

1 **MultiGWAS: A tool for GWAS analysis on**
2 **tetraploid organisms by integrating the results**
3 **of four GWAS software**

4 L. Garreta¹, I. Cerón-Souza¹, M.R. Palacio², and P.H. Reyes-Herrera¹

5 ¹Corporación Colombiana de Investigación Agropecuaria
6 (AGROSAVIA), CI Tibaitatá, Kilómetro 14, Vía a Mosquera, 250047,
7 Colombia

8 ²Corporación Colombiana de Investigación Agropecuaria
9 (AGROSAVIA), CI El Mira, Kilómetro 38, Vía Tumaco Pasto,
10 Colombia

11 June 10, 2020

12 **Abstract**

13 **Summary:** The Genome-Wide Association Studies (GWAS) are essential to
14 determine the association between genetic variants across individuals. One way
15 to support the results is by using different tools to validate the reproducibility of
16 the associations. Currently, software for GWAS in diploids is well-established
17 but for polyploids species is scarce. Each GWAS software has its characteris-
18 tics, which can cost time and effort to use them successfully. Here, we present
19 MultiGWAS, a tool to do GWAS analysis in tetraploid organisms by executing
20 in parallel and integrating the results from four existing GWAS software: two
21 available for polyploids (GWASpoly and SHEsis) and two frequently used for
22 diploids (PLINK and TASSEL). The tool deals with all the elements of the GWAS
23 process in the four software, including (1) the use of different control quality
24 filters for the genomic data, (2) the execution of two GWAS models, the full
25 model with control for population structure and individual relatedness and the
26 Naive model without any control. The summary report generated by MultiG-
27 WAS provides the user with tables and plots describing intuitively the significant
28 association found by both each one and across four software, which helps users
29 to check for false-positive or false-negative results.

30
31 MultiGWAS generates five summary results integrating the four tools. (1)
32 Score tables with detailed information on the associations for each tool. (2)
33 Venn diagrams of shared SNPs among the four tools. (3) Heatmaps of signifi-
34 cative SNP profiles among the four tools. (4) Manhattan and QQ plots for the
35 association found by each tool. And (5) Chord diagrams for the chromosomes

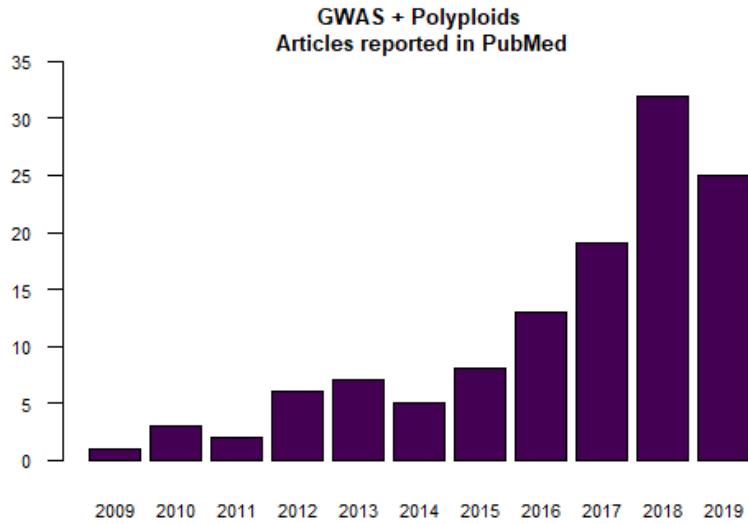


Figure 1: Timeline for articles reported for GWAS studies on polyploid species in PubMed. We present data for completed years.

36 vs. SNP by each tool. **Contact:** phreyes@agrosavia.co

37

38 **Keywords:** GWAS, tetraploids, SNPs,XXX

39 **1 Introduction**

40 The Genome-Wide Association Study (GWAS) is used to identify which variants
41 through the whole genome of a large number of individuals are associated with
42 a specific trait [6, 2]. This methodology started with humans and several model
43 plants, such as rice, maize, and *Arabidopsis* [20, 30, 7, 19, 15]. Because of the
44 advances in the next-gen sequencing technology and the decline of the sequencing
45 cost in recent years, there is an increase in the availability of genome sequences of
46 different organisms at a faster rate [10, 11]. Thus, the GWAS is becoming the stan-
47 dard tool to understand the genetic bases of either ecological or economic pheno-
48 typic variation for both model and non-model organisms. This increment in GWAS
49 includes complex species such as polyploids (Fig 1) [10, 26].

50 The GWAS for polyploid species has three related challenges. First, as all GWAS,
51 we should replicate the study as a reliable method to validate the results and recog-
52 nize real associations. This replication involves finding the same associations either
53 in several replicates from the study population using the same software or testing
54 different GWAS tools among the same study population. This approach involved

55 the use of different parameters, models, or conditions, to test how consistent the
56 results are [9, 17]. However, the performance of different GWAS software could
57 affect the results. For example, the threshold *pvalue* for SNP significance change
58 through four GWAS software (i.e., PLINK, TASSEL, GAPIT, and FaST-LMM) when
59 sample size varies [31]. It means that well-ranked SNPs from one package can be
60 ranked differently in another.

61 Second, although there are many GWAS software available to repeat the anal-
62 ysis under different conditions [14], most of them are designed exclusively for the
63 diploid data matrix [4]. Therefore, it is often necessary to "diploidizing" the poly-
64 ploid genomic data in order to replicate the analysis.

65 Third, there are very few tools focused on the integration of several GWAS soft-
66 ware, to make comparisons under different parameters and conditions across them.
67 As far as we know, there is only two software with this service in mind, such as iPAT
68 and easyGWAS.

69 The iPAT allows running in a graphic interface three well-known command-
70 line GWAS software such as GAPIT, PLINK, and FarmCPU (Chen and Zhang, 2018).
71 However, the output from each package is separated. On the other hand, the easyG-
72 WAS allows running a GWAS analysis on the web using different algorithms. This
73 analysis could run independently of both the computer capacity and operating sys-
74 tem. However, it needs either several datasets available or a dataset with a large
75 number of individuals to make replicates in order to compare among algorithms.
76 Moreover, the output from different algorithms is separated [13]. Thus, for both
77 software iPAT and easyGWAS, the integrative and comparative outputs among soft-
78 ware or algorithms are missing.

79 To solve all the three challenges above, we developed the MultiGWAS tool that
80 performs GWAS analyses for tetraploid species using four software in parallel. Our
81 tool include GWASpoly [25] and the SHEsis tool [28] that accept polyploid genomic
82 data, and PLINK [23] and TASSEL [5] with the use of a "diploidized" genomic ma-
83 trix. The tool deals with preprocessing data, running four GWAS tools in parallel,
84 and create comparative reports from the output of each software to help the user
85 to decide more intuitively the true or false associations.

86 2 Method

87 The MultiGWAS tool has three main consecutive steps: the adjustment, the multi
88 analysis, and the integration (Fig. 2). In the adjustment step, MultiGWAS processes
89 the configuration file. Then it cleans and filters the genotype and phenotype, and
90 MultiGWAS "diploidize" the genomic data. Next, during the multi analysis, each
91 GWAS tool runs in parallel. Subsequently, in the integration step, the MultiGWAS
92 tool scans the output files from the four packages (i.e., GWASPoly, SHEsis, PLink,
93 and TASSEL). Finally, it generates a summary of all results that contains score tables,
94 Venn diagrams, SNP profiles, and Manhattan plots.

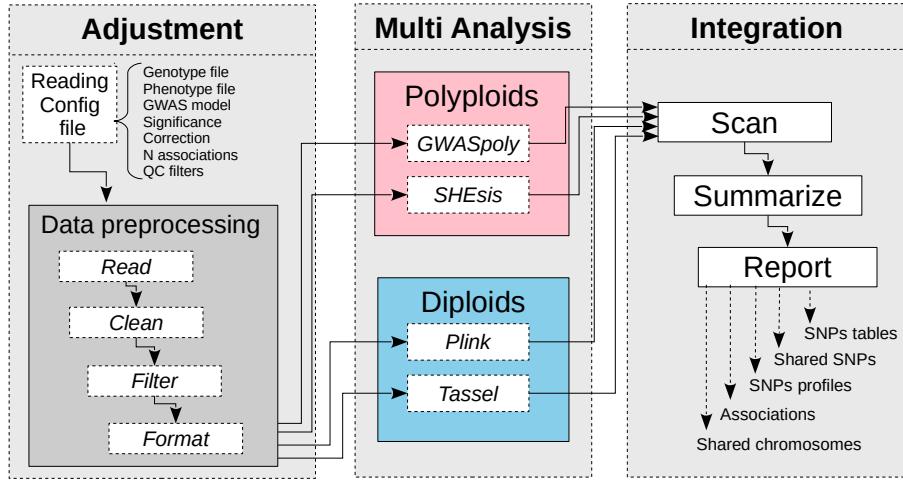


Figure 2: MultiGWAS flowchart has three consecutive steps: adjustment, multi analysis, and integration. The adjustment step manages the input data, reads the configuration file, and preprocesses the input genomic data (genotype and phenotype). The multi analysis step configures and runs the four GWAS packages in parallel. The integration step summarizes and reports results using different tabular and graphical visualizations.

95 2.1 Adjustment stage

96 MultiGWAS takes as input a configuration file where the user specifies the genomics
 97 data along with the parameters that will be used by the four tools. Once the configura-
 98 tion file is processed, MultiGWAS preprocess the data that is cleaning, filtering,
 99 and checking data quality. The output of this stage corresponds to the inputs for
 100 the four programs at the Multi Analysis stage.

101 **2.1.1 Reading configuration file**

102 The configuration file includes the following settings that we briefly describe:

103 **Input genotype and phenotype files:** Currently, MultiGWAS uses two input files,
 104 one for genotype and the other for the phenotype. Both data correspond to data
 105 matrices with column and row names (Figure 3). The genotype file uses SNP mark-
 106 ers in rows and samples in columns (Figure 3a). The phenotype file uses samples in
 107 rows and traits in columns (Figure 3b) with the first column corresponding to the
 108 sample name and the second column to trait value.

Marker,Chrom,Pos,Indiv01,Indiv02,Indiv03,...	Individual,Traitname Indiv01, 3.59 Indiv02, 4.07 Indiv03, 1.05 ...
a	b

Figure 3: MultiGWAS genotype and phenotype formats. Both files are in CSV format (Comma Separated Values) and contain as first row the header labels of the columns. Although the header labels are arbitrary, the column order is obligatory. **a.** Genotype file format, where “Marker”, “Chrom”, and “Pos”, correspond to the names for marker name, chromosome, and position in the first three columns respectively. The next columns names correspond to the individual names and the column content correspond to the genotype of each individual. **b.** Phenotype file format, where “Individual” and “Traitname” are the column for the individual ID and the column for the numerical value of the trait, respectively.

109 **GWAS model:** MultiGWAS is designed to work with quantitative phenotypes and
110 can run GWAS analysis using two types of statistical models that we have called *full*
111 and *naive* models. The *full model* is known in the literature as the Q+K model [32]
112 and includes a control for structure (Q) and relatedness between samples (K). In
113 contrast, the *naive model* does not include any correction. Both models are linear
114 regression approaches, and the four GWAS packages used by MultiGWAS imple-
115 mented variations of them. The *naive* is modeled with Generalized Linear Models
116 (GLMs, Phenotype + Genotype), and the *full* is modeled with Mixed Linear Models
117 (MLMs, Phenotype + Genotype + Structure + Kinship). The default model used by
118 MultiGWAS is the *full model* (Q+K) [32], following this equation:

$$y = X\beta + S\alpha + Q\nu + Z\mu + e$$

119 The vector y represents the observed phenotypes depends on the following fac-
120 tors: the fixed effect vector β , the SNP effects vector α , the population effect vector
121 ν , the polygene background effect vector μ , and, the residual effect vector e . The
122 Q , modeled as a fixed effect, refers to the incidence matrix for subpopulation co-
123 variates relating y to ν . Moreover, X , S , and Z are incidence matrices of ones and
124 zeros relating y to β , α , and μ , respectively.

125 **Genome-wide significance:** GWAS searches SNPs associated with a phenotype
126 trait in a statistically significant manner. A threshold or significance level α is spec-
127 ified and compared with the *p-value* derived for each association score. Standard
128 significance levels are 0.01 or 0.05 [14, 25], and MultiGWAS uses an α of 0.05 for
129 the four GWAS packages. However, the adjustment of the threshold is according
130 to each package. For example, GWASpoly and TASSEL calculate the SNP effect for
131 each genotypic class using different gene action models (see “Multi analysis stage”).
132 Therefore, the number of tested markers may be different in each model (see below)
133 that results in different *p-value* thresholds.

134 **Multiple testing correction:** Due to the massive number of statistical tests per-
135 formed by GWAS, it is necessary to perform a correction method for multiple hy-
136 pothesis testing and adjusting the *p-value* threshold accordingly. Two standard
137 methods for multiple hypothesis testing are the false discovery rate (FDR) and the

138 Bonferroni correction. The latter is the default method used by MultiGWAS because
139 it is one of the most rigorous. MultiGWAS adjust the threshold below which a *p*-
140 value is considered significant, that is α/m , where α is the significance level and m
141 is the number of tested markers from the genotype matrix.

142 **Number of reported associations:** Criticism has arisen in considering only sta-
143 tistically significant associations as the only possible correct associations [29, 18].
144 Many low *p-value* associations are closer to being significant, are discarded due to
145 the stringent significance levels, and, consequently, increase the number of false
146 negatives. To help to analyze both significant and non-significant associations,
147 MultiGWAS provides the option to specify the number of best-ranked associations
148 (lower *p-values*), adding the corresponding *p-value* to each association found. In this
149 way, it is possible to enlarge the number of results, and we can observe replicabil-
150 ity in the results for different programs. Nevertheless, MultiGWAS always presents
151 each associated SNP with its corresponding *p-value*.

152 **Quality control filters:** A control step is necessary to check the input data for
153 genotype or phenotype errors or poor quality that can lead to spurious GWAS re-
154 sults. MultiGWAS provides the option to select and define thresholds for the fol-
155 lowing filters that control the data quality: Minor Allele Frequency (MAF), individ-
156 ual missing rate (MIND), SNP missing rate (GENO), and HardyWeinberg threshold
157 (HWE):

- 158 • **MAF of x :** filters out SNPs with minor allele frequency below x (default 0.01);
- 159 • **MIND of x :** filters out all individuals with missing genotypes exceeding $x*100\%$
160 (default 0.1);
- 161 • **GENO of x :** filters out SNPs with missing values exceeding $x*100\%$ (default
162 0.1);
- 163 • **HWE of x :** filters out SNPs which have Hardy-Weinberg equilibrium exact test
164 *p-value* below the x threshold.

165 MultiGWAS does the MAF filtering and uses the PLINK package [14] for the other
166 three filters: MIND, GENO, and HWE.

167 2.1.2 Data preprocessing

168 Once the configuration file is processed, the genomic data is read and cleaned by
169 selecting individuals present in both genotype and phenotype. Then, based on pre-
170 vious selected quality-control filters and their thresholds, MultiGWAS remove indi-
171 viduals and SNPs with poor quality.

172 During this step, the format "ACGT" suitable for the polyploid software GWASpoly
173 and SHEsis, is "diploidized" for PLINK and TASSEL. The homozygous tetraploid
174 genotypes are converted to diploid thus: AAAA→AA, CCCC→CC, GGGG→GG, TTTT→TT.

175 Moreover, for tetraploid heterozygous genotypes, the conversion depends on the
176 reference and alternate alleles calculated for each position (e.g., AAAT→AT, ... ,CCCG→CG).
177 After this process, MultiGWAS convert the genotype and phenotype data to the
178 specific formats required for each of the four GWAS packages.

179 2.2 Multi analysis stage

180 MultiGWAS runs in parallel using two types of statistical models specified in the pa-
181 rameters file, the Full model (Q+K) and Naive (i.e., without any control) where Q
182 refers to population structure and K refers to relatedness, calculated by kinship co-
183 efficients across individuals [27]. The Full model (Q+K) controls for both population
184 structure and individual relatedness. For population structure, MultiGWAS uses the
185 Principal Component Analysis (PCA) and takes the top **five** PC as covariates. For re-
186 latedness, **MultiGWAS** uses kinship matrices that TASSEL and GWASpoly calculated
187 separately, and for PLINK and SHEsis, **relatedness depends on kinship coefficients**
188 **calculated with the PLINK 2.0 built-in algorithm** [8].

189 2.2.1 GWASpoly

190 GWASpoly [25] is an R package designed for GWAS in polyploid species used in
191 several studies in plants [3, 12, 27, 33]. GWASpoly uses a Q+K linear mixed model
192 with biallelic SNPs that account for population structure and relatedness. **Also, to**
193 **calculate the SNP effect for each genotypic class, GWASpoly provides eight gene**
194 **action models: general, additive, simplex dominant alternative, simplex dominant**
195 **reference, duplex dominant alternative, duplex dominant, diplo-general, and diplo-**
196 **additive. As a consequence, the number of statistical test performed can be different**
197 **in each action model and so thresholds below which the p-values are considered**
198 **significant.**

199 MultiGWAS is using GWASpoly version 1.3 **with all gene action models avail-**
200 **able to find associations. The MultiGWAS reports the top N best-ranked (the SNPs**
201 **with lowest p-values) that the user specified in the Ninput configuration file.** The *full*
202 model used by GWASpoly includes the population structure and relatedness, which
203 are estimated using the first five principal components and the kinship matrix, re-
204 spectively, both calculated with the GWASpoly built-in algorithms.

205 2.2.2 SHEsis

206 SHEsis is a program based on a linear regression model that includes single-locus
207 association analysis, among others. The software design includes polyploid species.
208 However, their use is mainly in diploids animals and humans [24, 21].

209 MultiGWAS is using version 1.0, which does not take account for population
210 structure or relatedness. Despite, MultiGWAS externally estimates relatedness for
211 SHEsis by excluding individuals with cryptic first-degree relatedness using the al-
212 gorithm implemented in PLINK 2.0 (see below).

213 **2.2.3 PLINK**

214 PLINK is one of the most extensively used programs for GWAS in humans and any
215 diploid species [22]. PLINK includes a range of analyses, including univariate GWAS
216 using two-sample tests and linear regression models.

217 MultiGWAS is using two versions of PLINK: 1.9 and 2.0. Linear regression from
218 PLINK 1.9 performs both naive and full model. For the full model, the software
219 calculates the population structure using the first five principal components calcu-
220 lated with a built-in algorithm integrated into version 1.9. Moreover, version 2.0
221 calculates the kinship coefficients across individuals using a built-in algorithm that
222 removes the close individuals with first-degree relatedness.

223 **2.2.4 TASSEL**

224 TASSEL is another standard GWAS program based on the Java software developed
225 initially for maize but currently used in several species [1, 34]. For the associa-
226 tion analysis, TASSEL includes the general linear model (GLM) and mixed linear
227 model (MLM) that accounts for population structure and relatedness. Moreover, as
228 [GWASPoly](#), TASSEL provides three-gene action models to calculate the SNP effect
229 of each genotypic class: general, additive, and dominant, and so the significance
230 threshold depends on each action model.

231 MultiGWAS is using TASSEL 5.0, [with all gene action models used to find the N best-ranked associations and reporting the top N best-ranked associations \(SNPs with lowest p-values\)](#). Naive GWAS uses the GLM, and full GWAS uses the MLM
232 with two parameters: population structure that uses the first five principal compo-
233 nents, and relatedness that uses the kinship matrix with centered IBS method, both
234 calculated with the TASSEL built-in algorithms.

237 **2.3 Integration stage.**

238 The outputs resulting from the four GWAS packages are scanned and processed to
239 identify both significant and best-ranked associations with *p-values* lower than and
240 close to a significance threshold, respectively.

241 **2.3.1 Calculation of *p-values* and significance thresholds**

242 GWAS packages compute *p-value* as a measure of association between each SNP
243 and the trait of interest. The real associations are those their *p-value* drops below
244 a predefined significance threshold. However, the four GWAS packages compute
245 differently *p-values* with the consequence to compute them too high or too low. If
246 *p-values* is too high, it would lead to false negatives or SNPs with real associations
247 with the phenotype, but that does not reach the significance threshold. Conversely,
248 if *p-values* are too low, then it would lead to false positives or SNPs with false asso-
249 ciations with the phenotype, but that reaches the significance threshold.

250 To overcome these difficulties, in the case of too high *p-values*, MultiGWAS iden-
251 tifies and reports both significant and best-ranked associations (the ones closest to

252 being statistically significant). Whereas, in the case of too low *p-values*, MultiGWAS provides two methods for adjusting *p-values* and significance threshold: the
253 false discovery rate (FDR) that adjust *p-values*, and the Bonferroni correction, that
254 adjusts the threshold.
255

256 By default, MultiGWAS uses the Bonferroni correction that uses the significance
257 level α/m (defined by the user in the configuration file), and m (the number of
258 tested markers) to adjust the significance threshold in the GWAS study. However,
259 the significance threshold can be different for each GWAS package as some of them
260 use several action models to calculate the SNP effect of each genotypic class. For
261 both PLINK and SHEsis packages, which use only one model, m is equal to the total
262 number of SNPs. However, for both GWASpoly and TASSEL packages, which use
263 eight and three gene action models, respectively, m is equal to the number of tests
264 performed in each model, which is different between models.

265 2.3.2 Selection of significant and best-ranked associations

266 After corrections, significant associations are selected as the ones with *p-values*
267 falling below a significant threshold on each GWAS package. Nevertheless, as de-
268 scribed above, it is equally important to know the best-ranked associations closer
269 to being statistically significant, as they may represent associations to consider for
270 posterior analysis.

271 In the case of GWAS packages with only one gene action model (PLINK and SHE-
272 SIS), the best-ranked associations are the top N identified by the package. However,
273 in GWAS packages with several gene action models (GWASpoly and TASSEL), the
274 best-ranked associations are selected as the top N from the “best action model”,
275 the one with more shared SNP associations. In other words, from the associations
276 identified in more than one model.

277 2.3.3 Integration of results

278 At this stage, [MultiGWAS integrates](#) the results to evaluate reproducible results
279 among tools (Fig 4). However, it still reports a summary of the results of each
280 tool:

- 281 • A Quantile-Quantile (QQ) plots for the resultant *p-values* of each tool and
282 the corresponding [inflation factor](#) λ to assess the degree of the test statistic
283 inflation.
- 284 • A Manhattan plot of each tool with two lower thresholds, one for the best-
285 ranked SNPs, and another for the significant SNPs.

286 To present the replicability, we use two sets: (1) the set of all the significative SNPs
287 provided by each tool and (2) the set of all the best-ranked SNPs. For each set,
288 we present a Venn diagram that shows SNPs predicted exclusively by one tool and
289 intersections that help to identify the SNPs predicted by one, two, three, or all the
290 tools. Also, we provide detailed tables for the two sets.

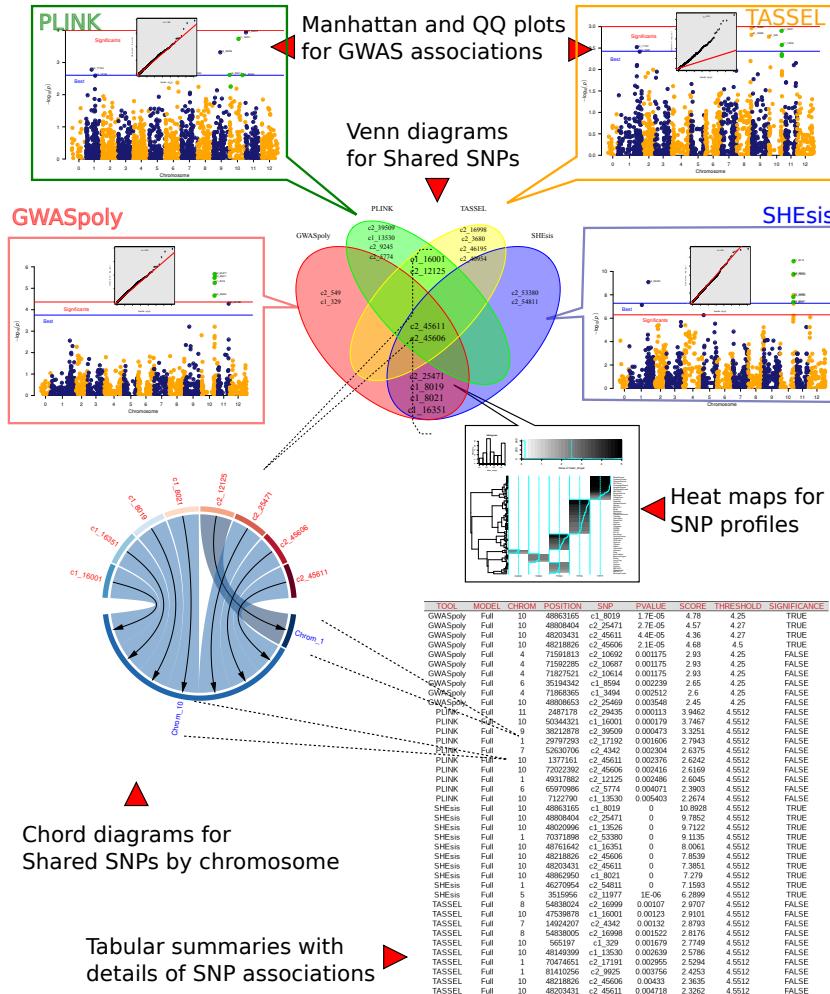


Figure 4: Reports presented by MultiGWAS. For each tool, first a QQ plot that assesses the resultant p-values. Second, a Manhattan plot for each tool with two lines, blue and red, respectively, is the lower limit for the best ranked and significative SNPs. We present two Venn diagrams, one for the significative SNPs and one for N best-ranked SNPs of each tool. We show the results for GWAspoly, PLINK, TASSEL, and SHEsis in red, green, yellow, and blue. For each SNP that is in the intersection, thus, that is predicted by more than one tool, we provide an SNP profile. SNPs by chromosome chord diagrams show that the strongest associations are limited to few chromosomes. Furthermore, we present tabular summaries with details of significant and best-ranked associations.

291 For each SNP identified more than once, we provide what we call the SNP pro-
292 file. That is a heat diagram for a specific SNP, where each column is a genotype
293 state AAAA, AAAB, AABB, ABBB, and BBBB. Moreover, each row corresponds to a
294 sample. Samples with close genotypes form together clusters. Thus to generate
295 the clusters, we do not use the phenotype information. However, we present the
296 phenotype information in the figure as the color. This figure visually provides in-
297 formation regarding genotype and phenotype information simultaneously for the
298 whole population. We present colors as tones between white and black for color
299 blind people.

300 MultiGWAS generates a report, one document with the content previously de-
301 scribed. Besides, there is a folder with the individual figures just in case the user
302 needs one. In the supplementary information, we include a report and a description
303 of the report content (**supplementary information XXX**)

304 In the following section, we present the results applied to a public dataset.

305 3 Results

306 Most of the GWAS packages used by MultiGWAS are based on a linear regression
307 approaches, but they often produce dissimilar association results for the same input.
308 For example, computed *p-values* for the same set of SNPs are different between
309 packages; SNPs with significant *p-values* for one package may be not significant
310 for the others; or well-ranked SNPs in one package may be ranked differently in
311 another.

312 To alleviate these difficulties, MultiGWAS produces five types of outputs using
313 different graphics and tabular views, these outputs are intended to help users to
314 compare, select, and interpret the set of possible SNPs associated with a trait of
315 interest. The outputs include:

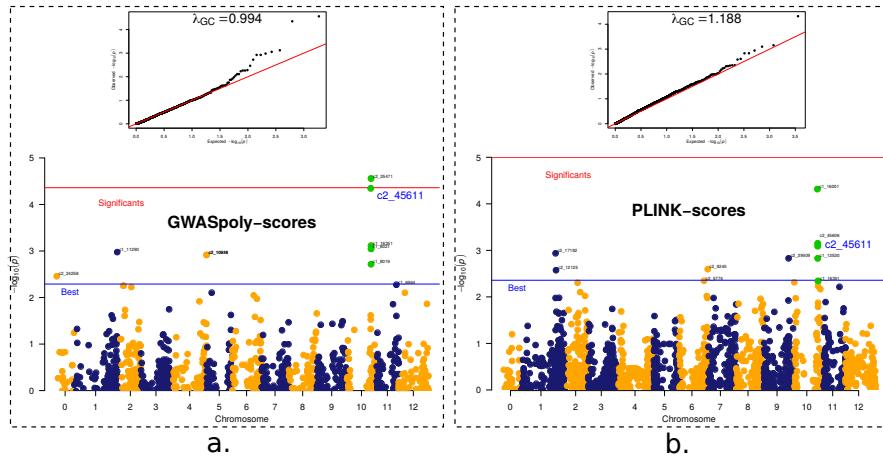
- 316 • Manhattan and Q-Q plots to show GWAS associations.
- 317 • Venn diagrams to show associations identified by single or several tools.
- 318 • Heat diagrams to show the genotypic structure of shared SNPs.
- 319 • Chord diagrams to show shared SNPs by chromosomes.
- 320 • Score tables to show detailed information of associations for both summary
321 results from MultiGWAS and particular results from each GWAS package

322 As an example of the functionality of the tool, here we show the outputs reported by
323 MultiGWAS in the tetraploid potato diversity panel, genotyped and phenotyped as
324 part of the USDA-NIFA Solanaceae Coordinated Agricultural Project (SolCAP) [16].
325 The complete report from MultiGWAS for the naive and full model is in the Supple-
326 mentary information (<https://github.com/agrosavia-bioinformatics/multiGWAS>)

328 3.1 Manhattan and QQ plots for GWAS associations

329 MultiGWAS uses classical Manhattan and Quantile–Quantile plots (QQ plots) to
 330 visualize the results of GWAS analysis from each package. In both plots, SNPs are
 331 represented by dots and their p -values are transformed to scores as $-\log_{10}(p\text{-value})$
 332 (see Figure 5). The Manhattan plot displays the SNP association strength (y-axis)
 333 distributed in their genomic location (x-axis), so the higher the score the stronger
 334 the association. Whereas the QQ plot is used to visually compare the expected
 335 distribution of p -values (y-axis) vs. the observed distribution (x-axis), so under the
 336 null hypothesis of no association of SNPs with the phenotype, both distributions
 337 should coincide, and most SNPs should lie on a diagonal line.
 338

339 MultiGWAS adds special marks to the Manhattan and QQ plots to help identify
 340 different types of SNPs: (a) In Manhattan plots, significant SNPs are above a red
 341 line, best-ranked SNPs are above a blue line, and shared SNPs (See Figure 6.b) are
 342 colored in green (b) In QQ plots, a red diagonal line indicates the expectation, so
 343 potential associations can be observed when the number of SNPs deviating from
 344 the diagonal is small, as in the case of monogenic traits, or when this number is
 345 somewhat higher, as in the case of truly polygenic traits. However, deviations for a
 346 high number of SNPs could reflect inflated p -values owing to population structure
 or cryptic relatedness.



347 **Figure 5: MultiGWAS visualization of associations.** MultiGWAS creates Manhattan and
 348 QQ plots for GWAS results of each GWAS packages. Here we show the plots for one tetraploid
 349 package, GWASpoly (a), and other diploid package, PLINK (b).

347 3.2 Tables and Venn diagrams for single and shared SNPs

348 MultiGWAS provides tabular and graphic views to report in an integrated way both
 349 the best-ranked and significant SNPs identified by the four GWAS packages (see
 350 Figure 6). Both p -values and significance levels have been scaled as $-\log_{10}(p\text{-value})$
 351 to give high scores to the best statistically evaluated SNPs.

352 First, best-ranked SNPs correspond to the top-scored N SNPs, wheter they were
 353 assesed significant or not by its package, and with N defined by the user in the
 354 configuration file. These SNPs are shown both in a SNPs table (Figure 6.a) and in a
 355 Venn diagram (Figure 6.b). The table lists them by package and sorts by decreasing
 356 score, whereas the Venn diagram shows them emphasizing if they were best-ranked
 357 either in a single package or in several at once (shared). And second, the significant
 358 SNPs correspond to the ones assesed statistically significant by each package, they
 359 are shown in a Venn diagram (Figure 6.c), and they are also shown in the SNPs
 360 table, marked with significance TRUE (T) in the table of the Figure6.a.

a.

TOOL	MODEL	GC	SNP	CHR	POS	PVALUE	SCR	THR	SGN
GWASPoly	additive	0.96	c2_25471	10	48808	0.000002	5.67	4.50	T
GWASPoly	additive	0.96	c2_45611	10	48203	0.000003	5.51	4.50	T
GWASPoly	additive	0.96	c1_8019	10	48863	0.000005	5.27	4.50	T
GWASPoly	additive	0.96	c2_45606	10	48218	0.000021	4.68	4.50	T
GWASPoly	additive	0.96	c2_22188	11	40777	0.000050	4.30	4.50	F
GWASPoly	additive	0.96	c2_549	9	16527	0.000589	3.23	4.50	F
GWASPoly	additive	0.96	c1_8021	10	48862	0.000589	3.23	4.50	F
GWASPoly	additive	0.96	c1_329	10	56519	0.001514	2.82	4.50	F
GWASPoly	additive	0.96	c1_16351	10	48761	0.001622	2.79	4.50	F
PLINK	additive	1.19	c1_16001	10	47539	0.000047	4.33	4.55	F
PLINK	additive	1.19	c2_45606	10	48218	0.000688	3.16	4.55	F
PLINK	additive	1.19	c2_45611	10	48203	0.000788	3.10	4.55	F
PLINK	additive	1.19	c2_17192	1	70472	0.001123	2.95	4.55	F
PLINK	additive	1.19	c2_39509	9	50174	0.001440	2.84	4.55	F
PLINK	additive	1.19	c1_13530	10	48149	0.001443	2.84	4.55	F
PLINK	additive	1.19	c2_9245	6	57953	0.002455	2.61	4.55	F
PLINK	additive	1.19	c2_12125	1	71450	0.002593	2.59	4.55	F
PLINK	additive	1.19	c2_5774	6	50345	0.004336	2.36	4.55	F
SHEsis	general	1.47	c1_8019	10	48863	0.000000	7.64	4.55	T
SHEsis	general	1.47	c1_13526	10	48020	0.000000	6.94	4.55	T
SHEsis	general	1.47	c2_25471	10	48808	0.000000	6.94	4.55	T
SHEsis	general	1.47	c2_53380	1	70371	0.000000	6.46	4.55	T
SHEsis	general	1.47	c1_16351	10	48761	0.000004	5.45	4.55	T
SHEsis	general	1.47	c2_45606	10	48218	0.000004	5.38	4.55	T
SHEsis	general	1.47	c2_45611	10	48203	0.000010	4.98	4.55	T
SHEsis	general	1.47	c1_8021	10	48862	0.000012	4.93	4.55	T
SHEsis	general	1.47	c2_54811	1	46270	0.000014	4.86	4.55	T
TASSEL	additive	0.86	c2_16999	8	54838	0.000247	3.61	3.89	F
TASSEL	additive	0.86	c2_16998	8	54838	0.000329	3.48	3.89	F
TASSEL	additive	0.86	c2_12125	1	71450	0.003287	2.48	3.89	F
TASSEL	additive	0.86	c1_16001	10	47539	0.006105	2.21	3.89	F
TASSEL	additive	0.86	c2_3680	11	39908	0.006701	2.17	3.89	F
TASSEL	additive	0.86	c2_46195	1	64259	0.007116	2.15	3.89	F
TASSEL	additive	0.86	c2_40954	1	63756	0.011097	1.95	3.89	F
TASSEL	additive	0.86	c2_45606	10	48218	0.011369	1.94	3.89	F
TASSEL	additive	0.86	c2_45611	10	48203	0.012091	1.92	3.89	F

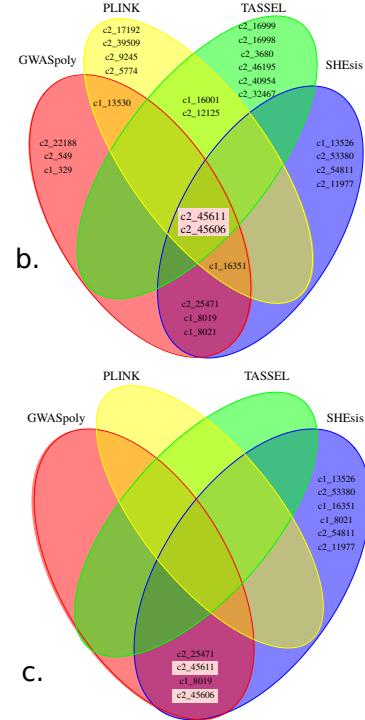


Figure 6: Shared SNPs Views. Tabular and graphical views of SNP associations identified by one or more GWAS packages (shared SNPs). SNPs identified by all packages are marker with red background in all figures (a) Table with details of the $N=9$ best-ranked SNPs from each GWAS package. Each row corresponds to a single SNP and the 9 columns are: tool name, model used by the tool, genomic control factor (inflation factor), SNP name, chromosome, position in the genome, p -value, score as $-\log_{10}(p\text{-value})$, significance threshold as $-\log_{10}(\alpha/m)$ where α is the significance level and m is the number of tested markers, and significance as true (T) or false (F) whether score $>$ threshold or not. (b) Venn diagram of the $N=9$ best-ranked SNPs. SNPs identified by all packages are located in the central intersection. Other SNPs identified by more than one packages are located in both upper central and lower central intersections. (c) Venn diagram of the significant SNPs (score $>$ threshold).

361 3.3 Heat diagrams for structure of shared SNPs

362 MultiGWAS creates a two-dimensional representation, called SNP profile, to visu-
 363 alize each trait by individuals and genotypes as rows and columns, respectively
 364 (Figure 7). At the left, the individuals are grouped in a dendrogram by their geno-
 365 type. At the right, there is the name or ID of each individual. At the bottom, the
 366 genotypes are ordered from left to right, starting from the major to the minor allele
 367 (i.e., AAAA, AAAB, AABB, ABBB, BBBB). At the top, there is a description of the
 368 trait based on a histogram of frequency (top left) and by an assigned color for each
 369 numerical phenotype value using a grayscale (top right). Thus, each individual ap-
 370 pears as a colored line by its phenotype value on its genotype column. For each
 371 column, there is a solid cyan line with the mean of each column and a broken cyan
 372 line that indicates how far the cell deviates from the mean.

373 Because each multiGWAS report shows one specific trait at a time, the histogram
 374 and color key will remain the same for all the best-ranked SNPs.

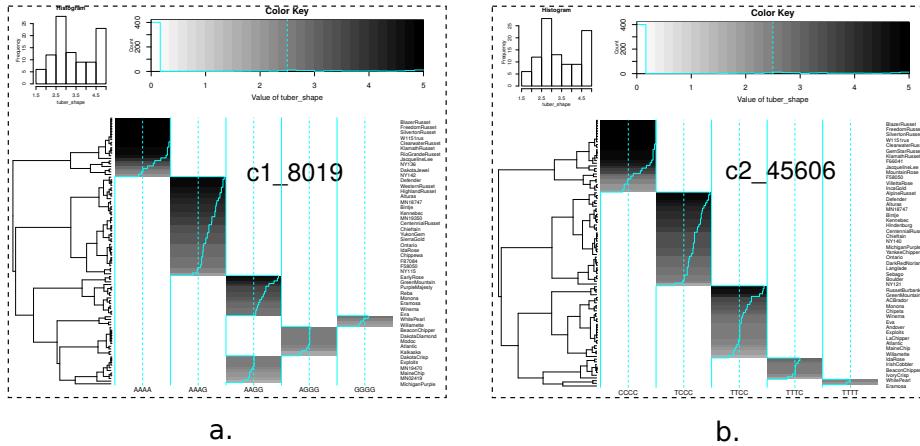


Figure 7: SNP profiles. SNP profiles for two of the best-ranked significant SNPs shown in figure 6.b. (a) SNP c2_45606 best-ranked by the four packages (central intersection of the Venn diagram Figure 6.b) (b) SNP c1_8019 best-ranked by the two tetraploid packages (Figure 6.b), and also identified as significant by the same packages (at the bottom of the Figure 6.a).

375 3.4 Chord diagrams for SNPs by chromosome

376 Generally, in a typical GWAS analysis the strongest associations are signaled by
 377 several nearby-correlated SNPs located in the same chromosome, as in manhattan
 378 plots, where these associations form neat peaks with several SNPs showing the same
 379 signal. Conversely, no peaks are shown when few SNPs correlate with a trait.

380 However, when the analysis is performed by several GWAS packages, as Multi-
 381 GWAS does, it can identify correlated SNPs between packages that show the same
 382 signal, what is presented by MultiGWAS through chord diagrams. For example, the
 383 Figure 8.a shows the chord diagram for the shared SNPs from the best-ranked asso-
 384 ciations previously described in the Figure 6.b. It can be observed that most SNPs

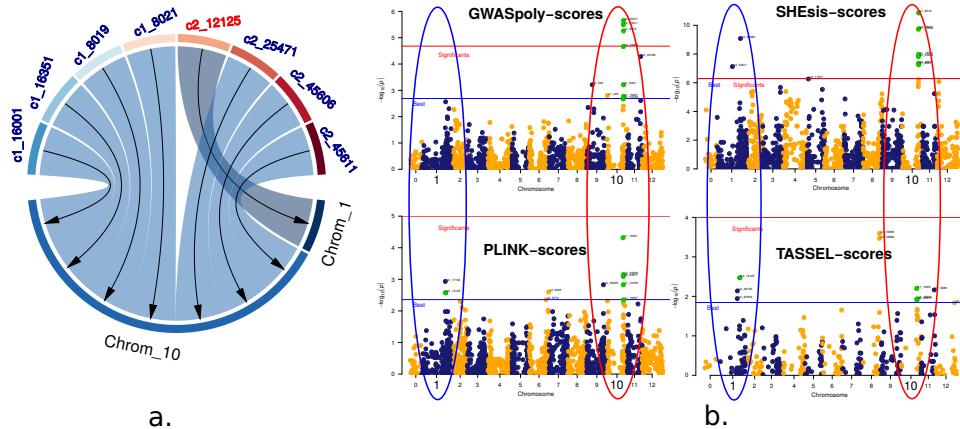


Figure 8: SNPs by chromosome. The figure shows how the best-ranked SNPs relate to chromosomes. (a) Chord diagram showing that most SNPs related to chromosome 10. SNPs are at the top of the diagram, chromosomes at the bottom, and associations are represented by arrows drawn from SNPs to their chromosomes. The more associations identified in one chromosome, the wider the space of its sector. (b) Manhattan plots from each GWAS package showing two important locations of associations: chromosome 1 and chromosome 10, marked with a blue and red ellipsis, respectively.

385 relate to chromosome 10 and only one to chromosome 1, which is also observed in
 386 the manhattan plots from each GWAS package (Figure 8.b).

387 4 Availability and Implementation

388 The core of the MultiGWAS tool was developed in R and users can interact with
 389 the tool by either a command line interface (CLI) developed in R or a graphical
 390 user interface (GUI) developed in Java (Figure 10). Source code, examples, doc-
 391 umentation and installation instructions are available at <https://github.com/>
 392 `agrosavia-bioinformatics/multiGWAS`.

393 4.1 Input parameters

394 MultiGWAS uses as the only input a simple configuration text file where users set
 395 the values for the main parameters that drives the GWAS process. The file can be
 396 created either using a general text editor or using the MultiGWAS GUI application
 397 (see below). In both cases, the file must have the structure shown in the Figure 9.a,
 398 where parameter names and values are separated by colon, filenames are enclosed
 399 in quotation marks, and TRUE or FALSE indicates whether filters are applied or not.
 400 In the second case, the user creates the config file in a simple and straightforward
 401 way using the input parameter view from the GUI application (see below).

```

default:
    genotypeFile      : "example-genotype.tbl"
    phenotypeFile     : "example-phenotype.tbl"
    significanceLevel : 0.05
    correctionMethod  : "Bonferroni"
    gwasModel         : "Full"
    nBest             : 10
    filtering         : TRUE
    MAF               : 0.01
    MIND              : 0.1
    GENO              : 0.1
    HWE               : 1e-10
    tools              : "GWASpoly SHEsis PLINK TASSEL"

```

Figure 9: Configuration file for MultiGWAS. The input parameters include: the output folder where results will be written, input genotype/phenotype filenames, genome-wide significance threshold, method for multiple testing correction, GWAS model, number of associations to be reported, filtering with TRUE or FALSE whether to use quality control filters or not. The filters are: minor allele frequency, individual missing rate, SNP missing rate, and Hardy-Weinberg threshold. At the end the tools parameter defines the GWAS packages to be used for the analysis.

402 4.2 Using the command line interface

403 The execution of the CLI tool is simple, it only needs to open a linux console, change
 404 to the folder where the configuration file was created, and type the name of the
 405 executable tool followed by the filename of the configuration file, like this:

406 multiGWAS Test01.config

407 Then, the tool starts the execution, showing information of the process in the
 408 console window, and when it finishes the results are saved to a new subfolder called
 409 “out-Test01. Results include a full html report containing the different views de-
 410 scribed in the results section, along with the original graphics and summary tables
 411 created by MultiGWAS and used to create the html report. Additionally, results
 412 include the preprocessed tables of the main outputs generated by the four GWAS
 413 packages used by MultiGWAS.

414 4.3 Using the graphical user interface

415 The MultiGWAS GUI can be executed by calling from a linux console the following
 416 command:

417 jmultiGWAS

418 After it opens, it shows a main frame with a tool bar at left and four tabs at the
 419 top (Figure 10). From the tool bar, users can select the GWAS packages to use in
 420 the analysis—two for tetraploids and two for diploids—, and start the analysis with
 421 the current parameters (or with parameters from a previous configuration). And,
 422 from the tabs, users can input the MultiGWAS parameters, and view the process
 423 and results of the analysis.

424 .

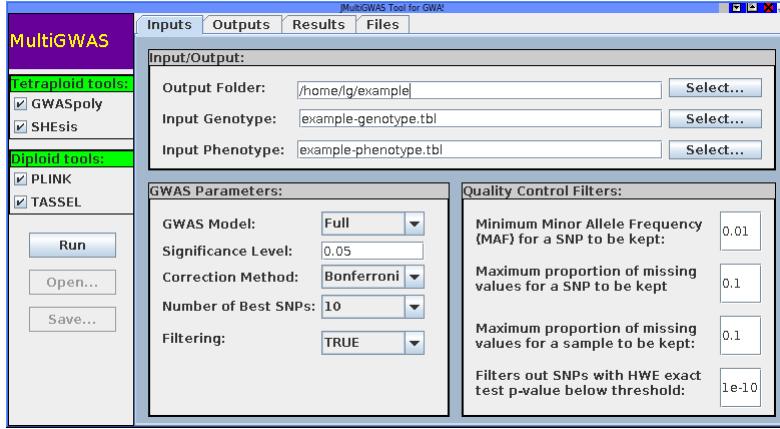


Figure 10: MultiGWAS GUI application. Main view of the MultiGWAS GUI application (“Inputs” view) where users can create the configuration file by setting values for input parameters. The GUI contains other three views: “Outputs” view shows the logs of the running process. “Results” view shows a report in html format with the tabular and graphics described in the results section. And, the “Files” view shows an embedded file manager pointing to the subfolder that contains the files created by MultiGWAS and used to create the report.

425 5 Discussion

426 XXXXXXXXXXXXXXXXXXXXXXXXX

427 References

- 428 [1] María F. Álvarez et al. “Identification of Novel Associations of Candidate
429 Genes with Resistance to Late Blight in Solanum tuberosum Group Phureja”.
430 In: *Frontiers in Plant Science* 8 (2017), p. 1040. ISSN: 1664-462X. DOI: 10.
431 3389/fpls.2017.01040. URL: <http://journal.frontiersin.org/article/10.3389/fpls.2017.01040/full>.
- 433 [2] Ferdouse Begum et al. “Comprehensive literature review and statistical con-
434 siderations for GWAS meta-analysis”. In: *Nucleic acids research* 40.9 (2012),
435 pp. 3777–3784.
- 436 [3] Jhon Berdugo-Cely et al. “Genetic diversity and association mapping in the
437 colombian central collection of solanum tuberosum L. Andigenum group us-
438 ing SNPs markers”. In: *PLoS ONE* 12.3 (2017). ISSN: 19326203. DOI: 10.
439 1371/journal.pone.0173039.
- 440 [4] Peter M. Bourke et al. “Tools for Genetic Studies in Experimental Popula-
441 tions of Polyploids”. In: *Frontiers in Plant Science* 9 (2018), p. 513. ISSN:
442 1664-462X. DOI: 10.3389/fpls.2018.00513. URL: <http://journal.frontiersin.org/article/10.3389/fpls.2018.00513/full>.

- 444 [5] Peter J Bradbury et al. “TASSEL: software for association mapping of complex
445 traits in diverse samples”. In: *Bioinformatics* 23.19 (2007), pp. 2633–2635.
446 ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btm308. URL: <https://doi.org/10.1093/bioinformatics/btm308>.
- 448 [6] Rita M Cantor, Kenneth Lange, and Janet S Sinsheimer. “Prioritizing GWAS
449 results: a review of statistical methods and recommendations for their ap-
450 plication”. In: *The American Journal of Human Genetics* 86.1 (2010), pp. 6–
451 22.
- 452 [7] Jun Cao et al. “Whole-genome sequencing of multiple *Arabidopsis thaliana*
453 populations”. In: *Nature genetics* 43.10 (2011), p. 956.
- 454 [8] Christopher C. Chang et al. “Second-generation PLINK: Rising to the chal-
455 lenge of larger and richer datasets”. In: *GigaScience* 4.1 (2015), pp. 1–16.
456 ISSN: 2047217X. DOI: 10.1186/s13742-015-0047-8. arXiv: 1410.4803.
- 458 [9] Rishika De, William S Bush, and Jason H Moore. “Bioinformatics Challenges
459 in Genome-Wide Association Studies (GWAS)”. In: *Clinical Bioinformatics*.
460 Ed. by Ronald Trent. New York, NY: Springer New York, 2014, pp. 63–81.
461 ISBN: 978-1-4939-0847-9. DOI: 10.1007/978-1-4939-0847-9_5. URL:
462 https://doi.org/10.1007/978-1-4939-0847-9_5.
- 463 [10] Robert Eklblom and Juan Galindo. “Applications of next generation sequenc-
464 ing in molecular ecology of non-model organisms”. In: *Heredity* 107.1 (2011),
465 pp. 1–15.
- 466 [11] Hans Ellegren. “Genome sequencing and population genomics in non-model
467 organisms”. In: *Trends in ecology & evolution* 29.1 (2014), pp. 51–63.
- 468 [12] Luís Felipe V. Ferrão et al. “Insights Into the Genetic Basis of Blueberry Fruit-
469 Related Traits Using Diploid and Polyploid Models in a GWAS Context”. In:
470 *Frontiers in Ecology and Evolution* 6 (2018), p. 107. ISSN: 2296-701X. DOI:
471 10.3389/fevo.2018.00107. URL: <https://www.frontiersin.org/articles/10.3389/fevo.2018.00107/full>.
- 473 [13] Dominik G Grimm et al. “easyGWAS: A Cloud-Based Platform for Comparing
474 the Results of Genome-Wide Association Studies”. In: *The Plant Cell* 29.1
475 (2017), pp. 5–19. ISSN: 1040-4651. DOI: 10.1105/tpc.16.00551. URL:
476 <http://www.plantcell.org/content/29/1/5>.
- 477 [14] Anja C Gumpinger et al. *Methods and Tools in Genome-wide Association Stud-
478 ies*. Vol. 1819. 2018. ISBN: 9781493986187.
- 479 [15] Bin Han and Xuehui Huang. “Sequencing-based genome-wide association
480 study in rice”. In: *Current opinion in plant biology* 16.2 (2013), pp. 133–138.
- 481 [16] Candice N. Hirsch et al. “Retrospective view of North American potato (*Solanum*
482 *tuberosum* L.) breeding in the 20th and 21st centuries”. In: *G3: Genes, Genomes,
483 Genetics* 3.6 (2013), pp. 1003–1013. ISSN: 21601836. DOI: 10.1534/g3.113.005595.

- 485 [17] “How to interpret a genome-wide association study”. In: *JAMA - Journal*
486 *of the American Medical Association* 299.11 (2008), pp. 1335–1344. ISSN:
487 00987484. DOI: 10.1001/jama.299.11.1335.
- 488 [18] Avjinder S Kaler and Larry C Purcell. “Estimation of a significance threshold
489 for genome-wide association studies”. In: *BMC Genomics* 20.1 (2019), p. 618.
490 ISSN: 1471-2164. DOI: 10.1186/s12864-019-5992-7. URL: <https://doi.org/10.1186/s12864-019-5992-7>.
- 492 [19] Arthur Korte and Ashley Farlow. “The advantages and limitations of trait
493 analysis with GWAS: a review”. In: *Plant methods* 9.1 (2013), p. 29.
- 494 [20] Gordan Lauc et al. “Genomics meets glycomics—the first GWAS study of hu-
495 man N-glycome identifies HNF1 α as a master regulator of plasma protein
496 fucosylation”. In: *PLoS genetics* 6.12 (2010).
- 497 [21] Jie Meng et al. “Genome-wide association analysis of nutrient traits in the
498 oyster *Crassostrea gigas*: Genetic effect and interaction network”. In: *BMC*
499 *Genomics* 20.1 (2019), pp. 1–14. ISSN: 14712164. DOI: 10.1186/s12864-
500 019-5971-z.
- 501 [22] Robert A. Power, Julian Parkhill, and Tilio De Oliveira. “Microbial genome-
502 wide association studies: lessons from human GWAS”. In: *Nature Reviews*
503 *Genetics* 18.1 (2016), pp. 41–50. ISSN: 14710064. DOI: 10.1038/nrg.2016.132.
- 505 [23] Shaun Purcell et al. “PLINK: A tool set for whole-genome association and
506 population-based linkage analyses”. In: *American Journal of Human Genetics*
507 81.3 (2007), pp. 559–575. ISSN: 00029297. DOI: 10.1086/519795.
- 508 [24] Hui Ping Qiao et al. “Genetic variants identified by GWAS was associated
509 with colorectal cancer in the Han Chinese population”. In: *Journal of Cancer*
510 *Research and Therapeutics* 11.2 (2015), pp. 468–470. ISSN: 19984138. DOI:
511 10.4103/0973-1482.150346.
- 512 [25] Umesh R. Rosyara et al. “Software for Genome-Wide Association Studies
513 in Autopolyploids and Its Application to Potato”. In: *The Plant Genome* 9.2
514 (2016), pp. 1–10. ISSN: 1940-3372. DOI: 10.3835/plantgenome2015.08.0073.
515 URL: <https://dl.sciencesocieties.org/publications/tpg/abstracts/9/2/plantgenome2015.08.0073>.
- 517 [26] Anna W Santure and Dany Garant. “Wild GWAS—association mapping in
518 natural populations”. In: *Molecular ecology resources* 18.4 (2018), pp. 729–
519 738.
- 520 [27] Sanjeev Kumar Sharma et al. “Linkage disequilibrium and evaluation of genome-
521 wide association mapping models in tetraploid potato”. In: *G3: Genes, Genomes,*
522 *Genetics* 8.10 (2018), pp. 3185–3202. ISSN: 21601836. DOI: 10.1534/g3.118.200377.
- 524 [28] Jiawei Shen et al. “SHEsisPlus, a toolset for genetic studies on polyploid
525 species”. In: *Scientific Reports* 6 (2016), pp. 1–10. ISSN: 20452322. DOI: 10.
526 1038/srep24095. URL: <http://dx.doi.org/10.1038/srep24095>.

- 527 [29] John R Thompson, John Attia, and Coseetta Minelli. “The meta-analysis of
528 genome-wide association studies”. In: *Briefings in Bioinformatics* 12.3 (2011),
529 pp. 259–269. ISSN: 1467-5463. DOI: 10.1093/bib/bbr020. URL: <https://doi.org/10.1093/bib/bbr020>.
- 531 [30] Feng Tian et al. “Genome-wide association study of leaf architecture in the
532 maize nested association mapping population”. In: *Nature genetics* 43.2 (2011),
533 pp. 159–162.
- 534 [31] Yan Y. Yan et al. “Effects of input data quantity on genome-wide association
535 studies (GWAS)”. In: *International Journal of Data Mining and Bioinformatics*
536 22.1 (2019), pp. 19–43. ISSN: 17485681. DOI: 10.1504/IJDMB.2019.
537 099286.
- 538 [32] J Yu et al. “A unified mixed-model method for association mapping that ac-
539 counts for multiple levels of relatedness.” In: *Nature genetics* 38.2 (2006),
540 pp. 203–208.
- 541 [33] Jiazheng Yuan et al. “Genome-Wide Association Study of Resistance to Potato
542 Common Scab”. In: *Potato Research* (2019). ISSN: 18714528. DOI: 10.1007/
543 s11540-019-09437-w.
- 544 [34] Shengkui Zhang et al. “Genome-wide association studies of 11 agronomic
545 traits in cassava (*Manihot esculenta crantz*)”. In: *Frontiers in Plant Science*
546 9.April (2018), pp. 1–15. ISSN: 1664462X. DOI: 10.3389/fpls.2018.
547 00503.