1  **ddRAD-seq-derived SNPs reveal novel association signatures for fruit-related traits in**

2  **peach**

3  **Running title: ddRAD-seq approach infers association signals in peach**

4  Najla Ksouri[1], Gerardo Sánchez[2], Carolina Font i Forcada[3], Bruno Contreras-Moreira[4*],

5  Yolanda Gogorcena[1*]

6  [1]Group of Genomics of Fruit Trees and Grapevine, Department of Pomology, Estación

7  Experimental de Aula Dei-Consejo Superior de Investigaciones Científicas, Avenida de

8  Montañana 1005, E50059 Zaragoza, Spain.

9  [2]Biotechnology Lab, Estación Experimental Agropecuaria (EEA) San Pedro, INTA, Ruta

10  N°9 km 170, B2930 San Pedro, Argentina.

11  [3]Department of Pomology, Estación Experimental de Aula Dei-Consejo Superior de

12  Investigaciones Científicas, Avenida de Montañana 1005, E50059 Zaragoza, Spain.

13  [4]Laboratory of Computational and Structural Biology, Department of Genetics and Plant

14  Production, Estación Experimental de Aula Dei-Consejo Superior de Investigaciones

15  Científicas, Avenida de Montañana 1005, E50059 Zaragoza, Spain.

16  *Senior authors

17  Emails:

18  Najla Ksouri: nksouri@eead.csic.es Tel: +34 976 716132

19  Gerardo Sánchez: sanchez.gerardo@inta.gob.ar Tel: +54 93329542117

20  Carolina Font i Forcada: carolffont@gmail.com Tel: +34 673329049

21  Bruno Contreras-Moreira*: bcontreras@eead.csic.es Tel: +34 976716089

22  Yolanda Gogorcena*: aoiz@eead.csic.es Tel: +34 976 716133

**Abstract**

Breeding for new peach cultivars with enhanced traits is a prime target in breeding programs. In this study, we used a discovery panel of 90 peach accessions in order to dissect the genetic architecture of 16 fruit-related traits. ddRAD-seq genotyping and the intersection between three variant callers yielded 13,045 high-confidence SNPs. These markers were subjected to an exhaustive association analysis by testing up to seven GWAS models. Blink was selected as the most adjusted, simultaneously balancing false positive and negative associations. Totally, we identified 16 association signals for six traits showing high broad-sense heritability: harvest date, fruit weight, flesh firmness, contents of flavonoids, anthocyanins and sorbitol. By assessing the allelic effect of significant markers on phenotypic attributes, nine SNP alleles were denoted favorable. A promising marker (SNC_034014.1_7012470) was found to be simultaneously associated with harvest date and fruit firmness conferring a positive allelic effect on both traits. We anticipate that this marker could be used to improve firmness in late harvested cultivars. Candidate causal genes were shortlisted when fulfilling the following criteria: i) position within the linkage disequilibrium block, ii) functional annotation and iii) expression pattern. A bibliographic review of previously reported QTLs mapping nearby the associated markers allowed us to benchmark the accuracy of our approach. Despite the moderate germplasm size, ddRAD-seq allowed us to produce an accurate representation of peach's genome resulting in SNP markers suitable for empirical association studies. Together with candidate genes, they lay the foundation for further genetic dissection of peach key traits.

**Key words**: lead SNP, prime candidate genes, haplotype blocks, fruit-related traits, linkage disequilibrium, *Prunus persica*

**Background**

Peach is one of the most economically valued fleshy fruits worldwide (FAO, http://faostat.fao.org). The advances in the peach industry largely rely on fruit quality improvement in response to the market and consumers' demands. The term quality may include all agronomical aspects and chemical compounds such as fruit size, firmness, sugar and acid concentration, etc. Some of those characteristics are thought to be monogenic,

52    controlled by a single gene (fruit shape, hairiness, flesh color, texture)[1–3] while others are

53    polygenic, such as sugar content, fruit firmness, antioxidant concentration[4].

54    Breeding for polygenic quantitative traits is far from being a straightforward task. Thus,

55    insights on genetic drivers controlling these traits and their inheritance are required to

56    bridge the phenotype-genotype gap[3,5]. For instance, the development of molecular markers

57    linked to desirable traits would considerably speed up the selection of superior plant

58    varieties through marker-assisted selection (MAS)[6]. Genome-wide association studies

59    (GWAS) have also revolutionized the breeding process by detecting the genetic loci

60    underlying trait variations at a relatively high resolution. This approach has been

61    successfully applied in many breeding programs. For instance, GWAS have provided

62    insight into fruit-related traits such as skin color in apple[7] and fruit firmness in sweet

63    cherry[8]. The power and prediction accuracy of GWAS critically depend on various

64    considerations, including phenotypic data quality, experimental sample size, linkage

65    disequilibrium (LD) between genetic variants and population structure. If not adjusted

66    properly, these factors may lead to spurious associations as well as masking the true ones.

67    Another key factor while performing GWAS is the density and chromosome distribution of

68    markers/SNPs along the reference genome.

69    Generally, genotyping methods fall into three categories; whole genome resequencing,

70    reduced representation sequencing, and SNP arrays□□[9]. Whole genome resequencing

71    returns the highest number of SNP calls if sequencing depth is sufficient, which is

72    expensive for large genomes. For this reason, SNP arrays are widely used, reducing the cost

73    and enabling the detection of thousands of SNPs in a single assay[9]. In peach, commercially

74    available arrays IPSC peach 9K[10] and IPSC peach 18K[11] have been used to explore the

75    genetic diversity and to assist the breeding process[1,12]. Despite their utility, the major

76    drawback of SNP genotyping arrays consists in their ascertainment bias[13]. In other words,

77    they narrow the discovery of novel variants other than those detected in the discovery panel

78    and used to build the respective array. This might distort subsequent genetic inferences.

79    Additionally, efficient SNP probes require a well-assembled reference genome and their

80    design and further optimization can be time consuming.

With the massive progress of high-throughput technologies, reduced representation sequencing such as restriction-associated DNA (RAD) sequencing and its derivative (ddRADseq) emerged to overcome both cost and ascertainment bias[14]. Double digest restriction-site associated DNA (ddRADseq) relies on the use of a pair of restriction enzymes to limit the sequencing effort to a subset of evenly distributed loci in the genome[14]. Moreover, by picking the best enzyme combination, repetitive DNA can be less targeted, thereby reducing the computational burden associated with aligning genomes with highly repetitive segments.

Unlike other genotyping methods, prior genomic information is strictly not required for ddRADseq[14]. Nevertheless, as shown in this work, it is most powerful when combined with a reference genome sequence. From a technical standpoint, a common shortcoming of ddRADseq is the high rate of missing calls which can be straightforwardly handled through genotype imputation.

Herein, we report the application of ddRADseq genotyping to identify high confidence SNPs in a discovery panel of 90 *Prunus persica* accessions. Consequently, GWAS was carried out to identify genomic loci associated with 16 fruit traits. To optimize the analysis and to overcome the limitations arising from the size of our peach germplasm, we considered the following aspects: 1) peach accessions were geographically distant in order to maximize the genetic variance, 2) SNPs were called using three variant detectors (BCFtools, Freebayes and GATK) and only those resulting from the intersection were retained for subsequent analysis, and 3) several statistical models were assessed to control the confounding effects.

Genotype-to-phenotype associations for agronomic and fruit-related traits have been widely tested in peach using different genotyping methods like SSRs[15], 9K SNP array[1,4,16], 18K SNP array[3,12] and high-throughput resequencing technology[17]. However, to the best of our knowledge this is the first report characterizing the genetic architecture of peach traits using ddRADseq-derived SNPs. In this study, we propose best practices for GWAS analysis mainly relying on a comparative approach for SNPs calling and statistical model assessment. Therefore, we demonstrate the utility of ddRAD-based genotyping in unveiling desirable alleles and genomic regions putatively responsible for trait variation. By

111    contrasting our findings with those previously reported using the peach 9K SNP array[16] we

112    confirm the accuracy of our approach.


113    **Results**


114    **Phenotypic analysis and heritability**

115    Broad sense heritability was estimated over three consecutive years and the results denote

116    that most of the traits were highly heritable (**Figure 1.A**). Hence, their phenotypic

117    variability among the individuals was mainly driven by the genetic effects. However,

118    contents of glucose, fructose, sucrose and total sugars (TS) were found to be lowly heritable

119    traits ($H^2 < 0.5$), denoting that their variability may be mostly due to the environmental

120    factors. These traits were therefore left out of the association analyses. Furthermore, normal

121    distribution fit tests conducted on averaged phenotypic measures, revealed that six out of 16

122    traits were found to be normally distributed (flesh firmness, soluble solids content (SSC),

123    ripening index, vitamin C, relative antioxidant capacity (RAC) and glucose). Source code,

124    documentation     and     detailed     results     can     be     accessed     at

125    https://github.com/najlaksouri/GWAS-Workflow. The remaining ones, skewed either

126    positively or negatively, were transformed accordingly. Likewise, the phenotypic

127    correlation was estimated and significant interactions between agronomical and fruit quality

128    traits were observed (**Figure 1.B**). For instance, harvest date (HvD) had the highest

129    heritability estimates ($H^2=0.94$) and exhibited strong positive correlations with flesh

130    firmness, sugar contents measured as (SSC, TS and sorbitol) and antioxidant activity

131    measured as (RAC, flavonoids and phenols). As expected, moderate positive interaction

132    was also reported between the HvD and fruit weight as well as between total and individual

133    sugars. Moreover, a strong positive correlation was also observed between total phenolics

134    and flavonoids. Indeed, flavonoids are the largest group of naturally occurring phenolic

135    compounds in plants. Both compounds showed a significant positive interaction with

136    (RAC) suggesting that they could be used as a good indicator of antioxidant properties in

137    peaches.


5

**SNP genotyping**

To construct an informative SNP panel, polymorphic sites were called in individual sample mode using three different algorithms. Raw calls were subjected to standardized quality thresholds in order to mitigate the effect of sequencing and/or alignment flaws. Post-filtered calls from each pipeline were merged together into multi-samples format (**Table 1**). According to our results, GATK-HaplotypeCaller (HC) outperformed both Freebayes and BCFTools in terms of computational time and sensitivity yielding a total of 233,535 SNP calls (see repository https://github.com/najlaksouri/GWAS-Workflow). Freebayes ranked second, followed by BCFTools, with 166,080 and 148,998 SNPs, respectively. For a robust variant detection, the intersection between multi-sample sets was computed. About 32% of SNPs were found to be commonly shared by the above-stated tools. Multi-allelic and scaffold variants were excluded and additional filters (missing call rate and MAF) were applied (**Table 1**). Finally, a set of 13,045 SNPs was kept for subsequent analysis.

Using VEP tool, polymorphic sites were found to be distributed along upstream (21%), downstream (9%), intronic (26%) and intergenic (8%) regions (**Figure S1**). Low proportions of SNPs were tagged as 3' UTR and 5' UTR variants. Within coding regions, 11% of SNPs were defined as synonymous while 13% were annotated as missense variants.

**SNP distribution and LD decay**

The distribution of polymorphic sites was calculated within adjacent windows of 1 Mbp and provided a genome-wide coverage estimate along the eight peach chromosomes. As illustrated in **Figure 2.A**, markers were unevenly partitioned throughout the genome with the highest number of mapped SNPs on chromosome 2 (4,440) and the lowest on chromosome 5 (1,768). Interestingly, SNPs accumulated within the short arms of chromosomes 2 and 4. In contrast, large gaps were observed towards the telomere of the long arm of chromosome 2. Similarly, several blank regions were located along chromosome 1. Gaps highlighted with asterisks correspond to predicted centromeric regions[18].

To determine the extent of LD decay in the diversity panel, we estimated the pairwise LD coefficient ($r^2$) at chromosomic level. LD decay was estimated for each chromosome by estimating the intersection of $r^2=0.25$ with the physical distance (**Figure 2.B** and **Figure**

168    **S2**). We found that LD dropped at short distance, ranging from 250 to 500 kbp along all

169    chromosomes, with the exception of chromosome 5 (ca. 4.7 Mbp). After LD pruning, a

170    total of 1,959 unlinked SNPs was kept for population structure and kinship estimations.

171    **Population structure**

172    PCA analysis separated the germplasm panel into 4 sub-populations based on the genetic

173    origin (landrace vs modern breeding line) and fruit type (peach vs nectarine) (**Figure S3**).

174    Clade 1 on the top left corner, grouped exclusively modern breeding lines of peach and

175    nectarine. This group seems to be driven by the geographical origin as most of the

176    accessions were originated from North America (**Table S1**). Clade 2 represents a diverse

177    genetic entity gathering both landrace and breed peach varieties. Genotypes within this

178    clade were originated from Spain and North America suggesting the presence of higher

179    admixture that could arise due to the exchange of the germplasm material. In contrast,

180    clades 3 and 4 contained only landrace peach accessions mostly from different regions of

181    Spain, Europe and South Africa. A neighbor joining (NJ) tree also identified four clear

182    clusters, as illustrated in **Figure S4.** Comparable results were obtained from

183    fastSTRUCTURE and are provided in the GitHub repository.

184    **Critical evaluation of GWAS models**

185    Genome wide association studies may be susceptible to bias in the presence of

186    measurement errors. False positive and negative associations arising from population

187    structure or/and family relatedness may lead to erroneous conclusions. The examination of

188    Q-Q plots can be used as a straight visual inspection to determine the appropriate statistical

189    method controlling the confounding effects. In fact, Q-Q plots illustrate the distribution of

190    markers under the null hypothesis, by plotting the observed $-\log_{10}$ $P$-values (y-axis) versus

191    the expected $-\log_{10}$ $P$-values (x-axis). If a sharp diagonal line is observed then the null

192    hypothesis is respected and no significant associations are reported. However, an upper

193    deviated tail from the diagonal line would likely indicate true associations. Upward

194    inflation close to the line's origin indicates suspicious false positives while downward

195    deflated tail suggests false negatives.

196    We empirically evaluated the adjustment of seven models to our data and in **Figure 3**, we

197    plot their Q-Q behavior for significantly associated traits. Despite yielding statistically

198     significant associations, represented as bigger dots, both single locus models GLM and

199     SUPER exhibited prominent inflation beyond the expected null line. This deviation starting

200     close to the origin indicates false positive predictions due to confounding effects

201     (population stratification or genotype relatedness). MLM and CMLM multi-locus models

202     showed matching *P*-value distributions, therefore their Q-Q plots were overlaid. Except for

203     harvest date, where the null hypothesis cannot be rejected with neither inflated nor deflated

204     *P*-values, MLM and CMLM unveiled downshifted line tails when assessed with the rest of

205     traits. Such a result may indicate that these tests were able to reduce false positive

206     associations, but likely yielded false negative ones. Another complex model (MLMM) was

207     found to follow the null hypothesis with both harvest date and flavonoids; nonetheless a

208     slightly downward tail was discerned for fruit weight and sorbitol content. Although being

209     the best-fitting model yielding marker-trait associations with harvest date and flavonoids,

210     FarmCPU did not show the same statistical power with other traits. Finally, the observed *P*-

211     values produced by Blink (green color) were lying on the diagonal line with clear deviated

212     tails toward the y-axis for all six aforementioned traits. All in all, Blink seems to be the best

213     calibrated model, appropriately controlling false positive and false negative effects. For

214     these reasons, we consider Blink as the most suitable model, best adjusted with all

215     phenotypic data and from here on the GWAS results are based on it.

216     **Marker-trait associations and identification of candidate genes**

217     GWAS analysis was conducted on phenotypic traits with moderate to high heritability ($H^2$

218     > 0.5). Consequently, contents of glucose, fructose, sucrose and total sugars were discarded

219     from the subsequent analysis. To sum it up, among the remaining 12 traits, only six were

220     found to be potentially influenced by polymorphic markers. Sixteen marker-trait

221     association peaks were scattered throughout all chromosomes except chr 7 (**Table 2**). In the

222     following sections we will discuss the results for each of these traits, namely harvest date,

223     fruit weight, flesh firmness, and contents of flavonoids, anthocyanins, and sorbitol. For ease

224     of interpretation, in the following paragraphs we summarize the lead SNPs and their

225     corresponding LD blocks. The annotation of 250 kbp regions centering the peak SNPs

226     resulted in a list of candidate causal genes provided in **Table S2**.

227       **Harvest date (HvD)**

228   The GWAS analysis resulted in five SNPs meeting the Bonferroni-adjusted threshold

229 (**Figure 4**). Two SNPs were located on chr 4 and tagged as (SNC_034012.1_10916234,

230 G/T) and (SNC_034012.1_14096987, A/C). Their allelic effect is summarized in **Figure**

231 **S5**, where it can be seen that the first one correlates with delayed harvest and the second

232 one with early one. Another associated marker was located on chr 5

233 (SNC_034013.1_13023165, T/A). Although covering the highest portion of %PVE, no

234 significant allelic effect was observed (**Table 2**). This lead SNP was mapped within the

235 first exon of *Prupe.5G138500*, a gene encoding a germin-like protein. One more significant

236 site was identified on chr 6 and labeled as (SNC_034014.1_7012470, A/T). Allelic effect

237 on phenotypic variation highlighted that both heterozygous and homozygous genotypes

238 carrying the alternate allele (T) were lately harvested with respectively 6 and 13-days of

239 delay (**Figure 4.C).** Similarly**,** the intergenic SNP located on chr 8

240 (SNC_034016.1_18841611, A/G), showed approximately 20-days delay in harvest date

241 with heterozygous accessions (**Figure S5**).

242 LD block analysis revealed various candidate genes, including cell wall modification

243 (*Prupe.8G197700*: galacturonosyltransferase and *Prupe.8G199700*: cell division control

244 protein), cytochrome P450 enzymes (*Prupe.8G196800*, *Prupe.8G196900*, *Prupe.8G197100*

245 and *Prupe.8G197300*), UV-photoreceptor (*Prupe.4G185200*) and ethylene-responsive

246 transcription factor (*Prupe.8G198700*).

247       **Fruit weight (FW)**

248 Significant marker-trait associations were detected on three chromosomes: chr 3

249 (SNC_034011.1_26371177, T/A), chr 6 (SNC_034014.1_1805059, A/G) and chr 8

250 (SNC_034016.1_16407694, A/C). The explained variance oscillated between 17 and 22%,

251 with SNC_034014.1_1805059 tagged as the lead intergenic marker (**Table 2**). The allelic

252 effect of this lead marker (A/G) was found to be unfavorable, with the allele G associated

253 with weight loss (~22 grams) in homozygous accessions (**Figure 5.C**). A similar negative

254 effect was observed with the SNP on chr 3 (T/A), with a significant reduction in fruit

255 weight of 53g. Only marker mapped on chr 8 (A/C) was found to have a positive effect in

256 heterozygous (**Figure S6**). Based on the LD block results, the lead SNP fell within the

257   fourth block, a small interval (84 bp) overlapping no genes (**Figure 5.B**). Nonetheless, the

258   associated SNPs did overlap protein-coding genes. Among them, genes encoding β-

259   galactosidase (*Prupe.3G298200*), α-galactosyltransferase (*Prupe.3G298800*), thymidylate

260   kinase (*Prupe.3G301400*) and transcription factors (GTE8: *Prupe.3G301300* and trihelix

261   GT-4: *Prupe.3G300500*) (**Table S2**).

262       **Flesh Firmness (FF)**

263   A single intergenic marker (SNC_034014.1_7012470; A/T) detected on chr 6 was

264   statistically linked to flesh firmness and explained 33.9% of the total phenotypic variance

265   (**Table 2**). This polymorphism showed a significant increase in the fruit firmness in both

266   heterozygous and alternate homozygous genotypes which underlined the favorable effect of

267   the alternative allele (T) on fruit firmness (**Figure 6.C**). It's noteworthy to mention that this

268   is the only marker simultaneously associated with two different traits (HvD and FF).

269   Moreover, peach accessions carrying the aforementioned allele (either homozygous or

270   heterozygous), were denoted late-harvested and firm peach accessions. Such a result may

271   justify the high correlation existing between both traits (**Figure 1.B**).

272   By examining 250 kbp upstream and downstream the lead marker, it was found to reside in

273   block 3, which makes it a relevant region to seek for candidate firmness-related genes. On

274   the basis of their functional annotation, six genes were selected as potential candidates,

275   including *Prupe.6G100500* encoding an E3 ubiquitin-protein ligase, *Prupe.6G101100*

276   corresponding to vegetative cell wall protein, *Prupe.6G101600* annotated as aquaporin

277   PIP2 and *Prupe.6G102300* encoding homeobox-leucine zipper transcription factor (**Table**

278   **S2).**

279       **Flavonoids (Flvs)**

280   The Manhattan plot displayed two peaks statistically associated with flavonoids content

281   (**Figure S7.A**). The first peak was identified within the intergenic region of chr 2 and

282   named as (SNC_034010.1_643430, T/C). The alternative allele (C) was marked as

283   favorable for heterozygous (TC) and homozygous alternate (CC) genotypes since they

284   showed approximately two-fold increase in the flavonoids content (**Figure S7.C**). The

285   second associated SNP (SNC_034014.1_3066620; G/T) was located on chr 6 and

286   physically mapped on the first exon of *Prupe.6G041500*; a candidate gene encoding a non-

287  specific lipid-transfer protein-like (**Table S2**). The average flavonoids content in alternative

288  homozygous peach accessions (TT) was significantly enhanced compared to the reference

289  homozygous individuals (GG) (**Figure S8**). Thus, the T allele can be considered as a

290  favorable one. Based on LD block results, we annotated a total of 14 genes (**Table S2**).

291  According to their biological function and tissue-specific expression, we narrowed the list

292  to a few promising ones, including two genes encoding transcription factors

293  (*Prupe.2G009100*, bHLH and *Prupe.6G041400*, bZIP).

294      **Anthocyanins (ACNs)**

295  Regarding the anthocyanins content, we detected a single peak signal on chr 5 exceeding

296  the threshold line (**Figure S9.A**). This locus tagged as (SNC_034013.1_12838635; G/T)

297  falls within exon 2 of *Prupe.5G134900*, encoding a B3 domain-containing transcription

298  factor. Thus, *Prupe.5G134900* was considered as a prime candidate gene. The identified

299  marker explained a large portion of the variation (53%), and was found to exert an

300  unfavorable effect on anthocyanins content (**Figure S9.C**). Indeed, pairwise comparisons of

301  SNP allelic effect showed a significantly lower anthocyanins content in the homozygous

302  alternate individuals (TT) compared to the reference homozygous (GG). Screening for

303  genes residing within LD block resulted in three further candidate genes involved in

304  different biological functions (*Prupe.5G134200*, *Prupe.5G134800* and *Prupe.5G135200*)

305  (**Table S2**).

306      **Sorbitol (SRB)**

307  Four significant association signals dispersed on different chromosomes were predicted to

308  affect the sorbitol content (**Table 2** and **Figure S10**). On chr 1, an intergenic SNP

309  (SNC_034009.1_2706825; T/C) explained the lowest proportion of phenotypic variation.

310  The SNP on chr 2 (SNC_034010.1_3682553; G/C), in the third intron of a gene encoding a

311  flowering time control protein (*Prupe.2G0303400*), explained 12% of the PVE. Similarly,

312  (SNC_034014.1_28343678; G/A) was located on chr 6 and mapped on the intronic region

313  of *Prupe.6G320000*, a gene encoding a serine/arginine rich factor. Both *Prupe.2G0303400*

314  and *Prupe.6G320000* are suggested as plausible sorbitol-related genes. The lead SNP

315  explaining the highest PVE (14%) was identified in an intergenic region of chr 8

316  (SNC_034016.1_18841643; G/A).

317    With the exception of (SNC_034014.1_28343678) the remaining loci were observed to

318    have desirable effect on sorbitol content (**Figure S11**). We identified 26 genes distributed

319    in 250 kbp on either side of each associated SNP. Among them, some were discovered to

320    be over-expressed in the fruit (Log$_2$FC > □3□), including genes encoding heavy metal-

321    associated    isoprenylated    proteins    (*Prupe.2G033600*,    *Prupe.2G033700*    and

322    *Prupe.6G321400*), pectinesterases (*Prupe.6G318500*), exonucleases (*Prupe.6G316100*),

323    dormancy-associated proteins (*Prupe.6G319600*), cell cycle checkpoint control proteins

324    (*Prupe.6G321300*) and the E3 ubiquitin-protein ligase RNF4 (*Prupe.8G199600*). A cluster

325    of four cytochrome P450 encoding genes was also identified. This plethora of genes may

326    shed light on several key processes that are subject to influence the sorbitol biosynthesis.


**Discussion**

328              **Performance of variant callers**

329    SNPs discovery in plant genomes has been a widely used strategy for developing molecular

330    markers useful for MAS, genomic selection, phylogenetic analysis, etc. In order to detect

331    and track these genetic variations, we performed a SNP discovery pipeline on paired-end

332    reads   mapped   to   a   diploid   genome   using   BCFtools,   Freebayes,   and   GATK-

333    HaplotypeCaller. SNP calling is known to be error prone. Spurious variants may have

334    several sources; errors associated with sample processing (library preparation, PCR

335    amplification), sequencing, as well as, computational analysis[19]. To remove likely false

336    positive variants, best practices and carefully chosen cut-offs are needed. In our analysis, a

337    SNP site was kept when passing the following filters: mapping and call quality, read depth,

338    as well as call rate and MAF. Though either calling tool can be adapted, we observed a

339    certain inconsistency in the number of high-quality SNPs revealed by each tool. Notably,

340    GATK-HC exhibited the highest sensitivity in SNPs calling, followed by Freebayes then

341    BCFtools. The outperformance of GATK-HC is actually not surprising as it heavily relies

342    on local *de-novo* assembly of haplotypes in active regions[20]. In other terms and unlike the

343    rest of tools, whenever GATK encounters regions with substantial evidence of variation

344    relative to the reference, it discards the existing mapping information and reassembles the

345    read mappings. Our results are in line with[21] concluding that in *Arabidopsis thaliana,*

346    GATK-HC was found to be more accurate compared to BCFtools. Additionally, GATK-

347   HC had the lowest proportion of false positives compared to both Freebayes and

348   BCFtools[22]. On the other hand, the variation in the number of detected SNPs may be partly

349   due to the underlying algorithms. Indeed, GATK-HC and Freebayes are Bayesian variant

350   detectors while BCFtools mpileup uses Hidden Markov Models. Although having an

351   extensive format requirement (e.g: read group specified in the input header), GATK-HC

352   seems to be more precise dealing with ddRAD-seq mapped reads in peach. Nevertheless, to

353   further increase confidence, in this study we only considered SNPs called by all three

354   approaches.

355   **Statistical model selection**

356   Choosing a statistically reliable model is another fundamental pillar for a successful

357   GWAS. Population structure and genetic relatedness are confounding factors increasing the

358   rate of ambiguous associations and decreasing the statistical power. When ignored, they

359   lead to substantial inflation of *P*-values as highlighted in the GLM model (**Figure 3**). In

360   spite of including PCA components and kinship as covariates, SUPER model had also a

361   large number of false positives. This may be explained by the fact that both GLM and

362   SUPER are single-locus approaches failing to catch true associations when dissecting

363   complex traits. Comparable inflated *P*-values were observed in *Arabidopsis thaliana* when

364   testing flowering time, a polygenic trait, with the naïve model (GLM)[23]. In contrast, two

365   other single-locus models, MLM and its compressed version (CMLM), were observed to

366   adjust for false positives at the cost of failing to find any significant marker. Similar results

367   were observed with MLMM, a multi-locus extension of MLM model (**Figure 3**). Overall,

368   we conclude that MLM-based methods are likely missing potentially important SNPs.

369   The inspection of Q-Q plots declared FarmCPU and Blink as the most sophisticated

370   algorithms yielding significant associations. Whereas FarmCPU returned significant

371   signatures with only two traits (HvD and Flvs), Blink consistently inferred associations

372   with six traits (HvD, FW, FF, Flvs, ACNs and SRB). FarmCPU and Blink have emerged to

373   prevent over-fitting and to control false positives simultaneously[24,25]. FarmCPU employs

374   iteratively the fixed-effect model (FEM) and random effect model (REM) to eliminate

375   confounding factors. FEM contains testing markers, one at a time, and associated markers

376   as covariates to control false positives. To circumvent model over-fitting in FEM, the

13

377   associated markers are estimated in REM and are used to derive the kinship[24]. Additionally,

378   FarmCPU relies on the binning approach, where the whole genome is equally divided into

379   bins and only the most significant marker is selected from each bin[24]. Despite its promise,

380   this model is hampered by two major pitfalls: REM is computationally demanding and the

381   assumption of bins rarely occurs in practice. As a consequence, Blink was designed to

382   optimize the computational burden by substituting the REM with FEM through

383   approximating maximum likelihood using the Bayesian Information Criterion and by

384   increasing the statistical power by replacing the bin approach with the LD method[25].

385   Overall, Blink seems to be the well-suited model for our set of data, balancing false

386   positives and false negatives. This statement is underpinned by the GAPIT team, which

387   already stated that Blink is statistically more powerful than FarmCPU[26].

### Marker-trait association for the target traits

389   Out of 16 studied traits, association mapping using ddRAD-derived-SNPs and Blink,

390   revealed association signals with six traits. Totally, 16 significant loci were inferred and

391   distributed as follows: harvest date (chr 4, 5, 6 and 8), fruit weight (chr 3, 6 and 8), flesh

392   firmness (chr 6), flavonoids (chr 2 and 6), anthocyanins (chr 5), and sorbitol (chr 1, 2, 6,

393   and 8). Promising candidate genes were selected when residing within the LD block

394   containing the significant loci, known to be related to the targeted trait and being over-

395   expressed in fruit tissue. Our results were further discussed in comparison with[16] which

396   studied the same phenotypic data and germplasm material, but genotyped using the 9K SNP

397   array instead.

### Harvest date

399   Peaches and nectarines are generally harvested at physiological maturity, then ripening off

400   the trees. Harvest date and maturity date are frequently used as synonyms and are expressed

401   in Julian days. HvD is defined as the day on which a certain percentage of peaches reach

402   maturity. Maturity date (MD) is defined as the interval of time from the first day of the

403   calendar year till the harvest date[27]. In our study, five association signals for HvD were

404   highlighted. Two were inferred on chromosome 4 and the rest were distributed on

405   chromosomes 5, 6 and 8.

14

406    As established by[28,29], major QTLs controlling maturity date have been reported on linkage

407    groups LG4 and LG6 (**Table S3**). Particularly, a major QTL on LG4 referred to as qMD4.1

408    showed a pleiotropic effect on fruit weight and firmness[28,30]. Interestingly, our marker

409    SNC_034012.1_10916234 mapped at (~10.91 Mbp), was overlapping the (qMD4.1_CA)

410    locus from C×A progeny spanning the interval between 10.87-12.09 Mbp[30]. This same

411    QTL from W×By progeny (qMD4.1_WB) was found 65 kbp from our marker (**Figure**

412    **4.D**). In the same vein, SNC_034012.1_10916234 was delimited by one downstream (HD-

413    EJ-4)[31] and two upstream quantitative loci (qP-MD4)[32] and (qMD4_1)[33] mapped

414    respectively at 0.5, 5.3 and 1245 kbp from the SNP's coordinate (**Figure 4.D**). Likewise,

415    the second marker on chr 4 (SNC_034012.1_14096987) mapped at (~14.09 Mbp) was

416    found within the genomic region of strong confidence QTL (qMD4_2) spanning the

417    interval (11.20 - 14.10 Mbp)[34]. Contrasting with associated SNP from the 9K assay[16], our

418    markers seems to be more confident as they are located within the QTL boundaries which

419    supports their reliability. Altogether, we anticipate that the aforementioned SNPs on chr 4

420    could be integrated as promising markers for HvD breeding goals. As well, we conclude

421    that LG4 seems to be a chromosomal hotspot hosting a cluster of major QTLs associated

422    with the maturity date. QTLs influencing maturity date were also detected on LG4 in

423    peach-related species, for instance; sweet cherry[35]. Therefore, we believe that this trait

424    could be controlled by orthologous loci within *Prunus* species.

425    Marker 'SNC_034013.1_13023165' mapped on chr 5 (~13.02 Mbp) was supported by an

426    adjacent locus (QTLMD5) spanning the region (14.38 - 17.64 Mbp)[36] and other distant

427    signals (qP-MD5 and qMD5)[32,34]. Significant markers from 9K array[16] were found to be

428    physically closer to the QTLs (**Figure 4.D** and **Table S3**). Finally, the significant SNP on

429    chr 6 'SNC_034014.1_7012470' was residing within two QTL intervals[36] QTLMD6.1 and

430    QTLMD6.1, supporting it. Similar findings were observed with 9K-associated markers.

431    Multiple candidate genes potentially influencing the harvest date were shortlisted (**Table**

432    **S2**). Most importantly, an ethylene-responsive transcription factor (*Prupe.8G198700*).

433    Ethylene-responsive elements are relevant in climacteric fruits and have been proposed as

434    candidate genes for fruit maturation date in different *Prunus* species[31,37]. We also identified

435    a cell wall remodeling gene encoding galacturonosyltransferase. This finding is in

436 consonance with[37] defining a galacturonosyltransferase as a candidate gene for late
437 harvested cultivars.

### Fruit weight

439 Fruit weight is a quantitative trait with great importance in peach breeding. Previous studies
440 in peach have divulged that FW is monitored by multiple QTLs distributed across all
441 chromosomes[34,38–40]. Using GWAS, we identified a significant SNP on chr 3 (~26.37 Mbp)
442 located respectively at 4.07 and 7.27 Mbp downstream of two QTLs qFRW.ZC_3 and
443 qFRW.WB (**Figure 5, Table S3**). On chr 6, another significant marker was predicted at
444 (~1.80 Mbp). This marker was delimited in near proximity by two reliable QTLs
445 (qFRW.ZC_6)[40] and (qFW6.1)[34], situated respectively at 387 and 1,358 bp. On chr 8,
446 SNC_034016.1_16407694, was localized at (~6 Mbp) downstream of marker flanking QTL
447 (FW 10-b)[39]. This is in contrast with[16] where no associated loci were reported for this trait
448 (**Table S3**). Such results support the relevance of our findings in dissecting the genetic
449 control of complex fruit traits and shed light on the effectiveness of ddRAD-seq genotyping
450 on inferring *novel* association signatures.
451 Candidate genes prediction revealed two transcription factors, trihelix GT-4
452 (*Prupe.3G300500*) and GTE-8 (*Prupe.3G301300*). Transcriptional regulators are abundant
453 in plant genomes and they are implicated in various biological processes. Interestingly,
454 trihelix genes are known to be photo-responsive proteins[41]. It's well documented that light
455 exposure affects fruit size, shape and quality[42]. Thus, we speculate that trihelix TF may
456 regulate the fruit weight in peaches. Moreover, cell wall enzymes such as β-galactosidase,
457 α-galactosyltransferase may act as key components of cell wall turnover during stone fruit
458 growth[43]. Finally, thymidylate kinase exhibited strong upregulation suggesting a possible
459 role in peach fruit development as validated in rice, barley and maize[44].

### Flesh firmness

461 Firmness is a key textural indicator of peach quality and directly influences their shelf life.
462 In our study, we identified a single firmness related locus SNC_034014.1_7012470 on chr
463 6. In the same LG6, a firmness loss QTL (qP-FL5d6) was described (**Figure 6.D** and **Table
464 S3**). Another stable QTL (qP-FF6.1[m]) was also detected over two years in related species,

16

465    particularly in sweet cherry[35]. Using 9K inferred SNPs and MLM model[16], no significant

466    association signals were found.

467    Four genes were selected as strong candidates encoding: ubiquitin-protein ligase

468    (*Prupe.6G100500*), vegetative cell protein (*Prupe.6G101100*), aquaporin PIP2

469    (*Prupe.6G101600*) homeobox-leucine zipper protein (*Prupe.6G102300*). E3 ligase genes

470    were found to be differentially expressed in either melting flesh or stony hard fruit during

471    the ripening[48]. Aquaporins are transmembrane water transporters and water uptake within

472    fruit is highly related with fruit firmness[45]. Thus, aquaporins could play a key role in

473    maintaining cell turgor in peach. Finally, homeobox-leucine zipper proteins were denoted

474    as potential biomarkers for the ripening process in peach[46].


475            **Flavonoid and anthocyanin contents**

476    Flavonoids are major polyphenol compounds playing a central role in fruit color and flavor.

477    Our analysis yielded two potential association signatures in chr 2 and 6. These results go

478    along with[47] affirming that the majority of lead SNPs linked with many flavonoid

479    metabolites in peach were located on chr 2. Herein, SNC_034010.1_643430 was supported

480    by two QTLs[39] identified in Venus × Bigtop progeny and named as 'FLV 10-a' and 'FLV

481    10-b' (**Figure S7.D** and **Table S3**). It's well documented that flavonoid biosynthesis is a

482    complex pathway, transcriptionally regulated by members of Myb and bHLH families[48].

483    Although no Myb encoding gene was found in our analysis, a highly up-regulated bHLH-

484    TF was inferred and may be considered as a promising candidate gene involved in

485    flavonoid regulation.

486    Anthocyanins constitute an important group of plant pigments belonging to the flavonoid

487    family. Their differential accumulation in peach results in the distinctive fruit and flesh

488    color[48]. Although there is strong evidence that their biosynthesis is mainly regulated by a

489    Myb10 transcription factor on LG3, many anthocyanin-related QTLs were identified on

490    LG4, LG5, and LG6[34,39,40]. Our analysis detected a single lead marker on chr 5 accounting

491    for ~53% of the PVE. Thus, 'SNC_034013.1_12838635' may be a preferential target for an

492    effective marker assisted selection. It was delineated on both downstream and upstream

493    sides by (qANT)[39], (qATCYN.ZC)[40] and (qPSC5)[34]. When genotyped with the 9K array[16],

494    no associated markers were detected on LG5. Remarkably, our polymorphic marker was

495    physically falling in the exonic region of *Prupe.5G134900*, a gene encoding a B3 domain-

496     containing transcription factor. Although the functional relevance of this prime gene

497     requires further validation, we hypothesize that the genetic control of anthocyanins may be

498     driven by B3 DNA-binding protein. Curiously, for both anthocyanins and flavonoids, a B3

499     family transcription was selected as candidate gene (respectively *Prupe.5G134900* and

500     *Prupe.6G041000*). This may be explained by the fact that anthocyanins are a class of water-

501     soluble flavonoids. Thus, it's plausible to hypothesize that genes involved in flavonoids and

502     anthocyanins regulation are in coordination.

503     **Sorbitol**

504     Sugar content is one of the most important quality traits perceived by the consumers. The

505     sweetness intensity depends on the overall sugar amount brought by sucrose, glucose,

506     fructose and sorbitol. These first three sugar types were discarded from our analysis as they

507     didn't meet the heritability cutoff. Regarding the sorbitol, association signatures were found

508     in chr 1 (~27.06 Mbp), chr 2 (~3.68 Mbp), chr 6 (~28.34 Mbp) and chr 8 (~18.84 Mbp).

509     Genetic mapping has been extensively carried out to identify key QTLs responsible for

510     sorbitol biosynthesis. A reliable QTL (qSOR_1) was mapped on the upper region of LG1,

511     nearly 17.5 Mb upstream of our associated marker (**Figure S10.D**). Compared to the 9K

512     association study[18], no significant association signal was detected on LG1 (**Table S3**). On

513     chr 2, we were able to find an adjacent QTL supporting the accuracy of our results[49].

514     Indeed, qSOR_2 was positioned at ~1.2 Mbp from our marker SNC_034010.1_3682553.

515

516     Finally, this work depicts ddRAD-seq genotyping as an efficient approach for SNPs

517     detection and association studies. Akin to the 9K SNP array, ddRAD-seq yielded valuable

518     markers strongly supported by stable QTLs. However, while SNP arrays are engineered to

519     specifically include polymorphic loci from genomic regions of interest and focus on

520     harboring SNPs known to be linked to commercially important traits, ddRAD-seq samples

521     the genome randomly, without prior knowledge of target regions. For this reason, ddRAD-

522     seq might be a better fit for analyses concerned with unexplored biological processes.

523     Concisely, we successfully used ddRAD-seq-derived SNPs to identify genomic regions and

524     genes influencing major fruit-related traits in peaches. The inferred associated SNPs

525     appeared to be reliable as they often explained a fairly high percentage of the total

526 phenotypic variance. The survey of candidate genes for these relevant polymorphic sites

527 rendered plenty of genes implicated in various processes. Genes harboring significant

528 markers may be considered as preferential targets for peach breeding. However, due to the

529 complexity of the examined traits, future functional validation would provide additional

530 hints to support the breeding efforts.


531 **Material and Methods**


532 **Plant material and phenotypic evaluation**

533 A total of 90 peach and nectarine accessions were used for double digest restriction-site

534 associated sequencing (ddRAD-seq) and subsequent GWAS analysis. The germplasm panel

535 comprises 73 landraces and 17 modern breeding lines originating from Spain, United

536 States, France, Italy, New Zealand, and South Africa. All genotypes were grown under

537 Mediterranean soil conditions at the Experimental Station of Aula Dei (CSIC) located at

538 Zaragoza, Spain (41.7245 °N, 0.8118 °W) and analyzed during three fruiting seasons

539 (2008-2010). Information about plant accessions is summarized in **Table S1**.


540 The phenotypic data previously reported by[16] were re-analyzed in the present study.

541 Briefly, 16 traits were evaluated by randomly harvesting 20 fruits from each cultivar at the

542 commercial maturity during three years. Traits were split into two categories. Agronomic

543 features included harvest date (HvD; Julian days), fruit weight (FW; grams), flesh firmness

544 (FF; Newton), soluble solids content (SSC; °Brix), titratable acidity (TA; grams malic

545 acid/100 g flesh weight) and ripening index (RI; SSC/TA). Besides, biochemical variables

546 comprised vitamin C (Vit C; mg of ascorbic acid/100 g flesh weight), total phenolics (Phen;

547 mg of gallic acid equivalents/100 g flesh weight), contents of flavonoid (Flv; catechin

548 equivalents/100 g flesh weight) and anthocyanin (ACNs; cyanidin-3-glucoside/kg flesh

549 weight), sucrose (Suc; g/kg flesh weight), glucose (Glu; g/kg flesh weight), fructose (Fruc;

550 g/kg flesh weight), sorbitol (SRB; g/kg flesh weight), and total sugars (TS; g/kg flesh

551 weight) and relative antioxidant capacity (RAC; µg TE/g flesh weight).


552 Variance components and broad sense heritability ($H^2$) were estimated using the variability

553 R package v0.1.0. Only traits with $H^2 > 0.5$ were considered for association analysis.


19

554     Distribution of averaged phenotypic data was checked in R using Shapiro-Wilk test. Non-

555     normal distributions were transformed using bestNormalize package (v1.8.3)[50].

**DNA extraction and enzyme evaluation**

557     Genomic DNA was extracted from leaves using the DNeasy Plant Mini Kit (Qiagen,

558     Dusseldorf, Germany) following the manufacturer's recommendations. DNA concentration

559     and quality were checked using PicoGreen®dye and measured in a fluorospectrometer.

560     Whole-genome genotyping was carried out using ddRAD-seq approach by combining low

561     and high frequency cutter to digest DNA; respectively *Pst1* and *Mbol* as described in peach

562     by[51]. This enzyme pair yielded the highest number of loci with a size range between 300

563     and 400 bp and prevented repetitive region sampling. Selected loci are those having the

564     sticky ends of both enzymes[52].

**ddRAD libraries preparation and sequencing**

566     DNA libraries were constructed at the Genomic Unit at IABiMo INTA-CONICET

567     (Argentina) following[51,52] recommendations. Shortly, digested DNA with *Pst1*/*Mbol* pair

568     were gel excised, eluted then ligated to barcoded adapters specific to each sample. Ligated

569     fragments from 24 samples were subsequently pooled together and were PCR amplified

570     with indexed primers to tag each pool. Finally, paired-end reads (250 bp) were generated on

571     an Illumina NovaSeq 6000 instrument at CIMMYT, Mexico. The raw sequencing data was

572     deposited in the European Nucleotide Archive (ENA) under the BioProject PRJEB62784.

**Data processing and alignment**

574     Raw reads were de-multiplexed and trimmed using the process-radtag module from

575     STACKS suite (v2.59)[53]. After quality assessment, paired-end reads were mapped to

576     *Prunus persica* reference genome v2 (GCF_000346465.2, retrieved from NCBI RefSeq[54]

577     using BWA-mem (v0.7.17)[55]. Redundant reads known as PCR duplicates were expurgated

578     from downstream analysis as described by[22]. The resultant de-duplicated files were sorted

579     and indexed using SAMtools[56] to be ready for variant calling.

**Variant discovery pipeline**

581     Variant calling was conducted in a single-sample mode testing the performance of three

582     variant callers: BCFtools (v1.7)[56], Freebayes (v1.0.0)[57] and GATK-HaplotypeCaller

20

583    (v4.2.3.0)[20]. Raw SNPs underwent standard quality filtering based on mapping quality (MQ

584    > 40), variant quality (QUAL < 30) and depth of reads (DP ≥ 5) to remove artifactual calls.

585    Consequently, clean SNPs from each calling method were merged by position and by

586    reference/alternative alleles into multi-sample VCF files. SNPs resulting from the

587    intersection of multi-samples VCFs were considered as highly accurate calls and were

588    inspected to remove multi-allelic variants and those assigned to scaffolds. Then, they were

589    filtered by call rate > 80% and residual missing genotypes were imputed with beagle's

590    default settings (v4.1)60. The imputation accuracy was evaluated in Tassel (v5.0)[58] by

591    masking 1% of the genotype and calculating the error rate. SNPs with (MAF > 0.05) were

592    selected as a final call set to determine the population structure and marker-trait

593    associations. SNP identifiers were created by concatenating their assigned chromosome and

594    their base pair position (eg: SNC_034014.1_7012470).

595    **Linkage disequilibrium and population structure**

596    Intra-chromosomal LD was calculated using Plink (v1.9)[59], as a measure of *Pearson*

597    correlation coefficient ($r^2$) between marker-pairs. For each chromosome, LD distribution

598    was plotted against its physical distance (Mbp). The LD decay curve was estimated as the

599    average of $r^2$ variation across 100 kbp bins and was fitted in R program. LD decay was

600    defined by setting $r^2$=0.25 as a threshold. LD decay extent was defined as the physical

601    genomic distance at which the $r^2$ decreased to half of its maximum value. Polymorphic sites

602    showing strong LD were pruned in Plink by delimiting a window of 10 SNPs, removing

603    one of the SNPs pair with $r^2$ > 0.25 and then shifting the window 5 SNPs forward

604    repeatedly. Genetic distance and kinship matrix between pairs of genotypes were computed

605    using the centered identity-by-state method implemented in Tassel.

606    LD-pruned SNPs were selected to infer the population stratification of the GWAS panel

607    using two complementary approaches. First, the Bayesian clustering algorithm

608    implemented in fastSTRUCTURE (v1.04)[60] was tested on predefined K subgroups ranging

609    from 1 to 10. The optimal K value was estimated based on the lowest cross validation error.

610    Then, principal component analysis was computed with SmartPCA (v1.1.0) R-package[61].

**Genome wide association study**

For association mapping, seven statistical models, ranging from single to multi-locus, were simultaneously tested in GAPIT (v3.1.0)[62] Single locus models include general linear model (GLM), mixed linear model (MLM), compressed MLM (CMLM), and settlement of MLMs under progressively exclusive relationship (SUPER). Multi-locus algorithms comprise multiple loci mixed linear model (MLMM), fixed and random model circulating probability unification (FarmCPU), and Bayesian-information and linkage-disequilibrium iteratively nested keyway (Blink). Except for GLM, where no genotype relatedness was involved, population structure and kinship were both fitted as covariates in all models. Indeed, the first four PCA components and kinship were introduced respectively as fixed and random effects to reduce the false positives. The statistical model best fitting the data was chosen based on the quantile-quantile plot and the number of significant markers. Significantly associated markers were shortlisted based on the Bonferroni correction ($-\log_{10}(0.05)/13045 = 5.42$) and Manhattan plots were generated accordingly using CMplot package (v4.2.0)[63]. Statistically significant SNPs explaining at least 10% of the phenotypic variance (%PVE) were considered as most promising predictions and used for LD block analysis. Moreover, markers accounting for the largest proportion of phenotypic variance are hereinafter referred as 'lead SNPs'.

**Annotation of SNP effects and identification of favorable alleles**

First, genomic coordinates of SNPs were used to query Ensembl Plants REST services in order to obtain annotations of their effect on neighbor genomic features. In particular we used the Ensembl Variant Effect Predictor (VEP)[64] and a modification of recipe R8[65].

Then, allelic effect of significant SNP loci on trait variation was estimated through pairwise comparisons between the phenotypic values of the different genotypes: homozygous reference (0/0), heterozygous (0/1) and homozygous alternative (1/1). An allele is defined as favorable when a significant increase of the phenotypic value was observed between the homozygous reference and the remaining genotypes. Pairwise comparisons were run using the Games-Howell test and *P*-values were corrected for multiple testing using the FDR method. Results were visualized using ggstatsplot R-package (v0.9.1)[66].

**LD-block analysis and identification of candidate genes**

Significant SNPs were examined to identify candidate genes. Initially, it was considered whether polymorphisms would be localized in genic regions. Thereby, SNPs were mapped based on their physical position to *Prunus persica* genome (GCF_000346465.2). SNP-anchored genes were called 'prime candidates'. Strong candidate genes were shortlisted when meeting three criteria: falling within the LD-block region harboring the significant SNPs, being functionally related to the trait of interest and being differentially expressed in fruit tissue. Expression information was retrieved from a recent study by[67] which defined modules of co-expressed genes across different peach tissues and under various experimental conditions. Differentially expressed genes were those outlined in fruit experiments, particularly under cold storage and chilling injury.

LD-blocks were identified within 250 kbp windows upstream and downstream the lead sites. Block boundaries were delimited using a solid spine partitioning approach from LDBlockShow tool[68]. A block is defined as a group of SNPs that are in strong LD (D' ≥ 0.7) with the first and last marker in the same block. A D' value of 0 denotes complete linkage equilibrium, which implies frequent pairwise recombination between markers. Conversely, a D' value of 1 indicates a complete linkage disequilibrium. Note that D' and $r^2$ are common measures of non-random association between two or more loci; while D' refers to the co-inheritance of two alleles, $r^2$ considers the allele frequency to distinguish between common and rare. Identified LD blocks were therefore scanned for candidate genes via NCBI genome data viewer[69].

**QTLs review for fruit quality traits in peach**

To benchmark the accuracy of our results, an exhaustive bibliographic review of previously reported QTLs mapped in the same linkage group as the associated markers was done. In a practical term, if an associated SNP is located nearby or within a QTL interval, then it's considered as highly accurate and likely segregate with the observed trait. In case that the QTL boundaries are not defined as physical intervals (in bp), we used the nearest or the co-localizing markers as reference. Herein, we refer to the nearest marker as the closest one with a maximum of 5 cM from the QTL hotspot while the co-localizing marker is the one mapped within the QTL boundaries. Moreover, we calculated the physical distance

670 separating our predicted associated markers from the QTLs. Finally, we compared these
671 distances with a previous study using the same phenotypic data and peach material,
672 although genotyped using the 9K SNP array[16].

**Data availability**

674 Raw sequence data and final variant call file (vcf) have been deposited in the European
675 Nucleotide Archive (ENA) under the BioProject PRJEB62784 (data will be released at the
676 publication date). Source code and documentation can be accessed at
677 https://github.com/najlaksouri/GWAS-Workflow.

**Conflict of interests**

690 The authors declare no conflict of interest.

**Contributions**

692 Y.G., and B.C-M. conceived the project and its components. C.F-i-F. collected the samples,
693 extracted DNA and performed the phenotyping. G.S helped to process genotyping. N.K
694 performed the bioinformatic analysis. Y.G., and B.C-M. assisted the analysis and discussed

695    the results. N.K wrote the manuscript. Y.G., and B.C-M. reviewed and edited the text. All

696    authors read and approved the article.

**References**

698    1.    Micheletti, D. *et al.* Whole-genome analysis of diversity and SNP-major gene
699          association in peach germplasm. *PLoS One* **10**, 1–19 (2015).

700    2.    Gogorcena, Y., Sánchez, G., Moreno-Vázquez, S., Pérez, S. & Ksouri, N. Genomic-
701          based breeding for climate-smart peach varieties. in Genome designing of climate-
702          smart fruit crops (ed. Kole, C.) 291–351 (Springer-Nature, 2020).
703          doi:https://doi.org/10.1007/978-3-319-97946-5_9.

704    3.    Cirilli, M. *et al.* Genetic and phenotypic analyses reveal major quantitative loci
705          associated to fruit size and shape traits in a non-flat peach collection (*P. persica* L.
706          Batsch). *Hortic. Res.* **8**, 1–17 (2021).

707    4.    Da Silva Linge, C. *et al.* Multi-locus genome-wide association studies reveal fruit
708          quality hotspots in peach genome. *Front. Plant Sci.* **12**, 1–18 (2021).

709    5.    Aranzana, M. J. *et al. Prunus* genetics and applications after *de novo* genome
710          sequencing: achievements and prospects. *Hortic. Res.* **6**, 1–25 at
711          https://doi.org/10.1038/s41438-019-0140-8 (2019).

712    6.    De Mori, G. & Cipriani, G. Marker-assisted selection in breeding for fruit trait
713          improvement: a review. *Int. J. Mol. Sci.* **24**, 1–39 (2023).

714    7.    Kumar, S. *et al.* GWAS provides new insights into the genetic mechanisms of
715          phytochemicals production and red skin colour in apple. *Hortic. Res.* **9**, 1–10 (2022).

716    8.    Holušová, K. *et al.* High-resolution genome-wide association study of a large Czech
717          collection of sweet cherry (*Prunus avium* L.) on fruit maturity and quality traits.
718          *Hortic. Res.* **10**, 1–15 (2023).

719    9.    Pavan, S. *et al.* Recommendations for choosing the genotyping method and best
720          practices for quality control in crop genome-wide association studies. *Front. Genet.*
721          **11**, (2020).

722    10.   Verde, I. *et al.* Development and evaluation of a 9K SNP array for peach by
723          internationally coordinated SNP detection and validation in breeding germplasm.
724          *PLoS One* **7**, e35668 (2012).

725    11.   Thurow, L. B., Gasic, K., Bassols Raseira, M. do C., Bonow, S. & Marques Castro,
726          C. Genome-wide SNP discovery through genotyping by sequencing, population
727          structure, and linkage disequilibrium in Brazilian peach breeding germplasm. *Tree*
728          *Genet. Genomes* **16**, 1–14 (2020).

729    12.   Mas-Gómez, J., Cantín, C. M., Moreno, M. Á. & Martínez-García, P. J. Genetic

diversity and genome-wide association study of morphological and quality traits in peach using two Spanish peach germplasm collections. *Front. Plant Sci.* **13**, 1–14 (2022).

13. Geibel, J. *et al.* How array design creates SNP ascertainment bias. *PLoS One* **16**, 1–23 (2021).

14. Peterson, B. K., Weber, J. N., Kay, E. H., Fisher, H. S. & Hoekstra, H. E. Double digest RADseq: An inexpensive method for *de novo* SNP discovery and genotyping in model and non-model species. *PLoS One* **7**, 1–11 (2012).

15. Font i Forcada, C., Oraguzie, N., Igartua, E., Moreno, M. Á. & Gogorcena, Y. Population structure and marker-trait associations for pomological traits in peach and nectarine cultivars. *Tree Genet. Genomes* **9**, 331–349 (2013).

16. Font i. Forcada, C., Guajardo, V., Chin-Wo, S. R. & Moreno, M. Á. Association mapping analysis for fruit quality traits in *Prunus persica* using SNP markers. *Front. Plant Sci.* **9**, 1–12 (2019).

17. Cao, K. *et al.* Genome-wide association study of 12 agronomic traits in peach. *Nat. Commun.* **7**, 1–10 (2016).

18. Verde, I. *et al.* The high-quality draft genome of peach (*Prunus persica*) identifies unique patterns of genetic diversity, domestication and genome evolution. *Nat. Genet.* **45**, 487–494 (2013).

19. Olson, N. D. *et al.* Best practices for evaluating single nucleotide variant calling methods for microbial genomics. *Front. Genet.* **6**, 1–15 (2015).

20. Van der Auwera, G. A. et al. From fastq data to high-confidence variant calls: The genome analysis toolkit best practices pipeline. Current Protocols in Bioinformatics (2013).

21. Schilbert, H. M., Rempel, A. & Pucker, B. Comparison of read mapping and variant calling tools for the analysis of plant NGS data. *Plants* **9**, 1–14 (2020).

22. Ksouri, N. *et al.* ddRAD-seq variant calling in peach and the effect of removing PCR duplicates. *Acta Hortic.* **1352**, 405–412 (2022).

23. Atwell, S. *et al.* Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred lines. *Nature* **465**, 627–631 (2010).

24. Liu, X., Huang, M., Fan, B., Buckler, E. S. & Zhang, Z. Iterative usage of fixed and random effect models for powerful and efficient genome-wide association studies. *PLoS Genet.* **12**, 1–24 (2016).

25. Huang, M., Liu, X., Zhou, Y., Summers, R. M. & Zhang, Z. BLINK: A package for the next level of genome-wide association studies with both individuals and markers in the millions. *Gigascience* **8**, 1–12 (2019).

766   26.   Wang, J., Tang, Y. & Zhang, Z. Performing genome-wide association studies with multiple models using GAPIT. in *Genome-Wide Association Studies* (eds. Torkamaneh, D. & Belzile, F.) 199–217 (Springer US, 2022). doi:10.1007/978-1-0716-2237-7_13.

770   27.   Veerappan, K., Natarajan, S., Chung, H. & Park, J. Molecular insights of fruit quality traits in peaches, *Prunus persica*. *Plants* **10**, 1–14 (2021).

772   28.   Eduardo, I. et al. QTL analysis of fruit quality traits in two peach intraspecific populations and importance of maturity date pleiotropic effect. *Tree Genet. Genomes* **7**, 323–335 (2011).

775   29.   Dirlewanger, E. *et al.* Comparison of the genetic determinism of two key phenological traits, flowering and maturity dates, in three *Prunus* species: Peach, apricot and sweet cherry. *Heredity* **109**, 280–292 (2012).

778   30.   Pirona, R. *et al.* Fine mapping and identification of a candidate gene for a major locus controlling maturity date in peach. *BMC Plant Biol.* **13**, 1–13 (2013).

780   31.   Romeu, J. F. *et al.* Quantitative trait loci affecting reproductive phenology in peach. *BMC Plant Biol.* **14**, 1–16 (2014).

782   32.   Serra, O. *et al.* Genetic analysis of the slow-melting flesh character in peach. *Tree Genet. Genomes* **13**, 1–13 (2017).

784   33.   Kalluri, N., Eduardo, I. & Arús, P. Comparative QTL analysis in peach 'Earlygold' F2 and backcross progenies. *Sci. Hortic.* **293**, 1–8 (2022).

786   34.   Hernández Mora, J. R. *et al.* Integrated QTL detection for key breeding traits in multiple peach progenies. *BMC Genomics* **18**, 1–15 (2017).

788   35.   Calle, A. & Wünsch, A. Multiple-population QTL mapping of maturity and fruit-quality traits reveals LG4 region as a breeding target in sweet cherry (*Prunus avium* L.). *Hortic. Res.* **7**, 1–13 (2020).

791   36.   Nuñez-Lillo, G. *et al.* High-density genetic map and QTL analysis of soluble solid content, maturity date, and mealiness in peach using genotyping by sequencing. *Sci. Hortic.* **257**, 1–11 (2019).

794   37.   Núñez-Lillo, G. *et al.* Transcriptome and gene regulatory network analyses reveal new transcription factors in mature fruit associated with harvest date in *Prunus persica*. *Plants* **11**, 1–17 (2022).

797   38.   Da Silva Linge, C. *et al.* Genetic dissection of fruit weight and size in an F2 peach (*Prunus persica* (L.) Batsch) progeny. *Mol. Breed.* **35**, 1–19 (2015).

799   39.   Zeballos, J. L. *et al.* Mapping QTLs associated with fruit quality traits in peach [*Prunus persica* (L.) Batsch] using SNP maps. *Tree Genet. Genomes* **12**, 1–7 (2016).

801   40.   Abdelghafar, A., Da Silva Linge, C., Okie, W. R. & Gasic, K. Mapping QTLs for

802      phytochemical compounds and fruit quality in peach. *Mol. Breed.* **40**, 1–18 (2020).

803  41.  Kaplan-Levy, R. N., Brewer, P. B., Quon, T. & Smyth, D. R. The trihelix family of
804      transcription factors - light, stress and development. *Trends Plant Sci.* **17**, 163–171
805      (2012).

806  42.  Reale, L. *et al.* The influence of light on olive (*Olea europaea* L.) fruit development
807      is cultivar dependent. *Front. Plant Sci.* **10**, 1–10 (2019).

808  43.  Canton, M. *et al.* Metabolism of stone fruits: reciprocal contribution between
809      primary metabolism and cell wall. *Front. Plant Sci.* **11**, 1–10 (2020).

810  44.  Zhu, L. *et al.* A shortest-path-based method for the analysis and prediction of fruit-
811      related genes in Arabidopsis thaliana. *PLoS One* **11**, 1–14 (2016).

812  45.  Quirante-Moya, F., Martinez-Alonso, A., Lopez-Zaplana, A., Bárzana, G. &
813      Carvajal, M. Water relations after Ca, B and Si application determine fruit physical
814      quality in relation to aquaporins in *Prunus*. *Sci. Hortic.* **293**, 1–14 (2022).

815  46.  Nilo-Poyanco, R., Moraga, C., Benedetto, G., Orellana, A. & Almeida, A. M.
816      Shotgun proteomics of peach fruit reveals major metabolic pathways associated to
817      ripening. *BMC Genomics* **22**, 1–29 (2021).

818  47.  Cao, K. *et al.* Combined nature and human selections reshaped peach fruit
819      metabolome. *Genome Biol.* **23**, 1–25 (2022).

820  48.  Wang, J. *et al.* Two MYB and three bHLH family genes participate in anthocyanin
821      accumulation in the flesh of peach fruit treated with glucose, sucrose, sorbitol, and
822      fructose *in vitro*. *Plants* **11**, 1–14 (2022).

823  49.  Desnoues, E. *et al.* Dynamic QTLs for sugars and enzyme activities provide an
824      overview of genetic control of sugar metabolism during peach fruit development. *J.
825      Exp. Bot.* **67**, 3419–3431 (2016).

826  50.  Peterson, R. A. Finding optimal normalizing transformations via bestNormalize.
827      *Contrib. Res. J.* **13**, 310–329 (2021).

828  51.  Aballay, M. M., Aguirre, N. C., Filippi, C. V., Valentini, G. H. & Sánchez, G. Fine-
829      tuning the performance of ddRAD-seq in the peach genome. *Sci. Rep.* **11**, 1–13
830      (2021).

831  52.  Aguirre, N. C. *et al.* Optimizing ddRADseq in non-model species: A case study in
832      Eucalyptus dunnii Maiden. *Agronomy* **9**, 1–21 (2019).

833  53.  Rochette, N. C., Rivera-Colón, A. G. & Catchen, J. M. Stacks 2: Analytical methods
834      for paired-end sequencing improve RADseq-based population genomics. *Mol. Ecol.*
835      **28**, 4737–4754 (2019).

836  54.  O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: Current status,
837      taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–D745

838       (2016).

839    55.  Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
840       transform. *Bioinformatics* **25**, 1754–1760 (2009).

841    56.  Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, 1–4
842       (2021).

843    57.  Garrison, E. & Marth, G. Haplotype-based variant detection from short-read
844       sequencing. *ArXiv* **1207**, 1–5 (2012).

845    58.  Bradbury, P. J. *et al.* TASSEL: Software for association mapping of complex traits
846       in diverse samples. *Bioinformatics* **23**, 2633–2635 (2007).

847    59.  Purcell, S. *et al.* PLINK: A tool set for whole-genome association and population-
848       based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).

849    60.  Raj, A., Stephens, M. & Pritchard, J. K. FastSTRUCTURE: Variational inference of
850       population structure in large SNP data sets. *Genetics* **197**, 573–589 (2014).

851    61.  Herrando-Pérez, S., Tobler, R. & Huber, C. D. Smartsnp, an R package for fast
852       multivariate analyses of big genomic data. *Methods Ecol. Evol.* **12**, 2084–2093
853       (2021).

854    62.  Wang, J. & Zhang, Z. GAPIT Version 3: Boosting power and accuracy for genomic
855       association and prediction. *Genomics, Proteomics Bioinformatics.* **19**, 629–640
856       (2021).

857    63.  Yin, L. *et al.* rMVP: A memory-efficient, visualization-enhanced, and parallel-
858       accelerated tool for genome-wide association study. *Genomics, Proteomics*
859       *Bioinformatics.* **19**, 619–628 (2021).

860    64.  McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 1–14
861       (2016).

862    65.  Contreras-Moreira, B. *et al.* Scripting analyses of genomes in Ensembl Plants. In*:*
863       Edwards, D. (eds) *Plant Bioinformatics. Methods in Molecular Biology* (ed.
864       Edwards, D.), **2443**. 27-55 (New York, NY, 2022). https://doi.org/10.1007/978-1-
865       0716-2067-0_2 65.

866    66.  Patil, I. Visualizations with statistical details: The 'ggstatsplot' approach. *J. Open*
867       *Source Softw.* **6**, 1–5 (2021).

868    67.  Ksouri, N. *et al.* Tuning promoter boundaries improves regulatory motif discovery in
869       nonmodel plants□: the peach example. *Plant Physiol.* **185**, 1242–1258 (2021).

870    68.  Dong, S. S. *et al.* LDBlockShow: A fast and convenient tool for visualizing linkage
871       disequilibrium and haplotype blocks based on variant call format files. *Briefings in*
872       *Bioinformatics.* 1–6 (2020) doi:10.1093/bib/bbaa227.

873    69.  Rangwala, S. H. *et al.* Accessing NCBI data using the NCBI sequence viewer and

874     genome data viewer (GDV). *Genome Res.* **31**, 159–169 (2021).

875

876     **Figure legends**

877     **Figure 1**. (**A**): Broad sense heritability estimates and (**B**): phenotypic correlation among 16

878     peach agronomical and fruit quality traits. Dashed horizontal line corresponds to heritability

879     threshold ($H^2 = 0.5$). Correlation between traits was estimated using *Pearson* correlation.

880     Significant positive and negative correlations are displayed in red and blue respectively (*P*

881     $< 0.05$). Color intensity and size of the circles are proportional to the correlation

882     coefficients. Abbreviations are as follows: harvest date (HvD), fruit weight (FW), flesh

883     firmness (FF), soluble solids content (SSC), titratable acidity (TA), ripening index (RI),

884     content of vitamin C (Vit C), total phenolics (Phen), anthocyanins (ACNs), sucrose (Suc),

885     glucose (Glu), fructose (Fruc), sorbitol (SRB) total sugars (TS) and relative antioxidant

886     capacity (RAC).

887     **Figure 2**. SNPs density plot and intra-chromosomique linkage disequilibrium decay. (**A**):

888     SNPs density across the eight peach chromosomes. The horizontal axis shows the

889     chromosome length in (Mbp) and the different colors reveal the SNP density per window of

890     1 Mbp. Underlined numbers correspond to the total number of polymorphic sites per

891     chromosome. The asterisks highlight the putative position of centromeres predicted as

892     follows: Chr 1=NC_034009.1: (~21 Mbp), Chr 2=NC_034010.1: (~8 Mbp), Chr 3

893     =NC_034011.1: (~12 Mbp), Chr 4=NC_034012.1: (~24 Mbp), Chr 5=NC_034013.1: (~7

894     Mbp), Chr 6=NC_034014.1: (~15 Mbp), Chr 7=NC_034015.1: (~7 Mbp) and Chr

895     8=NC_034015.1: (~10 Mbp). (**B**): chromosome wide LD decay of $r^2$ (y-axis) over the

896     physical distance in Mbp (x-axis). Each colored line represents a smoothed r2 for all

897     marker pairs on each chromosome. The horizontal dashed red line indicates a cut-off

898     $r^2$=0.25.

899     **Figure 3.** Q-Qplot comparison between the GWAS models implemented in GAPIT:

900     General Linear Model (GLM), Mixed Linear Model (MLM), Compressed MLM (CMLM),

901     Settlement of MLM under Progressively Exclusive Relationship (SUPER), Multiple Loci

902     Mixed Linear Model (MLMM), Fixed and random model Circulating Probability

903    Unification (FarmCPU) and Bayesian-information and Linkage-disequilibrium Iteratively

904    nested keyway (BLINK). Note that MLM and CMLM models are overlaid. For each SNP,

905    the expected -$\log_{10}$ transformed *P*-value (x-axis) is plotted against the -$\log_{10}$ the observed

906    *P*-value (y-axis). The red dashed diagonal line corresponds to the expected Q-Q trendline

907    under the null hypothesis (no association with the phenotype). Larger size dots refer to

908    SNPs statistically associated with a trait. For clarity, only phenotypic traits with significant

909    associations were represented.

910    **Figure 4**. Genome Wide Association and LD block analysis for harvest date (HvD). (**A**):

911    Circular Manhattan plot and association signals based on Blink model. Black dashed

912    circular line corresponds to the Bonferroni adjusted threshold (-$\log_{10}(P)$=5.42). Red and

913    large size dots correspond to statistically associated SNPs. Degradation from blue to red

914    indicates the SNP density per 1 Mbp window across peach chromosomes. (**B**): Locus-

915    specific Manhattan plot (upper panel) and LD heatmap (bottom panel) within 250 Kbp on

916    either side of the lead SNP (SNC_034013.1_13023165). The prime candidate gene is

917    represented as a blue box which in this case contains a single coding exon, where blue

918    fragment refers to the exon. Pairwise LD measurements are displayed as D' values with a

919    color transition from yellow to red. (**C**): Boxplot depicting allelic effect of significant SNP

920    on trait variation. Herein we highlight the SNP commonly affected Harvest date and fruit

921    firmness. Mean value for each genotype is indicated by red circle and ** indicates

922    significant pairwise comparison calculated by Games Howel test ($P \leq 0.05$). (**D**): Genomic

923    distribution of significant ddRAD-derived SNPs (red), reviewed QTLs in the literature

924    (blue) and 9K array derived SNPs (green).

925    **Figure 5**. Genome Wide Association and LD block analysis for fruit weight (FW). (**A**):

926    Circular Manhattan plot and association signals based on Blink model. Black dashed

927    circular line corresponds to the Bonferroni adjusted threshold (-$\log_{10}(P)$=5.42). Red and

928    large size dots correspond to statistically associated SNPs. Degradation from blue to red

929    indicates the SNP density per 1 Mbp window across peach chromosomes. (**B**): Locus-

930    specific Manhattan plot (upper panel) and LD heatmap (bottom panel) within 250 Kbp on

931    either side of the lead SNP. Pairwise LD measurements are displayed as D' values with a

932    color transition from yellow to red. (**C**): Boxplot depicting allelic effect of lead SNP on trait

933 variation. Mean value for each genotype is indicated by red circle and ** indicates

934 significant pairwise comparisons calculated by Games Howel test ($P \leq 0.05$). (**D**): Genomic

935 distribution of significant ddRAD-derived SNPs (red) and reviewed QTLs in the literature

936 (blue).

937 **Figure 6**. Genome Wide Association and LD block analysis for flesh firmness (FF). (**A**):

938 Circular Manhattan plot and association signals based on Blink model. Black dashed

939 circular line corresponds to the Bonferroni adjusted threshold ($-\log_{10}(P)=5.42$). Red and

940 large size dots correspond to statistically associated SNPs. Degradation from blue to red

941 indicates the SNP density per 1 Mbp window across peach chromosomes. (**B**): Locus-

942 specific Manhattan plot (upper panel) and LD heatmap (bottom panel) within 250 Kbp on

943 either side of the lead SNP. Pairwise LD measurements are displayed as D' values with a

944 color transition from yellow to red. (**C**): Boxplot depicting allelic effect of lead SNP on trait

945 variation. Mean value for each genotype is indicated by red circle and ** indicates

946 significant pairwise comparisons calculated by Games Howel test ($P \leq 0.05$). (**D**): Genomic

947 distribution of significant ddRAD-derived SNPs (red) and reviewed QTLs in the literature

948 (blue)

949

950 **Tables**

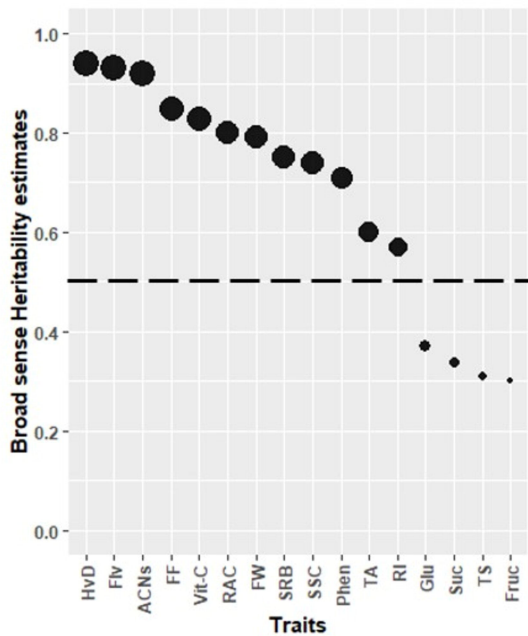951 **Table 1.** SNPs count and filtering steps.

| Applied filters | Retained SNPs |
|---|---|
| Clean multi-samples SNPs from GATK-HaplotypeCaller | 233,535 |
| Clean multi-samples SNPs from Freebayes | 166,080 |
| Clean multi-samples SNPs from BCFtools | 148,998 |
| Intersected set | 56,430 |
| Removing scaffold SNPs | 56,647 |
| Removing multi-allelic sites | 56,430 |
| Missing call rate < 20% | 26,188 |
| Minor Allele Frecuency > 0.05 | 13,045 |

952

953

954 **Table 2**. Information on significantly associated SNP markers with fruit-related traits in

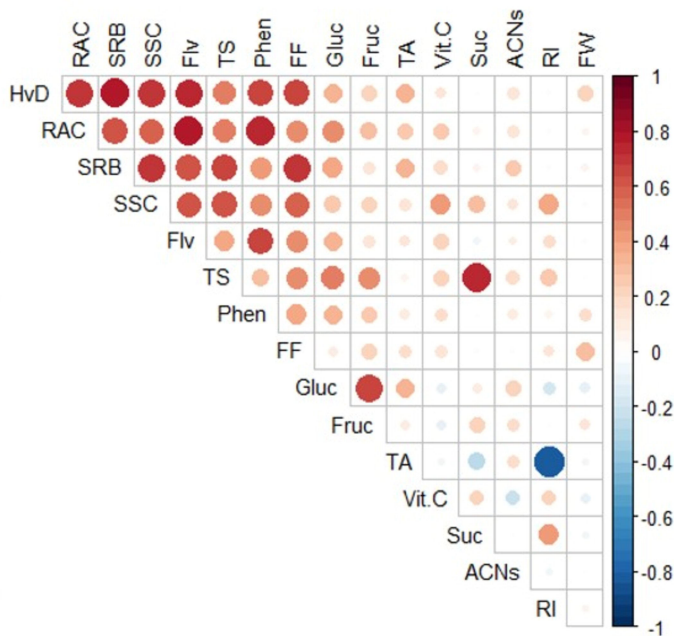955 *Prunus persica.* Alleles are shown on the forward strand as reference/alternate.

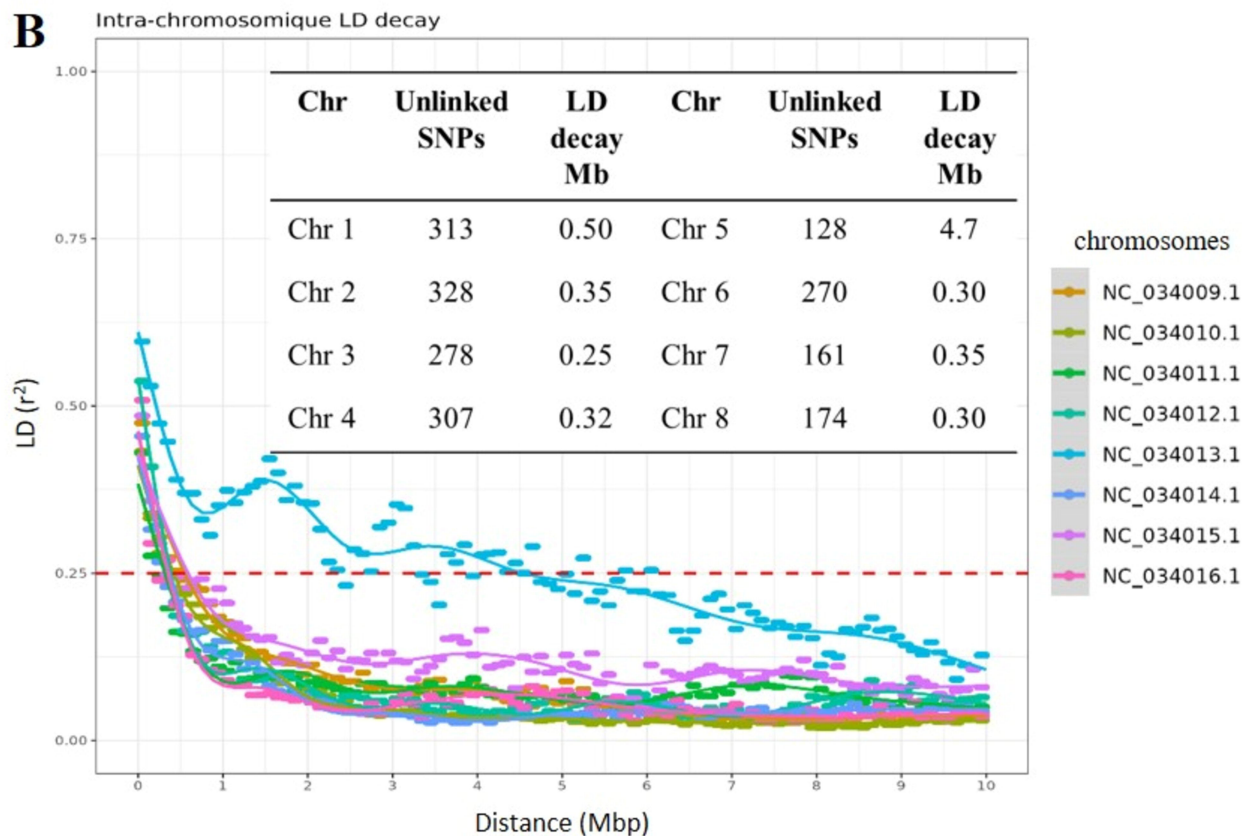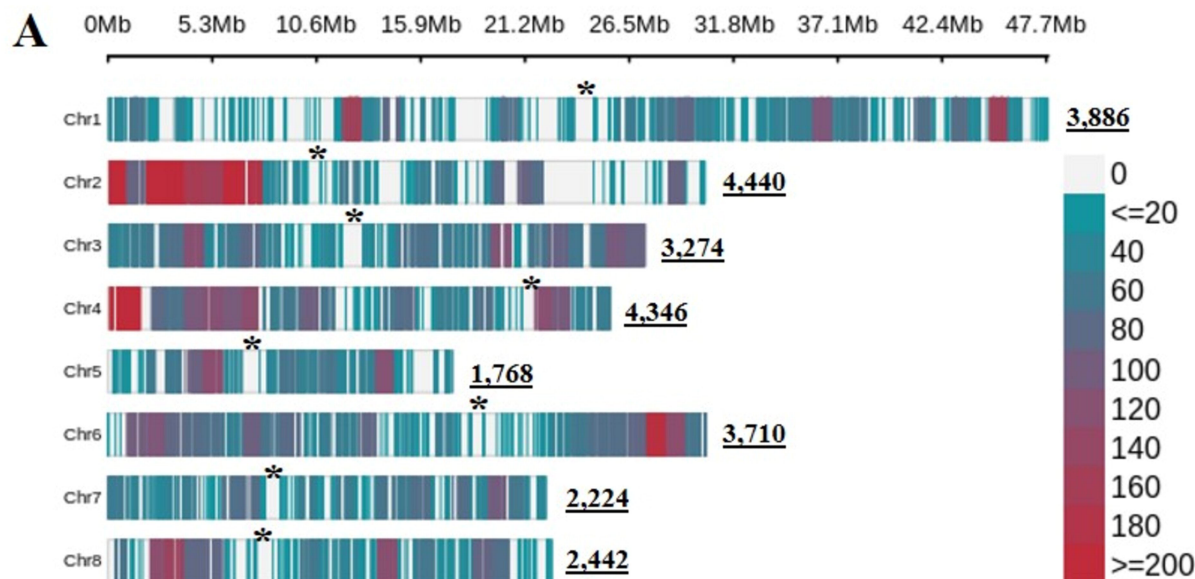| Traits | SNP identifier | Alleles | Chr | Position | %PVE | SNP location [effect] |
|--------|----------------|---------|-----|----------|------|-----------------------|
| | SNC_034012.1_10916234 | G/T | 4 | 10,916,234 | 10.7 | intergenic |
| | SNC_034012.1_14096987 | A/C | 4 | 14,096,987 | 24.5 | intronic |
| HvD | **SNC_034013.1_13023165** | T/A | 5 | 13,023,165 | 30.0 | exonic |
| | SNC_034014.1_7012470 | A/T | 6 | 7,012,470 | 2.8 | intergenic |
| | SNC_034016.1_18841611 | A/G | 8 | 18,841,611 | 10.2 | intergenic |
| | SNC_034011.1_26371177 | T/A | 3 | 26,371,177 | 16.9 | exonic |
| FW | **SNC_034014.1_1805059** | A/G | 6 | 1,805,059 | 22.0 | intergenic |
| | SNC_034016.1_16407694 | A/C | 8 | 16,407,694 | 18.7 | exonic |
| FF | **SNC_034014.1_7012470** | A/T | 6 | 7,012,470 | 33.9 | intergenic |
| | **SNC_034010.1_643430** | T/C | 2 | 643,430 | 35.7 | intergenic |
| FLVs | SNC_034014.1_3066620 | G/T | 6 | 3,066,620 | 14.5 | exonic [missense] |
| ACNs | **SNC_034013.1_12838635** | G/T | 5 | 12,838,635 | 52.9 | exonic [missense] |
| | SNC_034009.1_27061825 | T/C | 1 | 27,061,825 | 9.0 | exonic [missense] |
| SRB | SNC_034010.1_3682553 | G/C | 2 | 3,682,553 | 11.8 | intronic |
| | SNC_034014.1_28343678 | G/A | 6 | 28,343,678 | 10.4 | intronic |
| | **SNC_034016.1_18841643** | G/A | 8 | 18,841,643 | 14.0 | intergenic |

34

956    Variant in bold refers to 'lead SNP', explaining the highest proportion of phenotypic

957    variance (PVE). Chromosome (Chr), Harvest date (HvD), fruit weight (FW), flesh firmness

958    (FF), and contents of flavonoids (FLVs), anthocyanins (ACNs) and sorbitol (SRB).
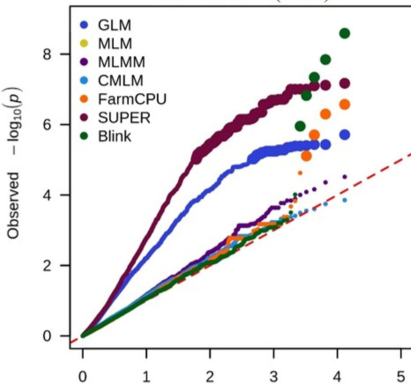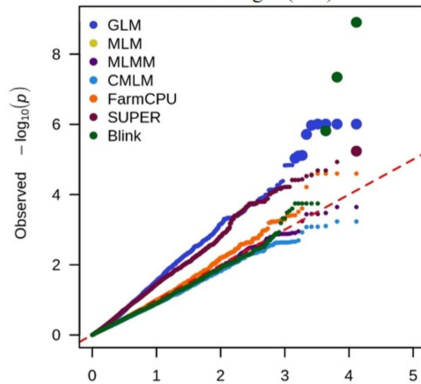
**A**

Chr1 — 3,886
Chr2 — 4,440
Chr3 — 3,274
Chr4 — 4,346
Chr5 — 1,768
Chr6 — 3,710
Chr7 — 2,224
Chr8 — 2,442

Scale: 0, <=20, 40, 60, 80, 100, 120, 140, 160, 180, >=200

**B** Intra-chromosomique LD decay

| Chr | Unlinked SNPs | LD decay Mb | Chr | Unlinked SNPs | LD decay Mb |
|---|---|---|---|---|---|
| Chr 1 | 313 | 0.50 | Chr 5 | 128 | 4.7 |
| Chr 2 | 328 | 0.35 | Chr 6 | 270 | 0.30 |
| Chr 3 | 278 | 0.25 | Chr 7 | 161 | 0.35 |
| Chr 4 | 307 | 0.32 | Chr 8 | 174 | 0.30 |

chromosomes

NC_034009.1
NC_034010.1
NC_034011.1
NC_034012.1
NC_034013.1
NC_034014.1
NC_034015.1
NC_034016.1

**A** Circular Manhattan plot showing GWAS results across Chr1–Chr8 with color scale from 0 to ≥200.

**B** Regional association plot for SNC_034014.1_7012470 (6 Region 394.42 kb, Start 6.800 – End 7.194 Mb) with D' color key (0.0–1.0, NA) and Block 3 (12 SNPs; 123.3 Kbp) LD heatmap.

**C** Boxplot of phenotypes by genotype for SNC_034014.1_7012470:
- A A (n = 34), $\hat{\mu}_{mean} = 32.84$
- A T (n = 31), $\hat{\mu}_{mean} = 39.05$
- T T (n = 25), $\hat{\mu}_{mean} = 46.54$

$p_{FDR-adj} = 4.61e\text{-}03$ **, $p_{FDR-adj} = 1.03e\text{-}06$ **, $p_{FDR-adj} = 4.61e\text{-}03$ **

**D** Chr 6 map showing SNC_034014.1_7012470 and qP-FLP5D6 positions from 5Mbp to 30Mbp.

Legend:
- 9K SNPs (green)
- QTLs (blue)
- ddRAD-seq SNPs (red)