

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 4, 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

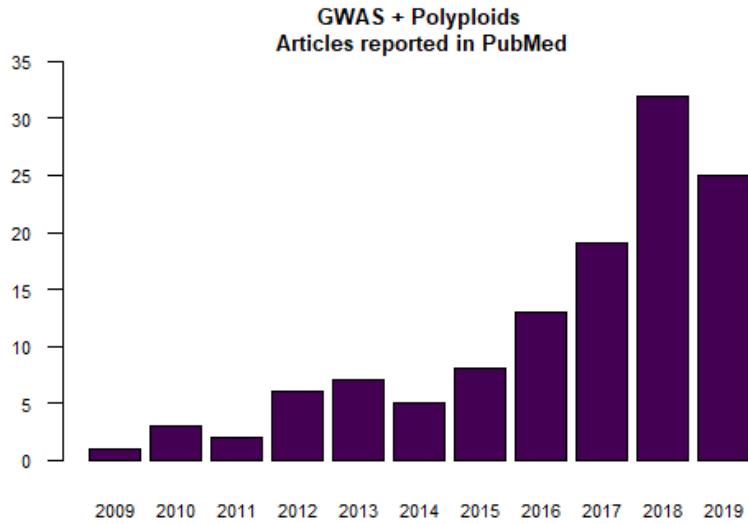


Figure 1: Timeline for articles reported for GWAS studies on polyploid species in PubMed. We present data for completed years.

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 analysis
62 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.

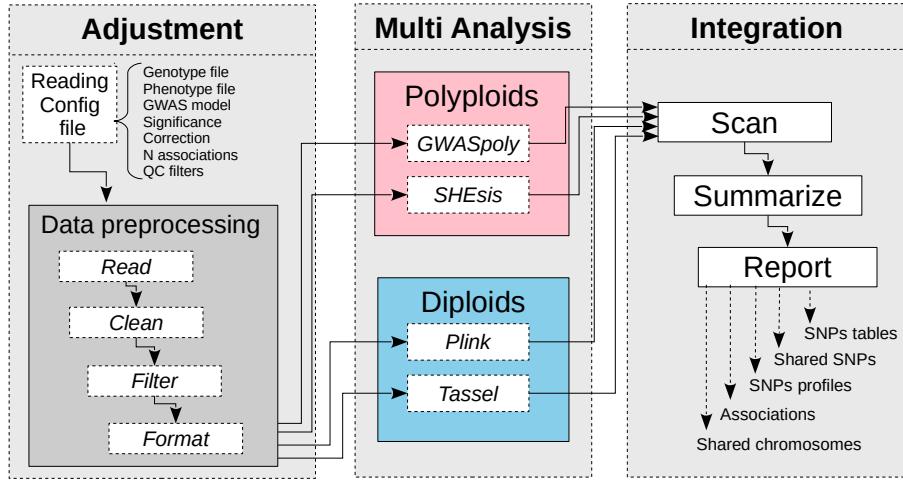


Figure 2: MultiGWAS flowchart has three consecutive steps: adjustment, multi analysis, and integration. The adjustment step manages the input data, reads the configuration file, and preprocesses the input genomic data (genotype and phenotype). The multi analysis step configures and runs the four GWAS 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, the genomic data, genotype and phenotype, are converted to
178 the 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
181 parameters file, the Full model (Q+K) and Naive (i.e., without any control) [27].
182 The Full model (Q+K) controls for both population structure and individual relatedness.
183 For population structure, MultiGWAS uses the Principal Component Analysis
184 (PCA) and takes the top five PC as covariates. For relatedness, MultiGWAS uses
185 kinship matrices that TASSEL and GWASpoly calculated separately, and for PLINK
186 and SHEsis, relatedness depends of kinship coefficients calculated with the PLINK
187 2.0 built-in algorithm [8].

188 As MultiGWAS implements two types of GWAS analysis, naive and full, each
189 tool is called in two different ways.

190 2.2.1 GWASpoly

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

199 MultiGWAS is using GWASpoly version 1.3, employing all gene action models
200 to find associations and reporting the top *N* best-ranked (the SNPs with lowest *p-*
201 values), where *N* is defined by the user in the input 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 another program designed for polyploid species that includes single locus
207 association analysis, among others. It is based on a linear regression model, and it
208 has been used in some studies of animals and humans [24, 21].

209 MultiGWAS is using the version 1.0 which does not take account for popula-
210 tion structure or relatedness, however MultiGWAS externally estimates relatedness
211 for SHEsis by excluding individuals with cryptic first-degree relatedness using the
212 algorithm 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 diploids species. It
215 was developed for humans but it is applicable to any species [22]. PLINK includes
216 a range of analysis, including univariate GWAS using two-sample tests and linear
217 regression models.

218 MultiGWAS is using two versions of PLINK: 1.9 and 2.0. Linear regression from
219 PLINK 1.9 is used to achieve both types of analysis, naive and full. For the full
220 analysis, population structure is estimated using the first five principal components
221 calculated with the PLINK 1.9 built in algorithm. But relatedness is estimated from
222 the kinship coefficients calculated with the PLINK 2.0 built in algorithm, removing
223 the close relatives or individuals with first-degree relatedness.

224 **2.2.4 TASSEL**

225 TASSEL is another common GWAS program based on the Java software. It was de-
226 veloped for maize and it has been used in several studies in plants [1, 34], but like
227 PLINK, it is applicable to any species. For association analysis, TASSEL includes the
228 general lineal model (GLM) and mixed linear model (MLM) that accounts for popu-
229 lation structure and relatedness. **And, in the same manner that GWASPoly, TASSEL**
230 **provides three gene action models to calculate the SNP effect of each genotypc class:**
231 **general, additive, and dominant, and so the significance threshold depends of each**
232 **action model.**

233 MultiGWAS is using TASSEL 5.0, **with all gene action models used to find the N**
234 **best-ranked associations and reporting the top N best-ranked associations (SNPs**
235 **with lowest p -values).** Naive GWAS is achieved by the GLM, and full GWAS is
236 achieved by the MLM with two parameters: population structure that uses the first
237 five principal components, and relatedness that uses the kinship matrix with cen-
238 tered IBS method, both calculated with the TASSEL built-in algorithms.

239 **2.3 Integration stage.**

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

243 **2.3.1 Calculation of p -values and significance thresholds**

244 GWAS packages compute p -value as a measure of association between each individ-
245 ual SNP and the trait of interest. The SNPs are considered statistically significant,
246 and so possible true associations, when their p -value drops below a predefined sig-
247 nificance threshold. But, most GWAS packages compute differently p -values with
248 the possibility to compute them too high or too low. If p -values are too high, then
249 it would lead to false negatives or SNPs with true associations with the phenotype
250 but that does not reach the significance threshold. Conversely, if p -values are too

251 low, then it would lead to false positives or SNPs with false associations with the
252 phenotype but that reaches the significance threshold.

253 To overcome these difficulties, in the case of too high *p-values*, MultiGWAS identifies
254 and reports both significant and best-ranked associations (the ones closer to
255 being statistically significant). Whereas, in the case of too low *p-values*, MultiG-
256 WAS provides two methods for adjusting *p-values* and significance threshold: the
257 false discovery rate (FDR) that adjust *p-values*, and the Bonferroni correction, that
258 adjusts the threshold.

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

268 2.3.2 Selection of significant and best-ranked associations

269 After corrections, significant associations are selected as the ones with *p-values*
270 falling below a significant threshold, which is calculated for each GWAS package.
271 But, as described above, it is equally important to know the best-ranked associa-
272 tions, closer to being statistically significant, as they may represent important asso-
273 ciations to consider for posterior analysis.

274 In the case of GWAS packages with only one gene action model (PLINK and
275 SHESIS), the best-ranked associations are selected from the top N identified by
276 the package. But, in the case of GWAS packages with several gene action mod-
277 els (GWASpoly and TASSEL), the best-ranked associations are selected as the top
278 N from the “best action model”, the one with more shared SNP associations, in
279 other words, from the action model that identifies more associations that are also
280 identified in the other models.

281 2.3.3 Integration of results

282 At this stage, MultiGWAS integrates the results to evaluate reproducible results
283 among tools (Fig 4). But, it still reports a summary for the results of each tool:

- 284 • A Quantile-Quantile (QQ) plots for the resultant *p-values* of each tool and
285 the corresponding inflation factor λ to asses the degree of the test statistic
286 inflation.
- 287 • A Manhattan plot of each tool with two lower thresholds, one for the best-
288 ranked SNPs, and another for the significant SNPs.

289 To present the replicability, we use two sets: (1) the set of all the significative SNPs
290 provided by each tool and (2) the set of all the best-ranked SNPs. For each set,

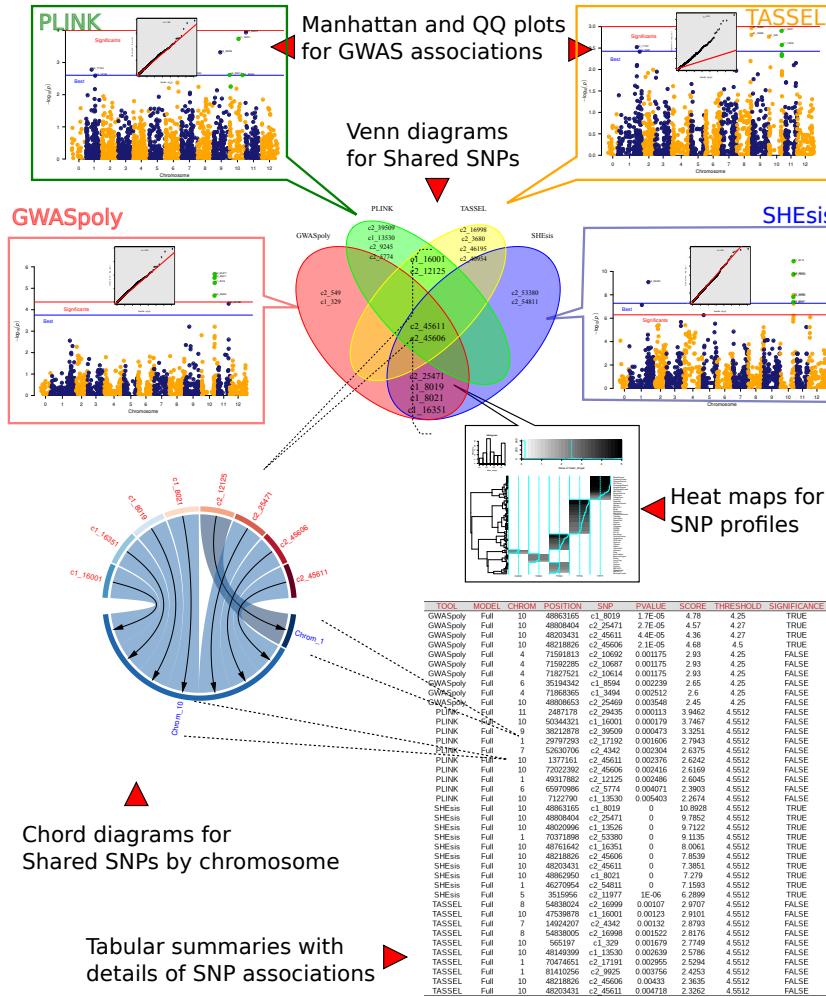


Figure 4: Reports presented by MultiGWAS. For each tool, first a QQ plot that assesses the resultant p-values. Second, a Manhattan plot for each tool with two lines, blue and red, respectively, is the lower limit for the best ranked and significative SNPs. We present two Venn diagrams, one for the significative 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, respectively. For each SNP that is in the intersection; thus, that is predicted by more than one tool we provide SNP profile. **SNPs by chromosome chord diagrams** shows how the strongest associations are mostly found on few chromosomes. And we also present tabular summaries with details of significant and best-ranked associations.

291 we present a Venn diagram that shows SNPs predicted exclusively by one tool and
292 intersections that help to identify the SNPs predicted by one, two, three, or all the
293 tools. In addition, we provide detailed tables for the two sets.

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

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

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

307 3 Results

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

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

- 318 • Manhattan and Q-Q plots to show GWAS associations resulting from each
319 GWAS package
- 320 • Venn diagrams to show SNP associations identified by both single or several
321 tools
- 322 • Heat diagrams to show the genotypic structure of the shared SNPs (SNP pro-
323 files)
- 324 • Chord diagrams to show how shared SNPs correlate with chromosomes
- 325 • Score tables to show detailed information of associations for both summary
326 results from MultiGWAS and particular results from each GWAS package

327 As an example of the functionality of the tool, here we show the outputs reported by
328 MultiGWAS in the tetraploid potato diversity panel, genotyped and phenotyped as
329 part of the USDA-NIFA Solanaceae Coordinated Agricultural Project (SolCAP) [16].

330 The complete report from MultiGWAS for the naive and full model is in the Supple-
 331 mentary information (<https://github.com/agrosavia-bioinformatic/multiGWAS>)

332 3.1 Visualization of Associations

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

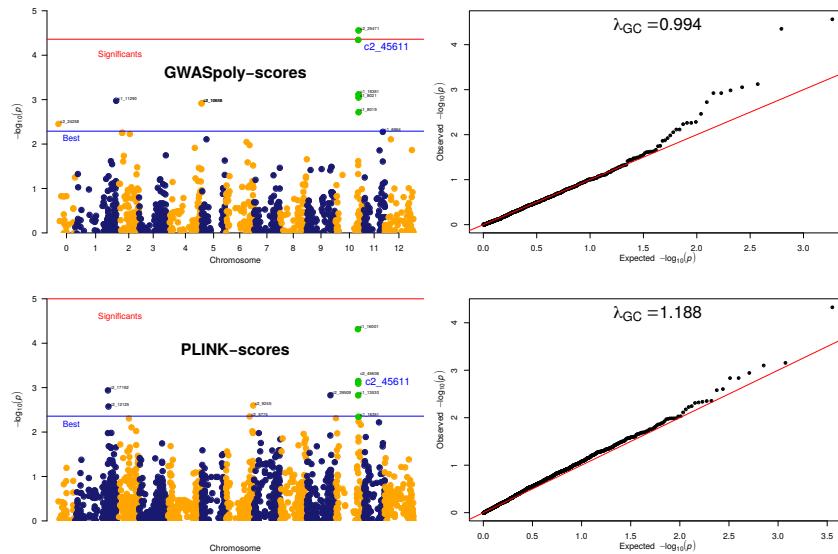


Figure 5: MultiGWAS visualization of associations. MultiGWAS adds special marks to the Manhattan and QQ plots to help identify different types of SNPs: (a) In Manhattan plots, significant SNPs are above a red line, best-ranked SNPs are above a blue line, and shared SNPs (See Figure 6.b) are colored in green (b) In QQ plots, a red diagonal line indicates the expectation, so potential associations can be observed when the number of SNPs deviating from the diagonal is small, as in the case of monogenic traits, or when this number is somewhat higher, as in the case of truly polygenic traits. However, deviations for a high number of SNPs could reflect inflated p -values owing to population structure or cryptic relatedness.

342 3.2 Shared SNPs view

343 MultiGWAS provides tabular and graphic views to report in an integrated way both
 344 the best-ranked and significant SNPs identified by the four GWAS packages (see Fig-
 345 ure 6). Both p -values and significance levels have been scaled as $-\log_{10}(p\text{-value})$
 346 to give high scores to the best statistically evaluated SNPs. First, the best-ranked

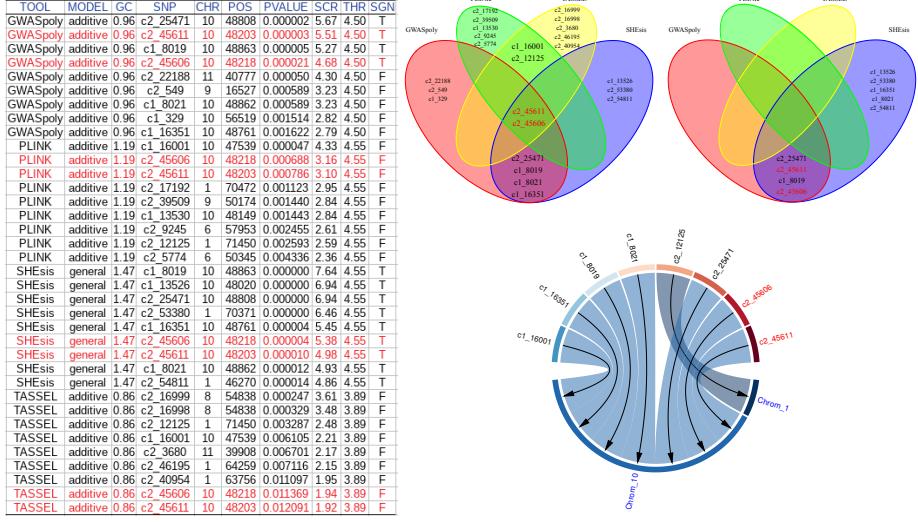


Figure 6: Shared SNPs Views. Tabular and graphical views of SNP associations identified by one or more GWAS packages (shared SNPs). SNPs identified by all packages are colored with red color in all figures **(a)** Table with details of the N=9 best-ranked SNPs from each GWAS package. Each row corresponds 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 and marked with red color. 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). **(d)** Chord diagram showing the connection of shared SNPs with chromosomes.

347 SNPs correspond to the top-scored N SNPs, wheter these SNPs were assesed signif-
 348 icant or not by the package reporting them, and with N defined by the user in the
 349 configuration file. These SNPs are shown both in a SNPs table (Figure 6.a) and in
 350 a Venn diagram (Figure 6.b). The table lists them by package and sorts by decreasing
 351 score, whereas the Venn diagram shows them emphasizing if these ones were
 352 best-ranked either in a single package or in several at once (shared). And second,
 353 the significant SNPs correspond to the ones assesed statistically significant by each
 354 package, they are shown in a Venn diagram (Figure 6.b), and they are also shown
 355 in the SNPs table, marked with significance TRUE and score greater than threshold,
 356 columns SGN, SCR, and THR, respectively in the table of the Figure6.a.

3.3 Visualization of common significative SNPs

358 MultiGWAS creates a two-dimensional representation, called SNP profile, to visu-
 359 alize each trait by individuals and genotypes as rows and columns, respectively

(Figure 7). At the left, the individuals are grouped in a dendrogram by their genotype. At the right, there is the name or ID of each individual. At the bottom, the genotypes are ordered from left to right, starting from the major to the minor allele (i.e., AAAA, AAAB, AABB, ABBB, BBBB). At the top, there is a description of the trait based on a histogram of frequency (top left) and by an assigned color for each numerical phenotype value using a grayscale (top right). Thus, each individual appears as a colored line by its phenotype value on its genotype column. For each column, there is a solid cyan line with the mean of each column and a broken cyan line that indicates how far the cell deviates from the mean.

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

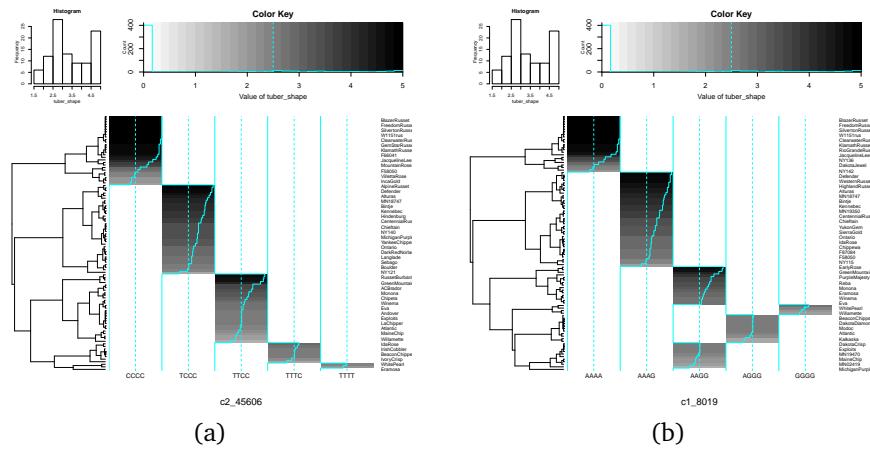


Figure 7: SNP profiles. SNP profiles for two of the best-ranked significant SNPs shown in the figure 6.b. (a) SNP c2_45606 best-ranked by the four packages (central intersection of the Venn diagram Figure 6.b) (b) SNP c1_8019 best-ranked by the two tetraploid packages (Figure 6.b), and also identified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

371 4 Availability and implementation:

The core of the MultiGWAS tool was developed in R and users can interact with the tool by either a command line interface (CLI) developed in R or a graphical user interface (GUI) developed in Java (Figure 8). Source code, examples, documentation and installation instructions are available at <https://github.com/agrosavia-bioinformatics/multiGWAS>.

377 4.1 Input parameters

MutiGWAS uses as the only input a simple configuration text file where users set the values for the main parameters that drives the GWAS process. The input parameters include: the output folder where results will be written, input genotype/phenotype

381 filenames, genome-wide significance threshold, method for multiple testing cor-
382 rection, GWAS model, number of associations to be reported, and TRUE or FALSE
383 whether to use quality control filters or not. The filters are: minor allele frequency,
384 individual missing rate, SNP missing rate, and Hardy-Weinberg threshold.

385 The configuration file can be created either using a general text editor or using
386 the GUI application. In the first case, the file must have the structure shown in the
387 Figure 8.a, where parameter names and values are separated by colon, filenames
388 are enclosed in quotation marks, and TRUE or FALSE indicates wheter filters are ap-
389 plied or not. Moreover examples for the config file <https://github.com/agrosavia-bioinformatics/MultiGWAS/tree/master/examples>

391 In the second case, the user creates the config file in a simple and straightforward
392 way using the input parameter view from the GUI application (Figure 8.b) and
393 clicking the “Save” button.

394 **4.2 Using the command line interface**

395 The execution of the tool in command line is simple, it only needs to open a linux
396 console, change to the folder where the configuration file was created, and type the
397 name of the executable tool followed by the filename of the configuration file, like
398 this:

399 `multiGWAS full.config`

400 Then, the tool starts the execution, showing information of the process in the
401 console window, and when it finishes the results are saved to a new subfolder called
402 “*outgwas/reports*”. Results include a full html report containing the different views
403 described in the results section, along with the original graphics and summary tables
404 created by MultiGWAS and used to create the html report. Additionally, results
405 include the preprocessed tables of the main outputs generated by the four GWAS
406 packages used by MultiGWAS.

407 **4.3 Using the graphical user interface**

408 The MultiGWAS GUI application can be executed either by running from a Linux
409 console the *jmultiGWAS* command or by clicking on the Java application file *JMulti-*
410 *GWAS.jar* located in the “*multiGWAS/sources*” subfolder. After it opens, it shows a
411 main frame with four tabs at the top (Figure 8b): “Inputs”, “Outputs”, “Results”, and
412 “Files”. The “Inputs” tab shows the form to create the configuration file and run the
413 application. The “Outputs” tab shows the messages from the running process after
414 it starts the execution. The “Results” tab shows the full html report described above.
415 And the “Files” tab shows an embedded file browser pointing to the subfolder that
416 contains the original files used in the html report and described above.

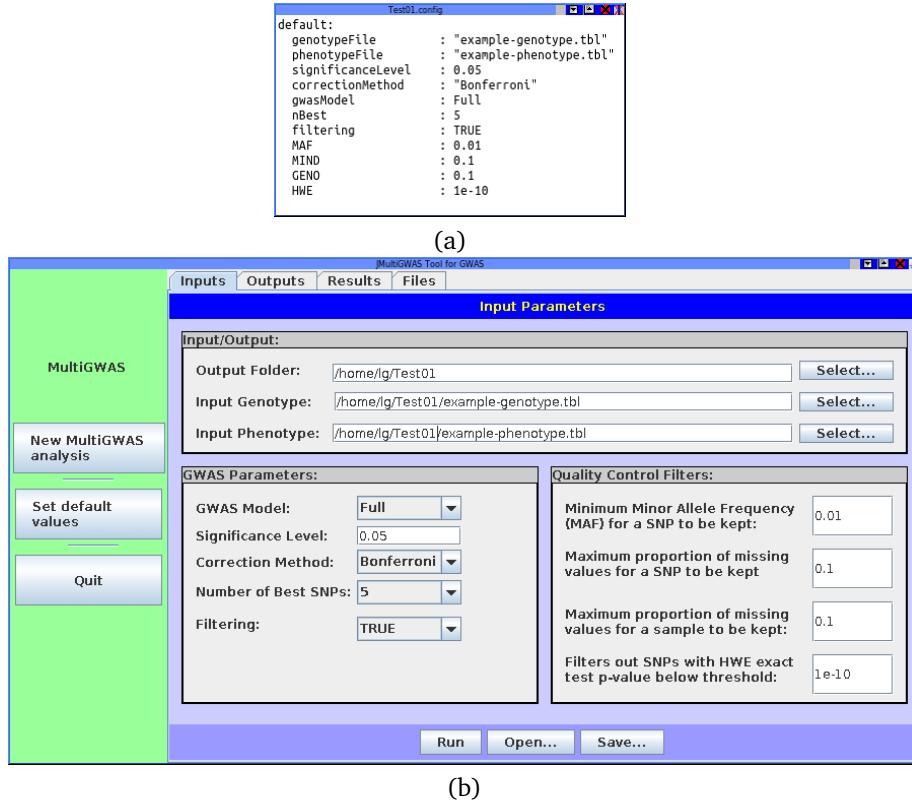


Figure 8: MultiGWAS inputs and interaction. MultiGWAS uses as input a simple configuration text file and can be executed using either a command line interface script in R (CLI) or a graphic user interface application in Java (GUI). (a) An example of a configuration text file named “*Test01.config*” including the parameters that drive the GWAS process. It can be created using a general text editor or using the GUI application (see below) (b) 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.

417 5 Discussion

418 XXXXXXXXXXXXXXXXXXXXXXXXX

419 References

- 420 [1] María F. Álvarez et al. “Identification of Novel Associations of Candidate
 421 Genes with Resistance to Late Blight in *Solanum tuberosum* Group Phureja”.
 422 In: *Frontiers in Plant Science* 8 (2017), p. 1040. ISSN: 1664-462X. DOI: 10 .

- 423 3389/fpls.2017.01040. URL: <http://journal.frontiersin.org/article/10.3389/fpls.2017.01040/full>.
- 424
- 425 [2] Ferdouse Begum et al. "Comprehensive literature review and statistical con-
426 siderations for GWAS meta-analysis". In: *Nucleic acids research* 40.9 (2012),
427 pp. 3777–3784.
- 428 [3] Jhon Berdugo-Cely et al. "Genetic diversity and association mapping in the
429 colombian central collection of solanum tuberosum L. Andigenum group us-
430 ing SNPs markers". In: *PLoS ONE* 12.3 (2017). ISSN: 19326203. DOI: 10.
431 1371/journal.pone.0173039.
- 432 [4] Peter M. Bourke et al. "Tools for Genetic Studies in Experimental Popula-
433 tions of Polyploids". In: *Frontiers in Plant Science* 9 (2018), p. 513. ISSN:
434 1664-462X. DOI: 10.3389/fpls.2018.00513. URL: <http://journal.frontiersin.org/article/10.3389/fpls.2018.00513/full>.
- 435
- 436 [5] Peter J Bradbury et al. "TASSEL: software for association mapping of complex
437 traits in diverse samples". In: *Bioinformatics* 23.19 (2007), pp. 2633–2635.
438 ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btm308. URL: <https://doi.org/10.1093/bioinformatics/btm308>.
- 439
- 440 [6] Rita M Cantor, Kenneth Lange, and Janet S Sinsheimer. "Prioritizing GWAS
441 results: a review of statistical methods and recommendations for their ap-
442 plication". In: *The American Journal of Human Genetics* 86.1 (2010), pp. 6–
443 22.
- 444 [7] Jun Cao et al. "Whole-genome sequencing of multiple Arabidopsis thaliana
445 populations". In: *Nature genetics* 43.10 (2011), p. 956.
- 446 [8] Christopher C. Chang et al. "Second-generation PLINK: Rising to the chal-
447 lenge of larger and richer datasets". In: *GigaScience* 4.1 (2015), pp. 1–16.
448 ISSN: 2047217X. DOI: 10.1186/s13742-015-0047-8. arXiv: 1410.
449 4803.
- 450 [9] Rishika De, William S Bush, and Jason H Moore. "Bioinformatics Challenges
451 in Genome-Wide Association Studies (GWAS)". In: *Clinical Bioinformatics*.
452 Ed. by Ronald Trent. New York, NY: Springer New York, 2014, pp. 63–81.
453 ISBN: 978-1-4939-0847-9. DOI: 10.1007/978-1-4939-0847-9_5. URL:
454 https://doi.org/10.1007/978-1-4939-0847-9%7B%5C_%7D5.
- 455 [10] Robert Eklblom and Juan Galindo. "Applications of next generation sequenc-
456 ing in molecular ecology of non-model organisms". In: *Heredity* 107.1 (2011),
457 pp. 1–15.
- 458 [11] Hans Ellegren. "Genome sequencing and population genomics in non-model
459 organisms". In: *Trends in ecology & evolution* 29.1 (2014), pp. 51–63.
- 460 [12] Luís Felipe V. Ferrão et al. "Insights Into the Genetic Basis of Blueberry Fruit-
461 Related Traits Using Diploid and Polyploid Models in a GWAS Context". In:
462 *Frontiers in Ecology and Evolution* 6 (2018), p. 107. ISSN: 2296-701X. DOI:
463 10.3389/fevo.2018.00107. URL: <https://www.frontiersin.org/articles/10.3389/fevo.2018.00107/full>.
- 464

- 465 [13] Dominik G Grimm et al. "easyGWAS: A Cloud-Based Platform for Comparing
466 the Results of Genome-Wide Association Studies". In: *The Plant Cell* 29.1
467 (2017), pp. 5–19. ISSN: 1040-4651. DOI: 10.1105/tpc.16.00551. URL:
468 <http://www.plantcell.org/content/29/1/5>.
- 469 [14] Anja C Gumpinger et al. *Methods and Tools in Genome-wide Association Stud-*
470 *ies*. Vol. 1819. 2018. ISBN: 9781493986187.
- 471 [15] Bin Han and Xuehui Huang. "Sequencing-based genome-wide association
472 study in rice". In: *Current opinion in plant biology* 16.2 (2013), pp. 133–138.
- 473 [16] Candice N. Hirsch et al. "Retrospective view of North American potato (*Solanum*
474 *tuberosum* L.) breeding in the 20th and 21st centuries". In: *G3: Genes, Genomes,*
475 *Genetics* 3.6 (2013), pp. 1003–1013. ISSN: 21601836. DOI: 10.1534/g3.
476 113.005595.
- 477 [17] "How to interpret a genome-wide association study". In: *JAMA - Journal*
478 *of the American Medical Association* 299.11 (2008), pp. 1335–1344. ISSN:
479 00987484. DOI: 10.1001/jama.299.11.1335.
- 480 [18] Avjinder S Kaler and Larry C Purcell. "Estimation of a significance threshold
481 for genome-wide association studies". In: *BMC Genomics* 20.1 (2019), p. 618.
482 ISSN: 1471-2164. DOI: 10.1186/s12864-019-5992-7. URL: <https://doi.org/10.1186/s12864-019-5992-7>.
- 484 [19] Arthur Korte and Ashley Farlow. "The advantages and limitations of trait
485 analysis with GWAS: a review". In: *Plant methods* 9.1 (2013), p. 29.
- 486 [20] Gordan Lauc et al. "Genomics meets glycomics—the first GWAS study of hu-
487 man N-glycome identifies HNF1 α as a master regulator of plasma protein
488 fucosylation". In: *PLoS genetics* 6.12 (2010).
- 489 [21] Jie Meng et al. "Genome-wide association analysis of nutrient traits in the
490 oyster *Crassostrea gigas*: Genetic effect and interaction network". In: *BMC*
491 *Genomics* 20.1 (2019), pp. 1–14. ISSN: 14712164. DOI: 10.1186/s12864-
492 019-5971-z.
- 493 [22] Robert A. Power, Julian Parkhill, and Tulio De Oliveira. "Microbial genome-
494 wide association studies: lessons from human GWAS". In: *Nature Reviews*
495 *Genetics* 18.1 (2016), pp. 41–50. ISSN: 14710064. DOI: 10.1038/nrg.2016.132.
- 497 [23] Shaun Purcell et al. "PLINK: A tool set for whole-genome association and
498 population-based linkage analyses". In: *American Journal of Human Genetics*
499 81.3 (2007), pp. 559–575. ISSN: 00029297. DOI: 10.1086/519795.
- 500 [24] Hui Ping Qiao et al. "Genetic variants identified by GWAS was associated
501 with colorectal cancer in the Han Chinese population". In: *Journal of Cancer*
502 *Research and Therapeutics* 11.2 (2015), pp. 468–470. ISSN: 19984138. DOI:
503 10.4103/0973-1482.150346.

- 504 [25] Umesh R. Rosyara et al. “Software for Genome-Wide Association Studies
505 in Autopolyploids and Its Application to Potato”. In: *The Plant Genome* 9.2
506 (2016), pp. 1–10. ISSN: 1940-3372. DOI: 10.3835/plantgenome2015.08.0073.
507 URL: <https://dl.sciencesocieties.org/publications/tpg/abstracts/9/2/plantgenome2015.08.0073>.
- 509 [26] Anna W Santure and Dany Garant. “Wild GWAS—association mapping in
510 natural populations”. In: *Molecular ecology resources* 18.4 (2018), pp. 729–
511 738.
- 512 [27] Sanjeev Kumar Sharma et al. “Linkage disequilibrium and evaluation of genome-
513 wide association mapping models in tetraploid potato”. In: *G3: Genes, Genomes,
514 Genetics* 8.10 (2018), pp. 3185–3202. ISSN: 21601836. DOI: 10.1534/g3.118.200377.
- 516 [28] Jiawei Shen et al. “SHEsisPlus, a toolset for genetic studies on polyploid
517 species”. In: *Scientific Reports* 6 (2016), pp. 1–10. ISSN: 20452322. DOI: 10.
518 1038/srep24095. URL: <http://dx.doi.org/10.1038/srep24095>.
- 519 [29] John R Thompson, John Attia, and Cosetta Minelli. “The meta-analysis of
520 genome-wide association studies”. In: *Briefings in Bioinformatics* 12.3 (2011),
521 pp. 259–269. ISSN: 1467-5463. DOI: 10.1093/bib/bbr020. URL: <https://doi.org/10.1093/bib/bbr020>.
- 523 [30] Feng Tian et al. “Genome-wide association study of leaf architecture in the
524 maize nested association mapping population”. In: *Nature genetics* 43.2 (2011),
525 pp. 159–162.
- 526 [31] Yan Y. Yan et al. “Effects of input data quantity on genome-wide association
527 studies (GWAS)”. In: *International Journal of Data Mining and Bioinformatics*
528 22.1 (2019), pp. 19–43. ISSN: 17485681. DOI: 10.1504/IJDMB.2019.
529 099286.
- 530 [32] J Yu et al. “A unified mixed-model method for association mapping that ac-
531 counts for multiple levels of relatedness.” In: *Nature genetics* 38.2 (2006),
532 pp. 203–208.
- 533 [33] Jiazheng Yuan et al. “Genome-Wide Association Study of Resistance to Potato
534 Common Scab”. In: *Potato Research* (2019). ISSN: 18714528. DOI: 10.1007/
535 s11540-019-09437-w.
- 536 [34] Shengkui Zhang et al. “Genome-wide association studies of 11 agronomic
537 traits in cassava (*Manihot esculenta crantz*)”. In: *Frontiers in Plant Science*
538 9.April (2018), pp. 1–15. ISSN: 1664462X. DOI: 10.3389/fpls.2018.
539 00503.