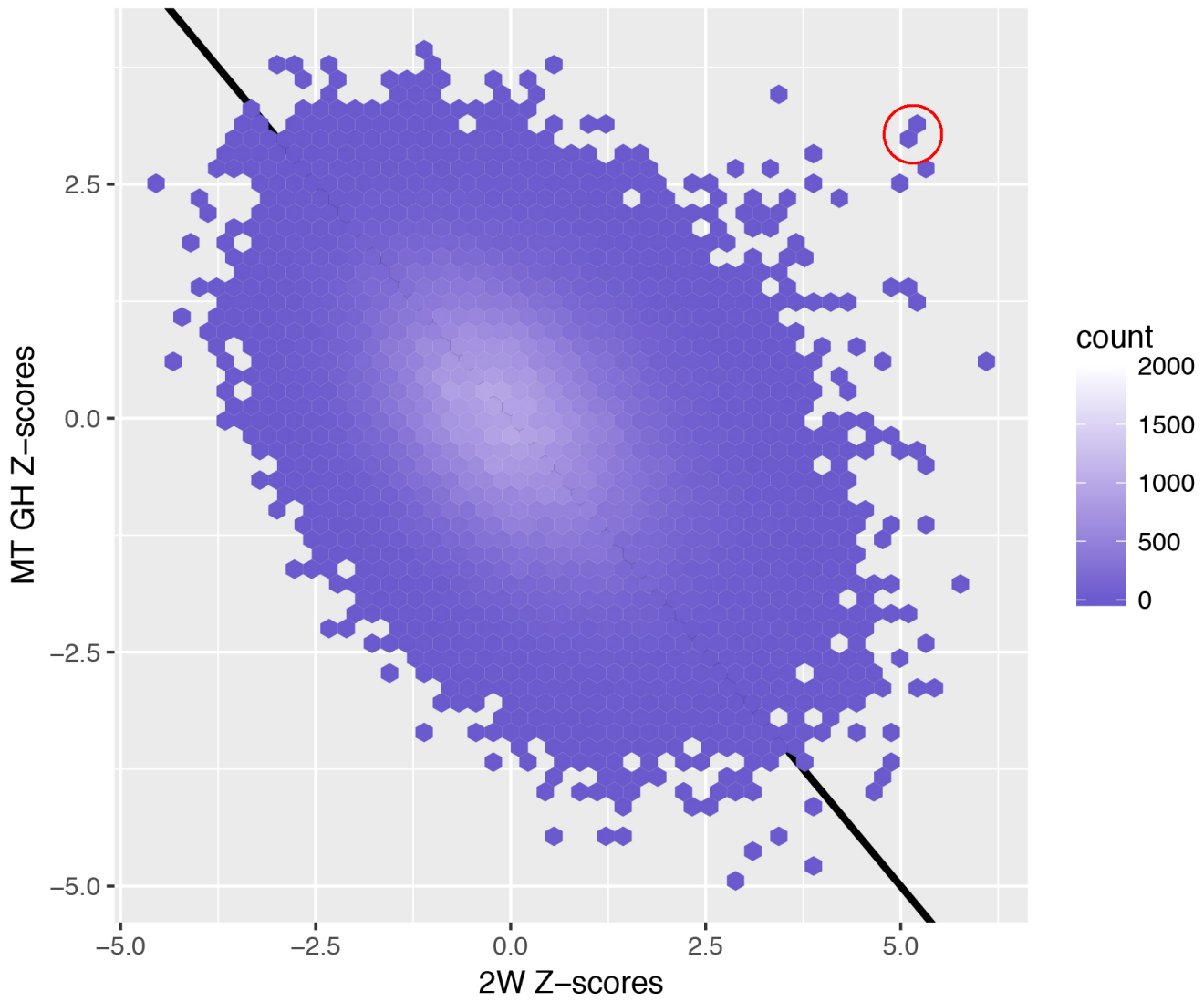
Gene	Top SNP ¹	MAF ²	P_{SMR}^3	P_{HEIDI}^4	P_{SMR}^5	P_{HEIDI}^6
<i>AT1G11560</i>	Chr1:3881093	0.34	6.8×10^{-3}	4.8×10^{-1}	3.2×10^{-2}	2.6×10^{-1}
<i>AT1G11655</i>	Chr1:3874970	0.39	4.1×10^{-2}	9.7×10^{-2}	5.9×10^{-1}	NA
<i>AT1G11690</i>	Chr1:4299126	0.04	3.7×10^{-1}	NA	9.4×10^{-1}	NA
<i>AT1G11590</i>	Chr1:3716355	0.11	5.0×10^{-1}	NA	2.2×10^{-2}	1.5×10^{-1}
<i>AT1G11482</i>	Chr1:3830013	0.63	8.2×10^{-1}	NA	1.5×10^{-1}	NA

325 **Table 3: Genes in flowering-time pathways whose expression are associated with**
 326 **the detected locus.** ¹p-value from a expression dataset generated from 648 accessions
 327 in the *A. thaliana* 1001-genomes project (*Kawakatsu et al., 2016*). ²FDR value computed
 328 from p-value¹. ³Replication p-value from another subset of 140 accessions (*Schmitz et al.,*
 329 *2013*).

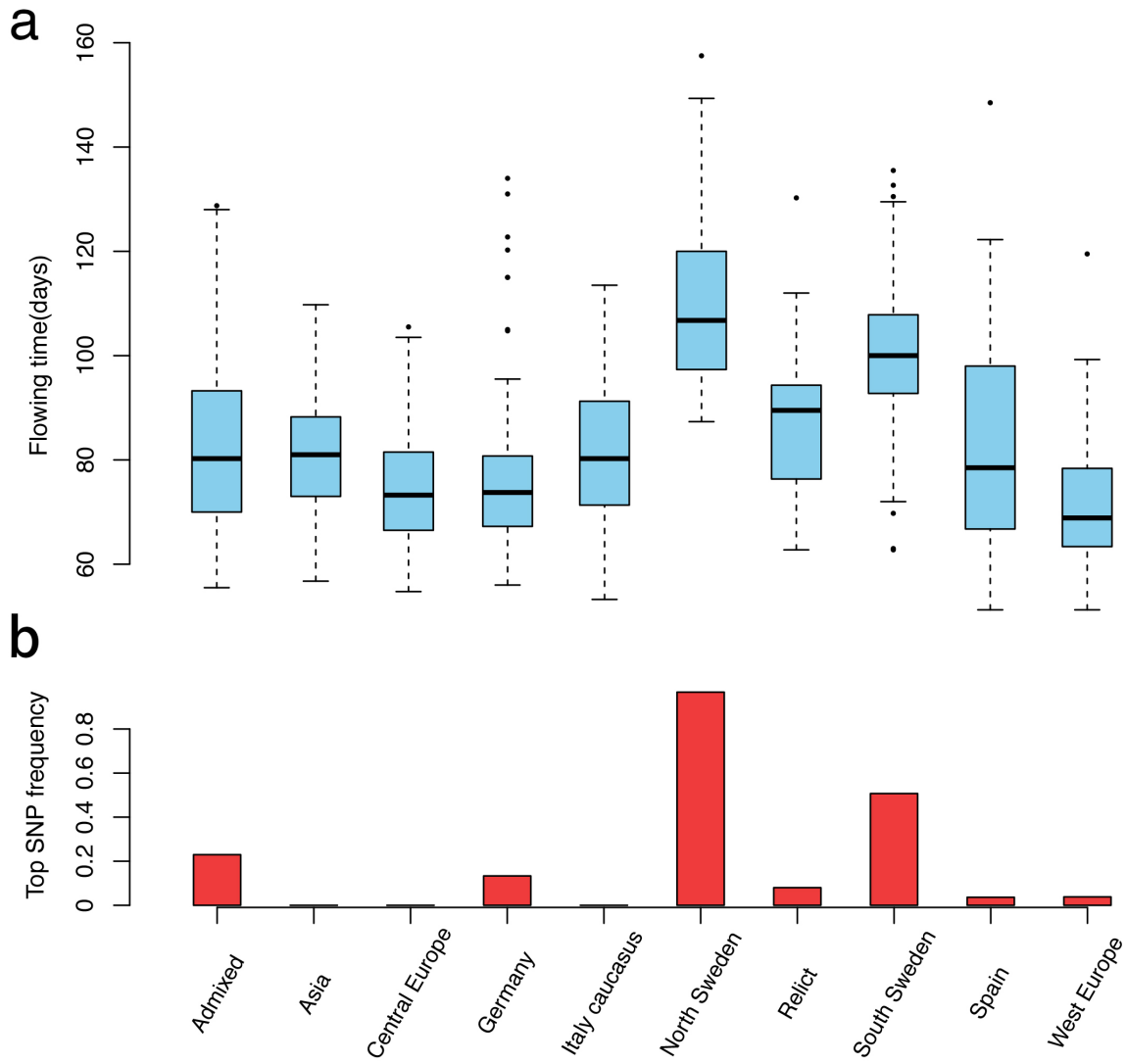
Locus ID	Gene Name	p-value ¹	q-value ²	Replication p-value ³
AT1G17590	<i>NF-YA8</i>	1.6×10^{-7}	2.3×10^{-5}	1.7×10^{-2}
AT5G53360	<i>AT5G53360</i>	5.8×10^{-7}	5.7×10^{-5}	3.2×10^{-4}
AT3G57920	<i>SPL15</i>	7.9×10^{-4}	7.8×10^{-3}	1.7×10^{-2}
AT5G62165	<i>AGL42</i>	1.2×10^{-3}	1.1×10^{-2}	6.3×10^{-3}
AT5G10140	<i>FLC</i>	1.5×10^{-3}	1.3×10^{-2}	5.7×10^{-4}
AT2G45660	<i>AGL20</i>	1.8×10^{-3}	1.4×10^{-2}	1.2×10^{-3}



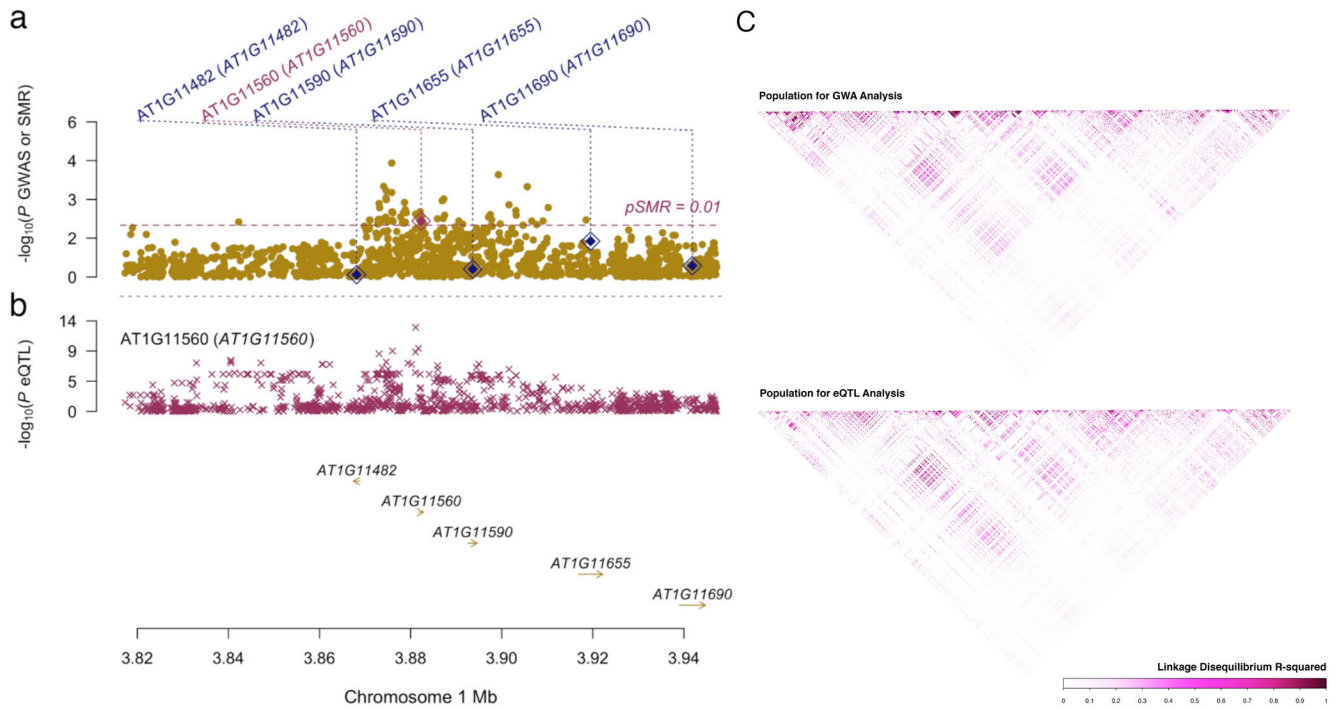
332 Figure 1



333 Figure 2



334 Figure 3



335 Figure 4

336 MATERIALS & METHODS

337 **Genome-wide 250k SNP array genotype data and phenotype data for** 338 **199 natural *Arabidopsis thaliana* accessions**

339 We downloaded a public dataset on collection of 199 natural *Arabidopsis thaliana* inbred
340 lines contains 107 phenotypes and corresponding genotypes (*Atwell et al., 2010*). Those
341 files are publicly available at https://github.com/Gregor-Mendel-Institute/atpolydb/blob/master/miscellaneous_data/phenotype_published_raw.tsv, and https://github.com/Gregor-Mendel-Institute/atpolydb/blob/master/250k_snp_data/call_method_75.tar.gz. 214,051 SNPs were avail-
342 able. After filtering out the variants with minor allele frequency less than 0.10, 173,220
343 SNPs remained.
344
345

346 **Whole genome re-sequencing and RNA-seq data for a population of** 347 **1,135 natural *A. thaliana* accessions**

348 1,135 natural *Arabidopsis thaliana* accessions have been collected and sequenced for the
349 whole genome and transcriptome (*1001 Genomes Consortium, 2016; Kawakatsu et al.,*
350 *2016*). We downloaded this sequencing dataset and removed the accessions with no
351 measured phenotype and SNPs with minor allele frequency below 0.05 and a call-rate
352 below 0.95. The final dataset includes 1001 individuals with 2,222,379 SNPs and measured
353 flowering time at 10°C. To scan for candidate genes, we also downloaded the transcriptome
354 dataset of a subset of this collection ($n = 728$) (*Kawakatsu et al., 2016*). The final eQTL
355 scan dataset contains RNA-seq derived RPKM-values for 24,150 genes in 648 accessions
356 whose phenotypic and genotypic data are both available.

357 **Whole genome re-sequencing derived SNP genotype and RNA-sequencing** 358 **derived transcriptome data for a population of 144 natural *A. thaliana*** 359 **accessions**

360 In an earlier study, Schmitz et al. (*Schmitz et al., 2013*) RNA-sequenced a collection
361 of 144 natural *A. thaliana* accessions. We downloaded this data together with their
362 corresponding whole-genome SNP genotypes available as a part of the 1001 Genomes
363 project (*1001 Genomes Consortium, 2016; Kawakatsu et al., 2016*) to replicate our SMR
364 findings. Following the quality control procedure in (*Zan et al., 2016*), we removed two
365 accessions from the data (Alst_1 and Ws_2) due to missing genotype data and two
366 accessions (Ann_1 and Got_7) due to their low transcript call rate (16,861 and 18,693
367 genes with transcripts as compared to the range of 22,574 to 26,967 for the other the
368 accessions). The final dataset used for eQTL mapping included 1,347,036 SNPs with
369 MAF above 0.05 and call-rate above 0.95 for 140 accessions, and corresponding RNA-seq

370 derived FPKM-values for 33,554 genes.

371 **Single-trait analysis for flowering time trait**

372 For all available traits in this dataset, we first performed a mixed model based single
373 trait genome wide association analysis to generate single trait summaries statistics.
374 Those summaries statistics were used as input for double trait analysis described in the
375 following section. To replicate our signal, we also performed a single trait genome wide
376 association analysis using a collection generated in 1001-genomes project (**1001 Genomes**
377 **Consortium, 2016**). To correct for the population structure in these *A. thaliana* accessions,
378 single-trait genome wide scan was performed based on linear mixed models, using the
379 polygenic and mmscore procedure in GenABEL (**Aulchenko et al., 2007**).

380 **Double-trait genome-wide association analysis**

381 We performed double-trait genome scans using our recently developed multivariate
382 analysis method implemented in the MultiABEL package (**Shen et al., 2017**). The method
383 takes the whole-genome summary statistics to infer shrinkage phenotypic correlation
384 coefficients and conducts MANOVA analysis. The shrinkage phenotypic correlation co-
385 efficient of two traits can be unbiasedly estimated by the correlation of genome-wide
386 Z-scores, which is proportional to the phenotypic correlation on the liability scale, with a
387 shrinkage factor of the square root of sample overlapping proportion. Bivariate p-values
388 are reported. In this way, the bivariate MANOVA analysis is carried out on the liability
389 scale, on partially overlapping sample.

390 **eQTL and SMR analysis**

391 We screened for candidate genes by analyzing the expression data in a subset of the
392 1001-genomes collection containing 140 accessions. Expression values for 19 genes
393 around 20kb up/downstream of the top associated SNP were extracted from (**Schmitz**
394 **et al., 2013**). 14 genes did not pass Kolmogorov-Smirnov test (ks test statistics < 0.8) were
395 filtered out due to potential unreliable measurement mentioned in (**Zan et al., 2016**).
396 The remaining five genes were subsequently passed onto eQTL mapping using qtscore
397 procedure in GenABEL (**Aulchenko et al., 2007**). Output were reformatted according to
398 the description in (**Zhu et al., 2016**). Together with the flowering time single-trait scan
399 results (**1001 Genomes Consortium, 2016**), these were further passed onto SMR analysis
400 scanning for association between individual gene expression and flowering time. The SMR
401 analysis were repeated for 5 top candidates, in an independent gene expression dataset
402 containing 648 accessions (**Kawakatsu et al., 2016**) following the same procedure.

References

- 403
404 **1001 Genomes Consortium.** 1,135 Genomes Reveal the Global Pattern of Polymorphism in
405 *Arabidopsis thaliana*. *Cell*. 2016 Jul; 166(2):481–491.
- 406 **Alexander DH,** Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated
407 individuals. *Genome research*. 2009 Sep; 19(9):1655–1664.
- 408 **Atwell S,** Huang YS, Vilhjálmsson BJ, Willems G, Horton M, Li Y, Meng D, Platt A, Tarone AM, Hu TT,
409 Jiang R, Mulyati NW, Zhang X, Amer MA, Baxter I, Brachi B, Chory J, Dean C, Debieu M, de Meaux
410 J, et al. Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred lines.
411 *Nature*. 2010 Jun; 465(7298):627–631.
- 412 **Aulchenko YS,** Ripke S, Isaacs A, van Duijn CM. GenABEL: an R library for genome-wide association
413 analysis. *Bioinformatics (Oxford, England)*. 2007 May; 23(10):1294–1296.
- 414 **Brachi B,** Faure N, Horton M, Flahauw E, Vazquez A, Nordborg M, Bergelson J, Cuguen J, Roux F.
415 Linkage and association mapping of *Arabidopsis thaliana* flowering time in nature. *PLoS genetics*.
416 2010 May; 6(5):e1000940.
- 417 **Burgess S,** Scott RA, Timpson NJ, Davey Smith G, Thompson SG, EPIC- InterAct Consortium. Using
418 published data in Mendelian randomization: a blueprint for efficient identification of causal risk
419 factors. *European journal of epidemiology*. 2015 Mar; 30:543.
- 420 **Crawley MJ.** *Plant Ecology*. John Wiley & Sons; 2009.
- 421 **Dittmar EL,** Oakley CG, Ågren J, Schemske DW. Flowering time QTL in natural populations of
422 *Arabidopsis thaliana* and implications for their adaptive value. *Molecular ecology*. 2014 Sep;
423 23(17):4291–4303.
- 424 **Hirschhorn JN,** Daly MJ. Genome-wide association studies for common diseases and complex
425 traits. *Nature reviews Genetics*. 2005 Feb; 6(2):95–108.
- 426 **Huang X,** Wei X, Sang T, Zhao Q, Feng Q, Zhao Y, Li C, Zhu C, Lu T, Zhang Z, Li M, Fan D, Guo Y,
427 Wang A, Wang L, Deng L, Li W, Lu Y, Weng Q, Liu K, et al. Genome-wide association studies of 14
428 agronomic traits in rice landraces. *Nature genetics*. 2010 Nov; 42(11):961–967.
- 429 **Kawakatsu T,** Huang SSC, Jupe F, Sasaki E, Schmitz RJ, Urich MA, Castanon R, Nery JR, Barragan
430 C, He Y, Chen H, Dubin M, Lee CR, Wang C, Bemm F, Becker C, O'Neil R, O'Malley RC, Quarless
431 DX, 1001 Genomes Consortium, et al. Epigenomic Diversity in a Global Collection of *Arabidopsis*
432 *thaliana* Accessions. *Cell*. 2016 Jul; 166(2):492–505.
- 433 **Korte A,** Farlow A. The advantages and limitations of trait analysis with GWAS: a review. *Plant*
434 *methods*. 2013; 9:29.
- 435 **Schmitz RJ,** Schultz MD, Urich MA, Nery JR, Pelizzola M, Libiger O, Alix A, McCosh RB, Chen H, Schork
436 NJ, Ecker JR. Patterns of population epigenomic diversity. *Nature*. 2013 Mar; 495(7440):193–198.

- 437 **Shen X**, De Jonge J, Forsberg SKG, Pettersson ME, Sheng Z, Hennig L, Carlborg O. Natural CMT2
438 variation is associated with genome-wide methylation changes and temperature seasonality.
439 PLoS genetics. 2014 Dec; 10(12):e1004842.
- 440 **Shen X**, Klarić L, Sharapov S, Mangino M, Ning Z, Wu D, Trbojević-Akmačić I, Pučić-Baković M, Rudan
441 I, Polasek O, Hayward C, Spector TD, Wilson JF, Lauc G, Aulchenko YS. Multivariate discovery
442 and replication of five novel loci associated with Immunoglobulin G N-glycosylation. Nature
443 communications. 2017 Sep; 8(1):447.
- 444 **Shen X**, Pettersson M, Rönnegård L, Carlborg O. Inheritance beyond plain heritability: variance-
445 controlling genes in *Arabidopsis thaliana*. PLoS genetics. 2012; 8(8):e1002839.
- 446 **Visscher PM**, Yang J. A plethora of pleiotropy across complex traits. Nature genetics. 2016 Jun;
447 48(7):707–708.
- 448 **Wang B**, Li Z, Xu W, Feng X, Wan Q, Zan Y, Sheng S, Shen X. Bivariate genomic analysis identifies a
449 hidden locus associated with bacteria hypersensitive response in *Arabidopsis thaliana*. Scientific
450 reports. 2017 Mar; 7:45281.
- 451 **Wellenreuther M**, Hansson B. Detecting Polygenic Evolution: Problems, Pitfalls, and Promises.
452 Trends in genetics : TIG. 2016 Mar; 32(3):155–164.
- 453 **Yang J**, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, Madden PA, Heath AC, Martin
454 NG, Montgomery GW, Goddard ME, Visscher PM. Common SNPs explain a large proportion of
455 the heritability for human height. Nature genetics. 2010 Jul; 42(7):565–569.
- 456 **Yang J**, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and pitfalls in the application of
457 mixed-model association methods. Nature genetics. 2014 Feb; 46(2):100–106.
- 458 **Zan Y**, Shen X, Forsberg SKG, Carlborg O. Genetic Regulation of Transcriptional Variation in Natural
459 *Arabidopsis thaliana* Accessions. G3 (Bethesda, Md). 2016; 6(8):2319–2328.
- 460 **Zhu Z**, Zhang F, Hu H, Bakshi A, Robinson MR, Powell JE, Montgomery GW, Goddard ME, Wray NR,
461 Visscher PM, Yang J. Integration of summary data from GWAS and eQTL studies predicts complex
462 trait gene targets. Nature genetics. 2016 May; 48(5):481–487.

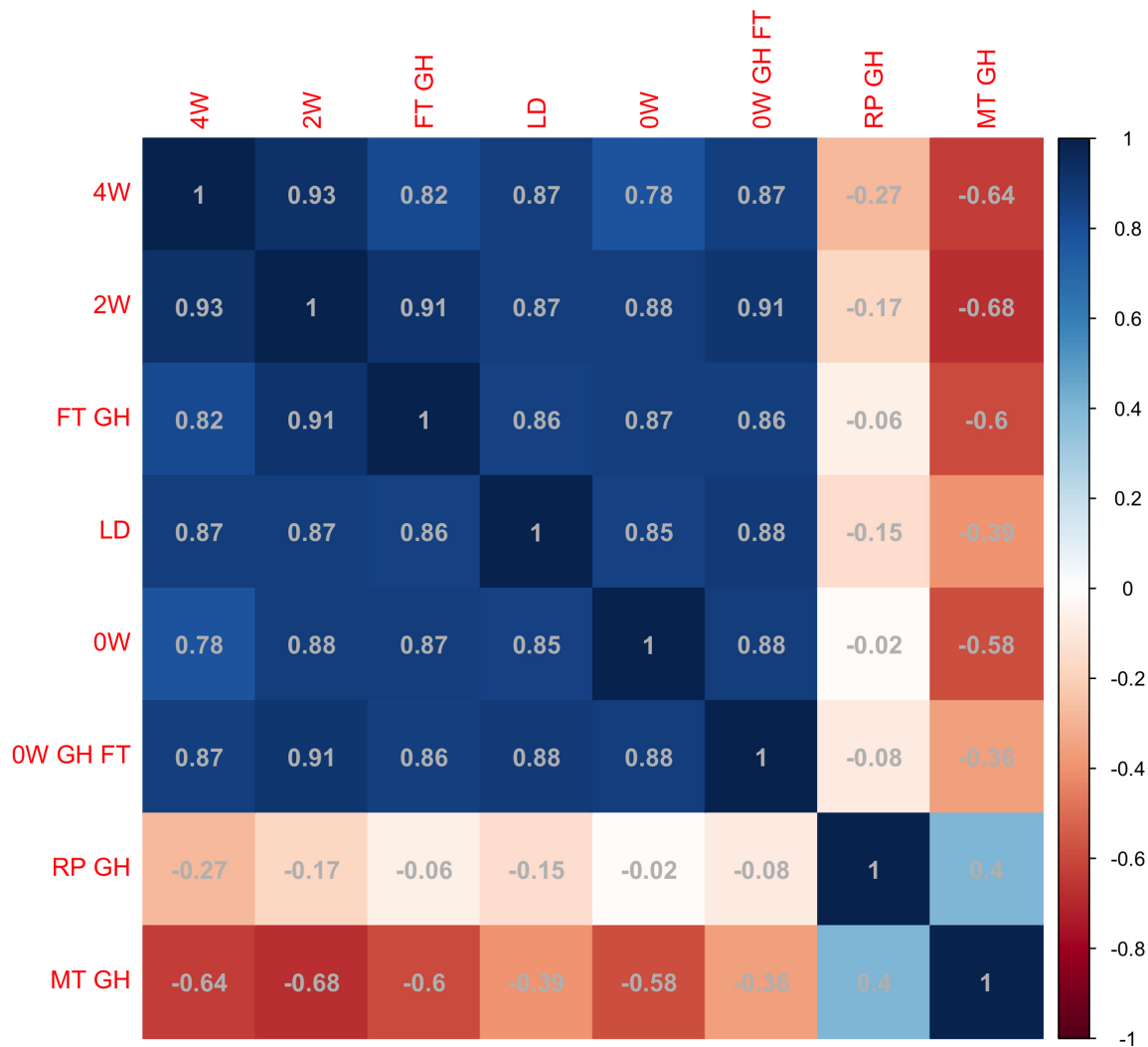
463 Supplementary Table 1: Phenotypes included in the bivariate analyses.

464 Details about phenotyping can be referred to *Atwell et al. (2010)*.

Phenotype	Description	Number of Accessions
LD	Days to flowering time (FT) under Long Day (LD)	167
LDV	Days to flowering time (FT) under Long Day (LD) (5 wks vernalization)	168
SD	Days to flowering time (FT) under Short Day (SD)	162
SDV	Days to flowering time (FT) under Short Day (SD) (5 wks vernalization)	159
0W	Days to FT under LD without vernalization	137
2W	Days to FT under LD with 2wks vernalization	152
4W	Days to FT under LD with 4wks vernalization	119
8W	Days to FT under LD with 8wks vernalization	155
FLC	FLC gene expression	167
FRI	FRI gene expression	164
FT10	Flowering time (FT), 10°C	194
FT16	Flowering time (FT), 16°C	193
FT22	Flowering time (FT), 22°C	193
LN10	leaf number at flowering time (LN), 10°C	177
LN16	leaf number at flowering time (LN), 16°C	176
LN22	leaf number at flowering time (LN), 22°C	176
8W GH FT	Days to FT with 8 wks vernalization	162
8W GH LN	LN at FT with 8 wks vernalization	163
0W GH FT	Days to FT without vernalization	153
0W GH LN	LN at FT without vernalization	135
FT Field	Days to flowering of plants grown in the field	180
FT Diameter Field	Plant diameter at flowering (field)	180
FT GH	Days to flowering (greenhouse)	166
LES	Presence or absence of lesioning	95
YEL	Presence or absence of yellowing	95
LY	Presence or absence of either lesioning or yellowing	95
FW	Fresh weight of plants	95
DW	Dry weight of plants	95
Chlorosis 10	Visual chlorosis presence, 10°C	177
Chlorosis 16	Visual chlorosis presence, 16°C	176
Chlorosis 22	Visual chlorosis presence, 22°C	176

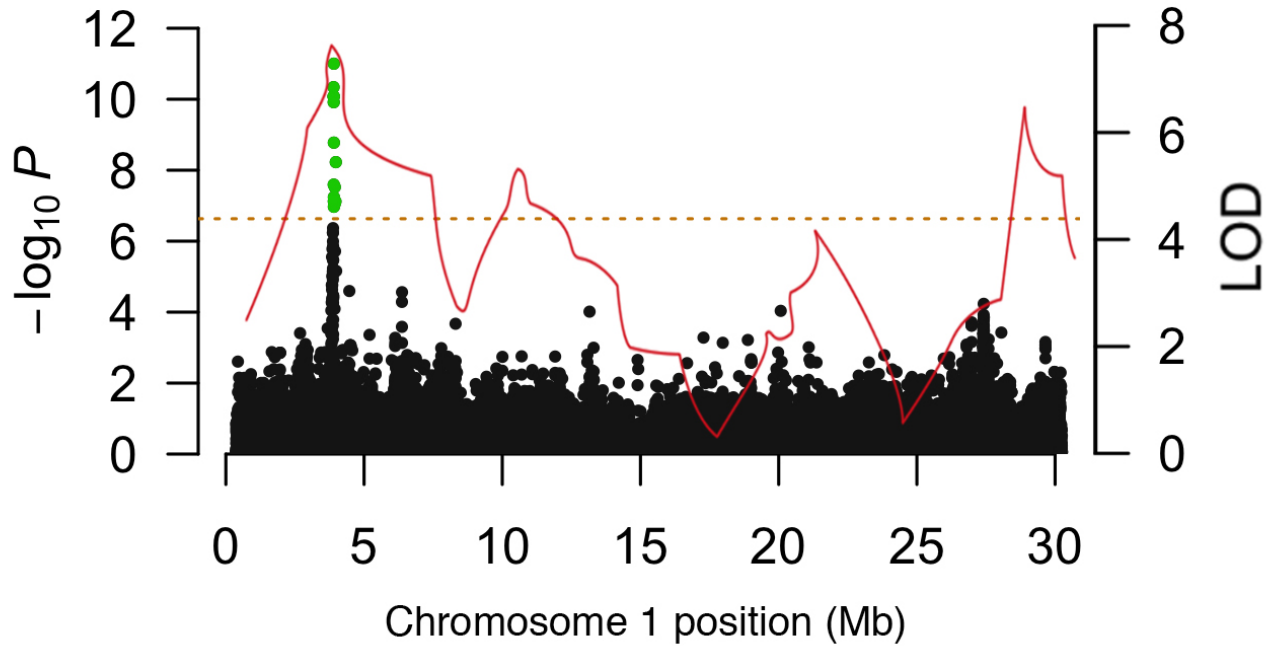
Anthocyanin 10	Visual anthocyanin presence, 10°C	177
Anthocyanin 16	Visual anthocyanin presence, 16°C	176
Anthocyanin 22	Visual anthocyanin presence, 22°C	177
Seed Dormancy	Seed dormancy level	83
Germ 10	Days to germination, 10°C	177
Germ 16	Days to germination, 16°C	176
Germ 22	Days to germination, 22°C	177
Seedling Growth	Seedling growth rate	100
Vern Growth	Vegetative growth rate during vernalization	110
After Vern Growth	Vegetative growth rate after vernalization	110
Secondary Dormancy	Decrease in germination rate after prolonged exposure to cold temperature	93
Germ in dark	Germination in the dark	93
DSDS50	Duration of seed dry storage required for 50% of the seeds to germinate	109
Seed bank 133-91	Non-monotonous dynamic of dormancy release	110
Storage 7 days	Primary dormancy, 7 days dry storage	110
Storage 28 days	Primary dormancy, 28 days dry storage	110
Storage 56 days	Primary dormancy, 56 days dry storage	110
Hypocotyl length	⁴⁶⁶ Hypocotyl length	89
Width 10	Plant diameter, 10°C	176
Width 16	Plant diameter, 16°C	175
Width 22	Plant diameter, 22°C	175
Leaf serr 10	Level of leaf serration, 10°C	174
Leaf serr 16	Level of leaf serration, 16°C	176
Leaf serr 22	Level of leaf serration, 22°C	176
Leaf roll 10	Leaf roll presence, 10°C	177
Leaf roll 16	Leaf roll presence, 16°C	176
Leaf roll 22	Leaf roll presence, 22°C	176
Rosette Erect 22	Presence of rosette erectness, 22°C	176
Silique 16	Silique length, 16°C	95
Silique 22	Silique length, 22°C	95
FT Duration GH	Flowering period duration	147
LC Duration GH	Life cycle period	147
LFS GH	Last flower senescence	148
MT GH	Maturation period	147
RP GH	Reproduction period	147

467



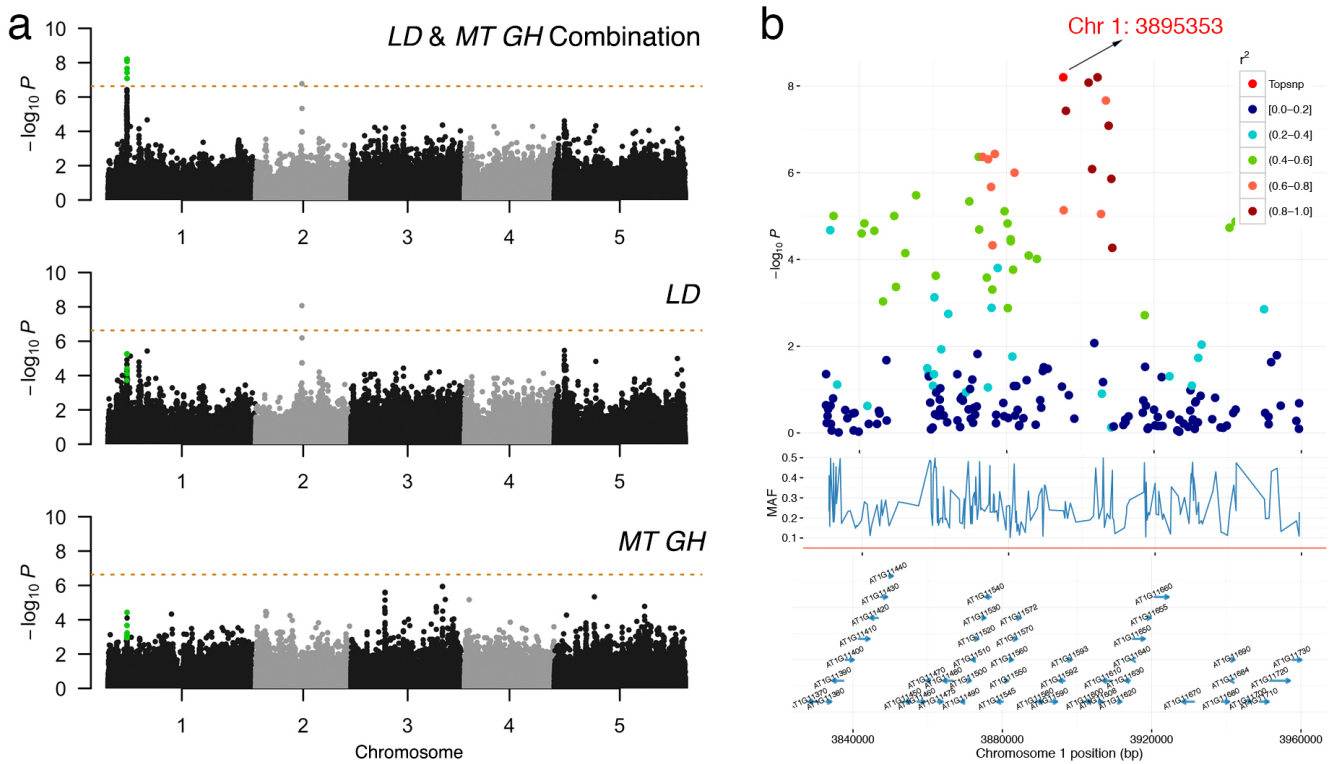
468 Supplementary Figure 1: Phenotypic correlations among flowering time
 469 related traits, maturation period and reproduction period phenotypes.

470 The flowering time related traits are: 4W: Days to flowering time (FT) under long day (LD)
 471 with vernalized for 4 wks at 5°C, 8hrs daylight; 2W: Days to flowering time (FT) under
 472 long day (LD) with vernalized for 2 wks at 5°C, 8hrs daylight; FT GH: Days to flowering
 473 (greenhouse); LD: Days to flowering time (FT) under Long Day (LD); OW: Days to flowering
 474 time (FT) under Long Day (LD) without vernalization; OW GH FT: Days to flowering time
 475 (FT).



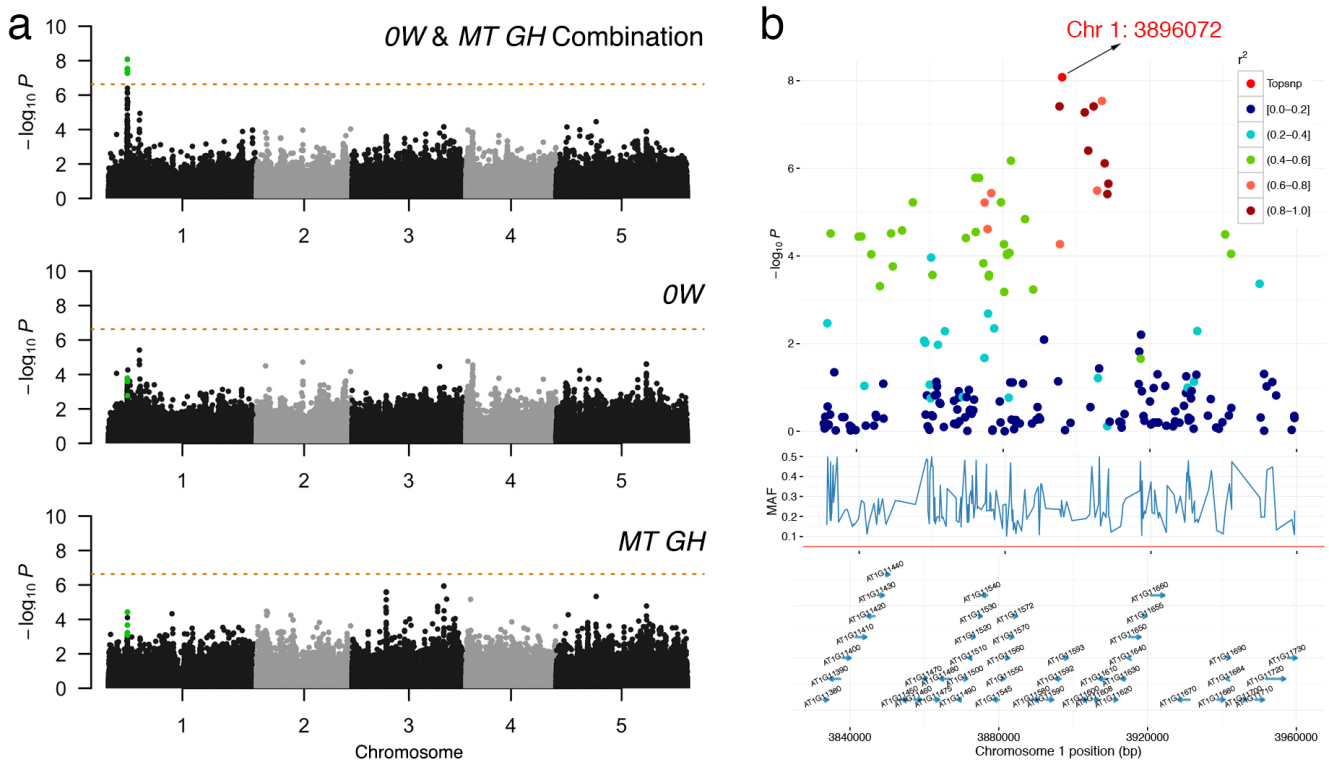
476 Supplementary Figure 2: Overlapping between QTL mapping and double-
477 trait GWAS result.

478 The curve shows stepwise LOD profiles in chromosome 1 that are generated from a QTL
479 mapping study using a cross between Italy and Sweden population analyzed by *Dittmar*
480 *et al. (2014)* (reproduced by depicting the curvature of Figure 3a therein). The Manhattan
481 plot shows chromosome 1 signal in our bivariate analysis.



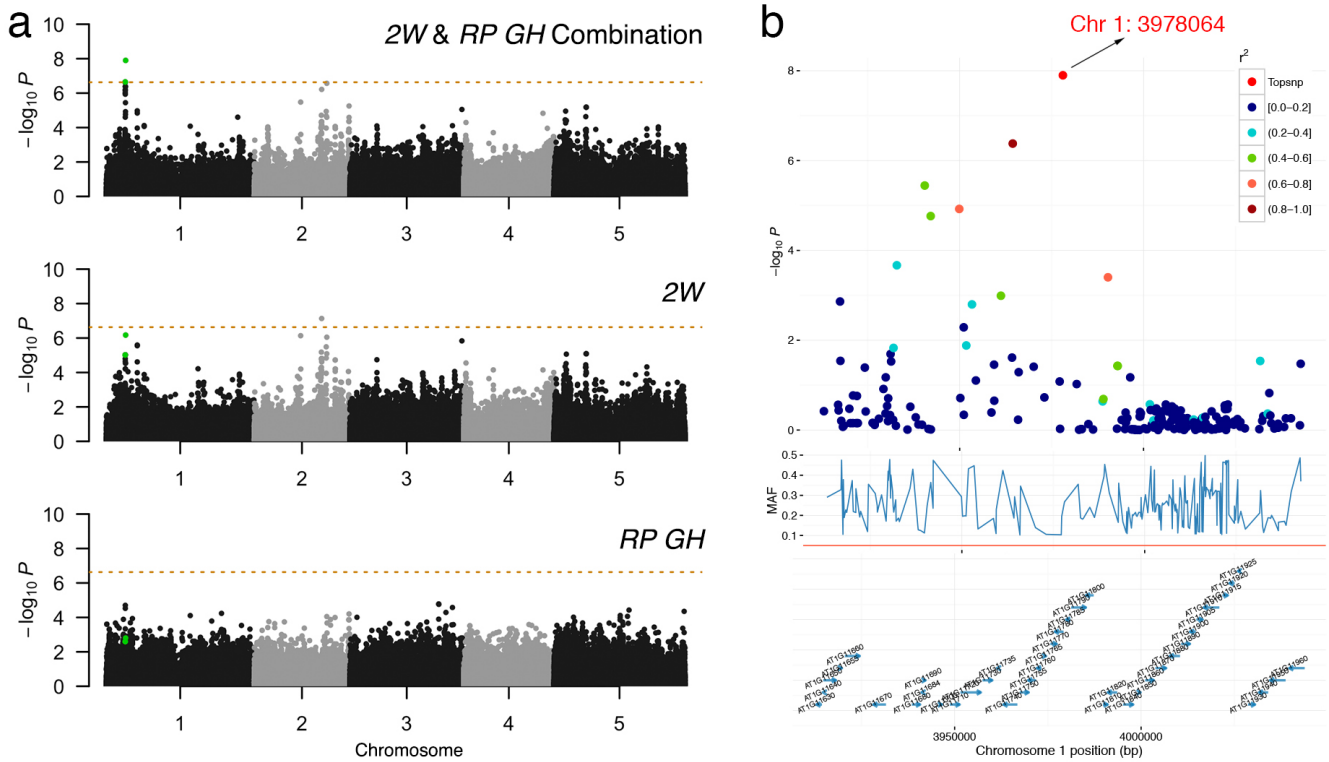
482 Supplementary Figure 3: Bivariate genome-wide association analysis of
483 two developmental trait, LD: Days to flowering time (FT) under Long Day
484 (LD), MT GH: Maturation period.

485 (a) Manhattan plots comparison of bivariate and univariate analysis results, where the
486 novel variants only discoverable when combining two phenotypes are shown in green.
487 The horizontal dashed line represents a 5% Bonferroni-corrected genome-wide significant
488 threshold. (b) Zooming in the novel locus detected using bivariate analysis. r : linkage
489 disequilibrium measured as correlation coefficient between the top variant and each
490 variant in the region.



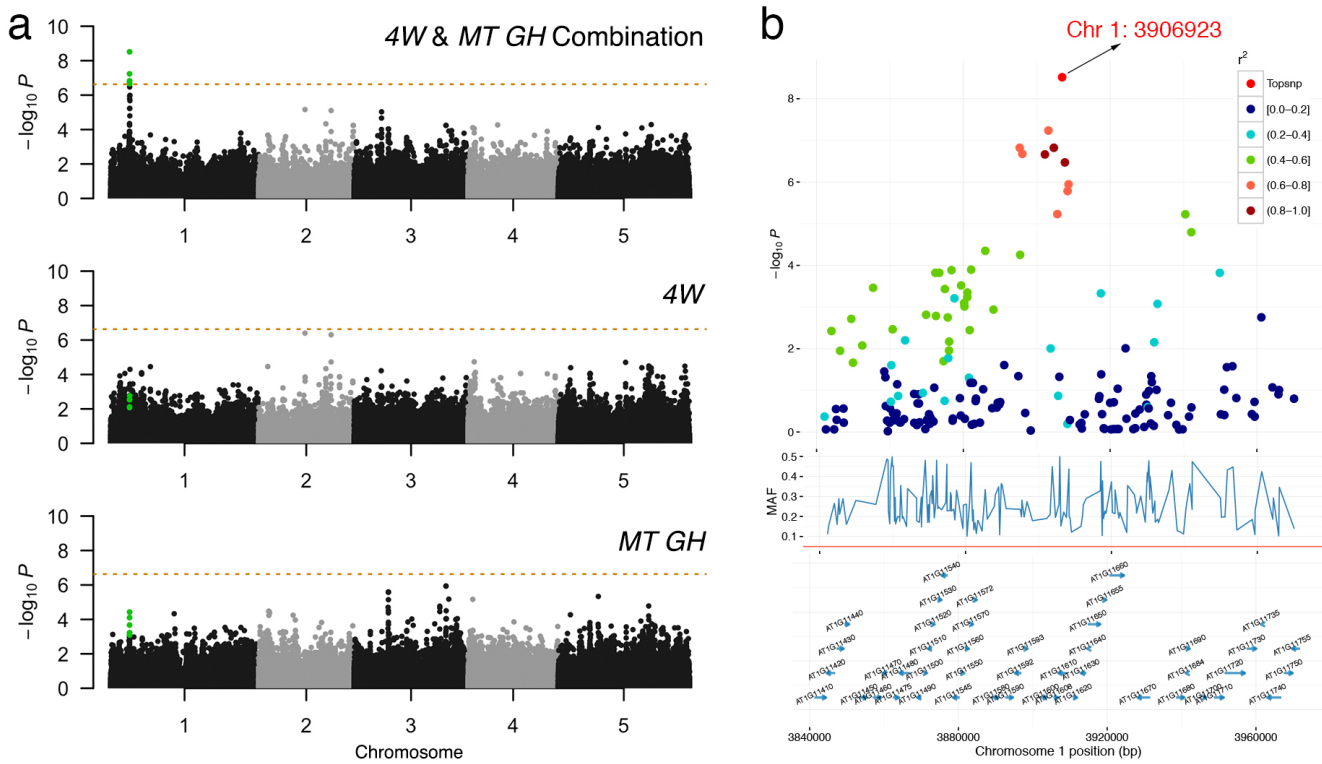
491 Supplementary Figure 4: Bivariate genome-wide association analysis of
492 two developmental trait, OW: Days to flowering time (FT) under Long Day
493 (LD) without vernalization, MT GH: Maturation period.

494 (a) Manhattan plots comparison of bivariate and univariate analysis results, where the
495 novel variants only discoverable when combining two phenotypes are shown in green.
496 The horizontal dashed line represents a 5% Bonferroni-corrected genome-wide significant
497 threshold. (b) Zooming in the novel locus detected using bivariate analysis. r : linkage
498 disequilibrium measured as correlation coefficient between the top variant and each
499 variant in the region.



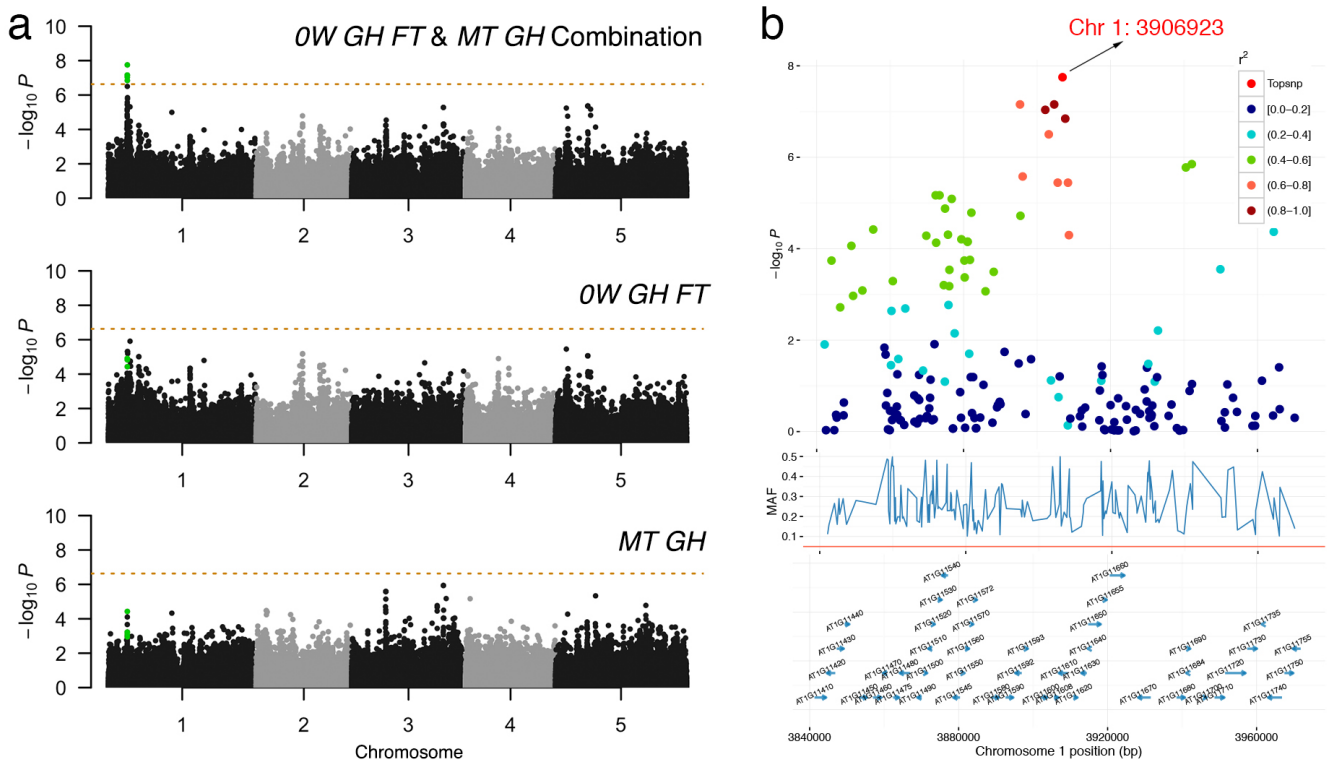
500 Supplementary Figure 5: Bivariate genome-wide association analysis of
501 two developmental trait, 2W: Days to flowering time (FT) under long day
502 (LD) with vernalized for 2 wks at 5°C, 8 hrs daylight, RP GH: Reproduction
503 period.

504 (a) Manhattan plots comparison of bivariate and univariate analysis results, where the
505 novel variants only discoverable when combining two phenotypes are shown in green.
506 The horizontal dashed line represents a 5% Bonferroni-corrected genome-wide significant
507 threshold. (b) Zooming in the novel locus detected using bivariate analysis. r : linkage
508 disequilibrium measured as correlation coefficient between the top variant and each
509 variant in the region.



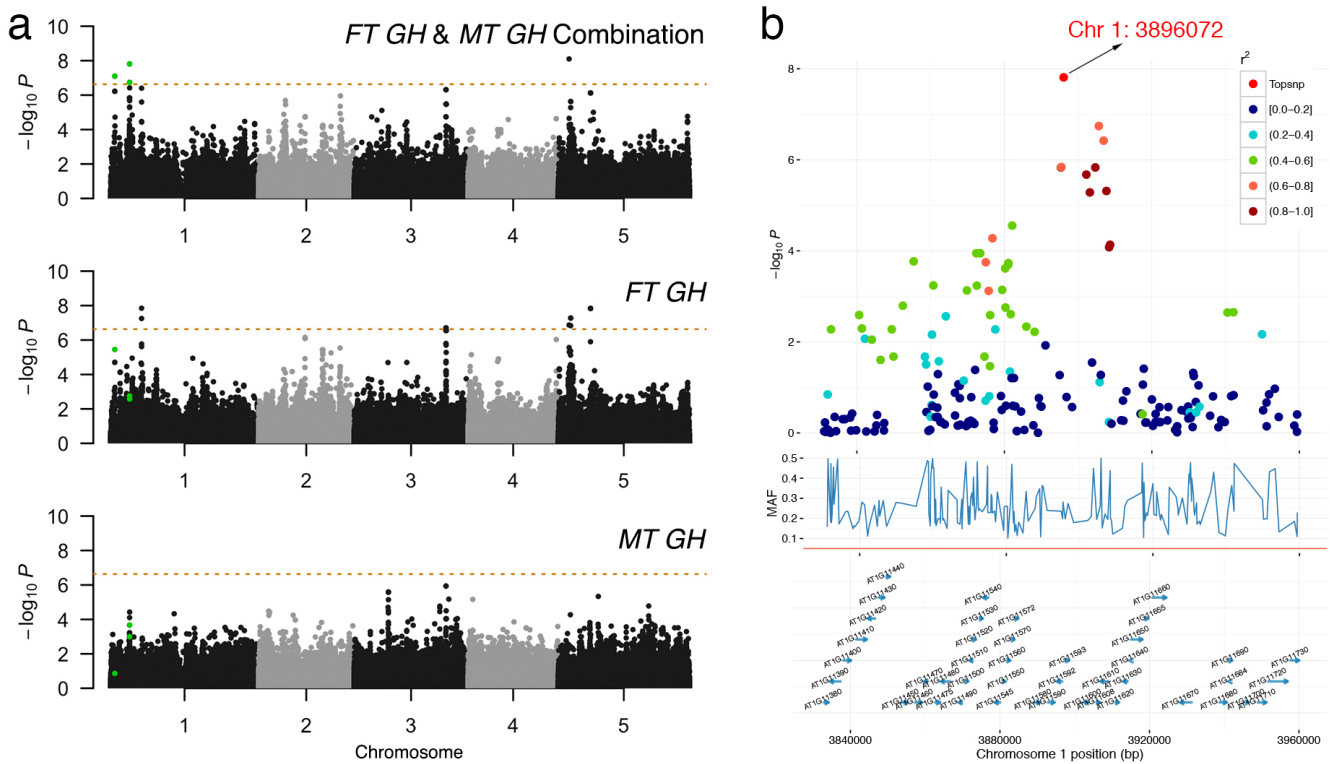
510 Supplementary Figure 6: Bivariate genome-wide association analysis of
511 two developmental trait, 4W: Days to flowering time (FT) under long day
512 (LD) with vernalized for 4 wks at 5°C, 8hrs daylight, MT GH: Maturation
513 period.

514 (a) Manhattan plots comparison of bivariate and univariate analysis results, where the
515 novel variants only discoverable when combining two phenotypes are shown in green.
516 The horizontal dashed line represents a 5% Bonferroni-corrected genome-wide significant
517 threshold. (b) Zooming in the novel locus detected using bivariate analysis. r : linkage
518 disequilibrium measured as correlation coefficient between the top variant and each
519 variant in the region.



520 Supplementary Figure 7: Bivariate genome-wide association analysis of
521 two developmental trait, 0W GH FT: Days to flowering: time (FT), MT GH:
522 Maturation period.

523 (a) Manhattan plots comparison of bivariate and univariate analysis results, where the
524 novel variants only discoverable when combining two phenotypes are shown in green.
525 The horizontal dashed line represents a 5% Bonferroni-corrected genome-wide significant
526 threshold. (b) Zooming in the novel locus detected using bivariate analysis. r : linkage
527 disequilibrium measured as correlation coefficient between the top variant and each
528 variant in the region.



529 Supplementary Figure 8: Bivariate genome-wide association analysis of
530 two developmental trait, FT GH: Days to flowering (greenhouse), MT GH:
531 Maturation period.

532 (a) Manhattan plots comparison of bivariate and univariate analysis results, where the
533 novel variants only discoverable when combining two phenotypes are shown in green.
534 The horizontal dashed line represents a 5% Bonferroni-corrected genome-wide significant
535 threshold. (b) Zooming in the novel locus detected using bivariate analysis. r : linkage
536 disequilibrium measured as correlation coefficient between the top variant and each
537 variant in the region.