

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

4      L. Garreta<sup>1</sup>, I. Cerón-Souza<sup>1</sup>, M.R. Palacio<sup>2</sup>, and P.H.  
5      Reyes-Herrera<sup>1</sup>

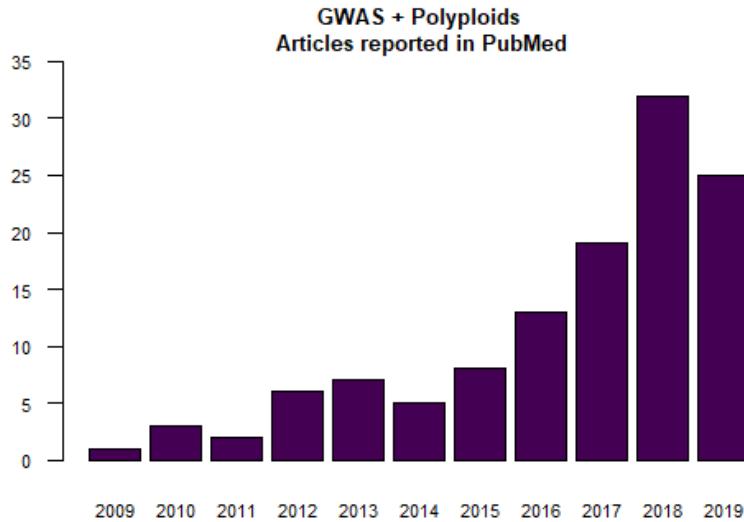
6      <sup>1</sup>Corporación Colombiana de Investigación Agropecuaria  
7      (AGROSAVIA), CI Tibaitatá, Kilómetro 14, Vía a Mosquera,  
8      250047, Colombia  
9      <sup>2</sup>Corporación Colombiana de Investigación Agropecuaria  
10     (AGROSAVIA), CI El Mira, Kilómetro 38, Vía Tumaco Pasto,  
11     Colombia

12     June 12, 2020

13     **Abstract**

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

32     MultiGWAS generates five summary results integrating the four tools.  
33     (1) Score tables with detailed information on the associations for each tool.  
34     (2) Venn diagrams of shared SNPs among the four tools. (3) Heatmaps



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

36 of significative SNP profiles among the four tools. (4) Manhattan and QQ  
 37 plots for the association found by each tool. And (5) Chord diagrams for  
 38 the chromosomes vs. SNP by each tool. **Contact:** phreyes@agrosavia.co

39  
 40 **Keywords:** GWAS, tetraploids, SNPs,XXX

## 41 1 Introduction

42 The Genome-Wide Association Study (GWAS) is used to identify which variants  
 43 through the whole genome of a large number of individuals are associated with  
 44 a specific trait [6, 2]. This methodology started with humans and several model  
 45 plants, such as rice, maize, and *Arabidopsis* [20, 30, 7, 19, 15]. Because of  
 46 the advances in the next-gen sequencing technology and the decline of the se-  
 47 quencing cost in recent years, there is an increase in the availability of genome  
 48 sequences of different organisms at a faster rate [10, 11]. Thus, the GWAS is  
 49 becoming the standard tool to understand the genetic bases of either ecolog-  
 50 ical or economic phenotypic variation for both model and non-model organisms.  
 51 This increment in GWAS includes complex species such as polyploids (Fig 1)  
 52 [10, 26].

53 The GWAS for polyploid species has three related challenges. First, as all  
 54 GWAS, we should replicate the study as a reliable method to validate the results

55 and recognize real associations. This replication involves finding the same associations either in several replicates from the study population using the same software or testing different GWAS tools among the same study population.  
56 This approach involved the use of different parameters, models, or conditions,  
57 to test how consistent the results are [9, 17]. However, the performance of  
58 different GWAS software could affect the results. For example, the threshold  
59 *pvalue* for SNP significance change through four GWAS software (i.e., PLINK,  
60 TASSEL, GAPIT, and FaST-LMM) when sample size varies [31]. It means that  
61 well-ranked SNPs from one package can be ranked differently in another.  
62

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

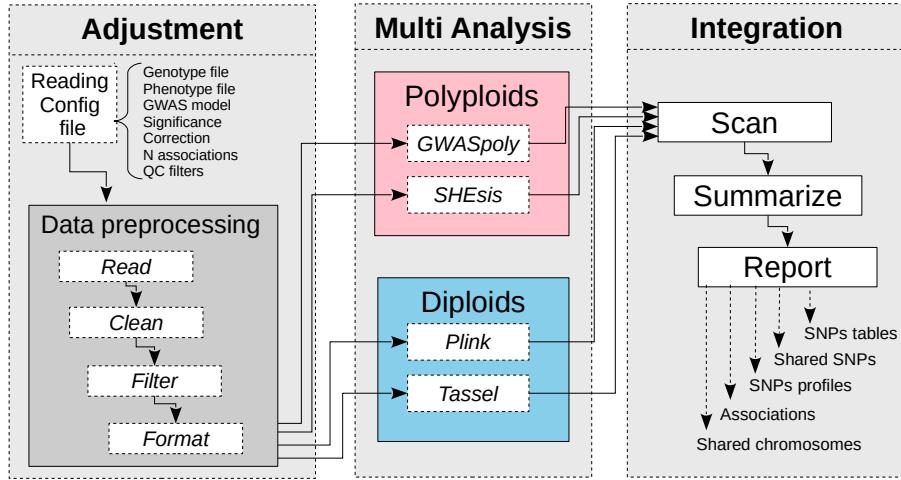
64 Third, there are very few tools focused on the integration of several GWAS software, to make comparisons under different parameters and conditions across them. As far as we know, there is only two software with this service in mind, such as iPAT and easyGWAS.

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

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

## 90 2 Method

91 The MultiGWAS tool has three main consecutive steps: the adjustment, the multi analysis, and the integration (Fig. 2). In the adjustment step, MultiGWAS processes the configuration file. Then it cleans and filters the genotype and phenotype, and MultiGWAS "diploidize" the genomic data. Next, during the multi analysis, each GWAS tool runs in parallel. Subsequently, in the integration step, the MultiGWAS tool scans the output files from the four packages (i.e., GWASPoly, SHEsis, PLink, and TASSEL). Finally, it generates a summary of all



**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.

98 results that contains score tables, Venn diagrams, SNP profiles, and Manhattan  
99 plots.

## 100 2.1 Adjustment stage

101 MultiGWAS takes as input a configuration file where the user specifies the ge-  
102 nomics data along with the parameters that will be used by the four tools. Once  
103 the configuration file is processed, MultiGWAS preprocess the data that is clean-  
104 ing, filtering, and checking data quality. The output of this stage corresponds to  
105 the inputs for the four programs at the Multi Analysis stage.

### 106 2.1.1 Reading configuration file

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

108 **Input genotype and phenotype files:** Currently, MultiGWAS uses two input  
109 files, one for genotype and the other for the phenotype. Both data correspond  
110 to data matrices with column and row names (Figure 3). The genotype file uses  
111 SNP markers in rows and samples in columns (Figure 3a). The phenotype file  
112 uses samples in rows and traits in columns (Figure 3b) with the first column  
113 corresponding to the sample name and the second column to trait value.

Marker,Chrom,Pos,Indiv01,Indiv02,Indiv03,.... c2_41437,0,805179,AAAG,AAGG,AAGG,... c2_24258,0,1252430,AAGG,AGGG,GGGG,... c2_21332,0,3499519,TTCC,TTCC,TTCC,... ...	a	Individual,Traitname Indiv01, 3.59 Indiv02, 4.07 Indiv03, 1.05 ...	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.

114 **GWAS model:** MultiGWAS is designed to work with quantitative phenotypes  
115 and can run GWAS analysis using two types of statistical models that we have  
116 called *full* and *naive* models. The *full model* is known in the literature as the  
117 Q+K model [32] and includes a control for structure (Q) and relatedness be-  
118 tween samples (K). In contrast, the *naive model* does not include any correction.  
119 Both models are linear regression approaches, and the four GWAS packages  
120 used by MultiGWAS implemented variations of them. The *naive* is modeled  
121 with Generalized Linear Models (GLMs, Phenotype + Genotype), and the *full*  
122 is modeled with Mixed Linear Models (MLMs, Phenotype + Genotype + Struc-  
123 ture + Kinship). The default model used by MultiGWAS is the *full model* (Q+K)  
124 [32], following this equation:

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

125 The vector  $y$  represents the observed phenotypes depends on the following  
126 factors: the fixed effect vector  $\beta$ , the SNP effects vector  $\alpha$ , the population ef-  
127 fect vector  $\nu$ , the polygene background effect vector  $\mu$ , and, the residual effect  
128 vector  $e$ . The  $Q$ , modeled as a fixed effect, refers to the incidence matrix for  
129 subpopulation covariates relating  $y$  to  $\nu$ . Moreover,  $X$ ,  $S$ , and  $Z$  are incidence  
130 matrices of ones and zeros relating  $y$  to  $\beta$ ,  $\alpha$ , and  $\mu$ , respectively.

131 **Genome-wide significance:** GWAS searches SNPs associated with a pheno-  
132 type trait in a statistically significant manner. A threshold or significance level  
133  $\alpha$  is specified and compared with the *p-value* derived for each association score.  
134 Standard significance levels are 0.01 or 0.05 [14, 25], and MultiGWAS uses an  
135  $\alpha$  of 0.05 for the four GWAS packages. However, the adjustment of the thresh-  
136 old is according to each package. For example, GWASpoly and TASSEL calculate  
137 the SNP effect for each genotypic class using different gene action models (see  
138 “Multi analysis stage”). Therefore, the number of tested markers may be differ-  
139 ent in each model (see below) that results in different *p-value* thresholds.

140 **Multiple testing correction:** Due to the massive number of statistical tests  
141 performed by GWAS, it is necessary to perform a correction method for mul-  
142 tiple hypothesis testing and adjusting the *p-value* threshold accordingly. Two

143 standard methods for multiple hypothesis testing are the false discovery rate  
144 (FDR) and the Bonferroni correction. The latter is the default method used  
145 by MultiGWAS because it is one of the most rigorous. MultiGWAS adjust the  
146 threshold below which a *p-value* is considered significant, that is  $\alpha/m$ , where  $\alpha$   
147 is the significance level and  $m$  is the number of tested markers from the geno-  
148 type matrix.

149 **Number of reported associations:** Criticism has arisen in considering only  
150 statistically significant associations as the only possible correct associations [29,  
151 18]. Many low *p-value* associations are closer to being significant, are discarded  
152 due to the stringent significance levels, and, consequently, increase the number  
153 of false negatives. To help to analyze both significant and non-significant asso-  
154 ciations, MultiGWAS provides the option to specify the number of best-ranked  
155 associations (lower *p-values*), adding the corresponding *p-value* to each associa-  
156 tion found. In this way, it is possible to enlarge the number of results, and  
157 we can observe replicability in the results for different programs. Nevertheless,  
158 MultiGWAS always presents each associated SNP with its corresponding *p-value*.

159 **Quality control filters:** A control step is necessary to check the input data for  
160 genotype or phenotype errors or poor quality that can lead to spurious GWAS  
161 results. MultiGWAS provides the option to select and define thresholds for the  
162 following filters that control the data quality: Minor Allele Frequency (MAF),  
163 individual missing rate (MIND), SNP missing rate (GENO), and HardyWeinberg  
164 threshold (HWE):

- 165 • **MAF of  $x$ :** filters out SNPs with minor allele frequency below  $x$  (default  
166 0.01);
- 167 • **MIND of  $x$ :** filters out all individuals with missing genotypes exceeding  
168  $x^*100\%$  (default 0.1);
- 169 • **GENO of  $x$ :** filters out SNPs with missing values exceeding  $x^*100\%$  (de-  
170 fault 0.1);
- 171 • **HWE of  $x$ :** filters out SNPs which have Hardy-Weinberg equilibrium exact  
172 test *p-value* below the  $x$  threshold.

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

### 175 2.1.2 Data preprocessing

176 Once the configuration file is processed, the genomic data is read and cleaned  
177 by selecting individuals present in both genotype and phenotype. Then, based  
178 on previous selected quality-control filters and their thresholds, MultiGWAS re-  
179 move individuals and SNPs with poor quality.

180 During this step, the format "ACGT" suitable for the polyploid software GWAS-  
181 poly and SHEsis, is "diploidized" for PLINK and TASSEL. The homozygous tetra-  
182 ploid genotypes are converted to diploid thus: AAAA→AA, CCCC→CC, GGGG→GG,  
183 TTTT→TT. Moreover, for tetraploid heterozygous genotypes, the conversion de-  
184 pends on the reference and alternate alleles calculated for each position (e.g.,  
185 AAAT→AT, ... , CCCG→CG).

186 After this process, MultiGWAS convert the genotype and phenotype data to  
187 the specific formats required for each of the four GWAS packages.

## 188 **2.2 Multi analysis stage**

189 MultiGWAS runs in parallel using two types of statistical models specified in  
190 the parameters file, the Full model (Q+K) and Naive (i.e., without any control)  
191 where Q refers to population structure and K refers to relatedness, calculated by  
192 kinship coefficients across individuals [27]. The Full model (Q+K) controls for  
193 both population structure and individual relatedness. For population structure,  
194 MultiGWAS uses the Principal Component Analysis (PCA) and takes the top five  
195 PC as covariates. For relatedness, MultiGWAS uses kinship matrices that TASSEL  
196 and GWASpoly calculated separately, and for PLINK and SHEsis, relatedness  
197 depends on kinship coefficients calculated with the PLINK 2.0 built-in algorithm  
198 [8].

### 199 **2.2.1 GWASpoly**

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

209 MultiGWAS is using GWASpoly version 1.3 with all gene action models avail-  
210 able to find associations. The MultiGWAS reports the top *N* best-ranked (the  
211 SNPs with lowest *p-values*) that the user specified in the *N* input configuration  
212 file. The *full* model used by GWASpoly includes the population structure and  
213 relatedness, which are estimated using the first five principal components and  
214 the kinship matrix, respectively, both calculated with the GWASpoly built-in al-  
215 gorithms.

### 216 **2.2.2 SHEsis**

217 SHEsis is a program based on a linear regression model that includes single-  
218 locus association analysis, among others. The software design includes poly-

219 ploid species. However, their use is mainly in diploids animals and humans [24,  
220 21].

221 MultiGWAS is using version 1.0, which does not take account for population  
222 structure or relatedness. Despite, MultiGWAS externally estimates relatedness  
223 for SHEsis by excluding individuals with cryptic first-degree relatedness using  
224 the algorithm implemented in PLINK 2.0 (see below).

### 225 2.2.3 PLINK

226 PLINK is one of the most extensively used programs for GWAS in humans and  
227 any diploid species [22]. PLINK includes a range of analyses, including univari-  
228 ate GWAS using two-sample tests and linear regression models.

229 MultiGWAS is using two versions of PLINK: 1.9 and 2.0. Linear regression  
230 from PLINK 1.9 performs both naive and full model. For the full model, the  
231 software calculates the population structure using the first five principal compo-  
232 nents calculated with a built-in algorithm integrated into version 1.9. Moreover,  
233 version 2.0 calculates the kinship coefficients across individuals using a built-in  
234 algorithm that removes the close individuals with first-degree relatedness.

### 235 2.2.4 TASSEL

236 TASSEL is another standard GWAS program based on the Java software devel-  
237 oped initially for maize but currently used in several species [1, 34]. For the as-  
238 sociation analysis, TASSEL includes the general linear model (GLM) and mixed  
239 linear model (MLM) that accounts for population structure and relatedness.  
240 Moreover, as GWASPoly, TASSEL provides three-gene action models to calculate  
241 the SNP effect of each genotypic class: general, additive, and dominant, and so  
242 the significance threshold depends on each action model.

243 MultiGWAS is using TASSEL 5.0, with all gene action models used to find  
244 the  $N$  best-ranked associations and reporting the top  $N$  best-ranked associations  
245 (SNPs with lowest  $p$ -values). Naive GWAS uses the GLM, and full GWAS uses the  
246 MLM with two parameters: population structure that uses the first five principal  
247 components, and relatedness that uses the kinship matrix with centered IBS  
248 method, both calculated with the TASSEL built-in algorithms.

## 249 2.3 Integration stage.

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

### 253 2.3.1 Calculation of $p$ -values and significance thresholds

254 GWAS packages compute  $p$ -value as a measure of association between each SNP  
255 and the trait of interest. The real associations are those their  $p$ -value drops  
256 below a predefined significance threshold. However, the four GWAS packages

257 compute differently *p*-values with the consequence to compute them too high  
258 or too low. If *p*-values is too high, it would lead to false negatives or SNPs with  
259 real associations with the phenotype, but that does not reach the significance  
260 threshold. Conversely, if *p*-values are too low, then it would lead to false pos-  
261 itives or SNPs with false associations with the phenotype, but that reaches the  
262 significance threshold.

263 To overcome these difficulties, in the case of too high *p*-values, MultiGWAS  
264 identifies and reports both significant and best-ranked associations (the ones  
265 closest to being statistically significant). Whereas, in the case of too low *p*-  
266 values, MultiGWAS provides two methods for adjusting *p*-values and significance  
267 threshold: the false discovery rate (FDR) that adjust *p*-values, and the Bonfer-  
268 roni correction, that adjusts the threshold.

269 By default, MultiGWAS uses the Bonferroni correction that uses the signifi-  
270 cance level  $\alpha/m$  (defined by the user in the configuration file), and  $m$  (the num-  
271 ber of tested markers) to adjust the significance threshold in the GWAS study.  
272 However, the significance threshold can be different for each GWAS package as  
273 some of them use several action models to calculate the SNP effect of each geno-  
274 typic class. For both PLINK and SHEsis packages, which use only one model,  $m$   
275 is equal to the total number of SNPs. However, for both GWASpoly and TAS-  
276 SEL packages, which use eight and three gene action models, respectively,  $m$   
277 is equal to the number of tests performed in each model, which is different  
278 between models.

### 279 2.3.2 Selection of significant and best-ranked associations

280 MultiGWAS selects two groups of associations from the results of each GWAS  
281 package: statistically significant and best-ranked. The latter equally important  
282 to the former as they are associations with lowest *p*-values not reaching the sig-  
283 nificance threshold but representing interesting associations for further analysis  
284 (possible false negatives).

285 The significant associations are selected from the ones with *p*-values falling  
286 below a significant threshold, calculated for each GWAS package; and the best-  
287 ranked associations as the  $N$  ones closer to being statistically significant, with  $N$   
288 defined by the user in the configuration file.

289 The selection of these groups takes into account whether the GWAS package  
290 uses only one gene action model, as PLINK and SHEsis do, or uses several ones,  
291 as GWASpoly and TASSEL do. In the first case, there is only one resulting set  
292 of associations and the selection is straightforward, as described above. How-  
293 ever, in the second case, there are several resulting sets of associations (one for  
294 each model), and MultiGWAS takes the significant and best-ranked associations  
295 from the best gene action model. Action models are scored taking into account  
296 the number of shared SNPs between models and the inflation factor, with the  
297 following equation:

$$score(M_i) = \frac{\sum_{all\ models} SharedSNPs(M_i)}{\sum_{all\ models} MaxSharedSNPs} + 1 - |1 - \lambda(M_i)|$$

298 where  $M_i$  is the gene action model  $i$ , with  $i$  from  $1..n$ , for a GWAS package  
 299 with  $n$  gene action models;  $SharedSNPs(M_i)$  is the number of SNPs identified  
 300 at once in both the model  $M_i$  and the other models (shared SNPs);  $MaxSharedSNPs$   
 301 is the max number of shared SNPs by each model; and  $\lambda(M_i)$  is the inflation fac-  
 302 tor for the model  $M_i$ .

303 The score is high when a model  $M_i$  both identifies a high number of shared  
 304 SNPs and has an inflation factor  $\lambda$  close to 1. Conversely, the score is low when  
 305 the model  $M_i$  both identifies a small number of shared SNPs and has an inflation  
 306 factor  $\lambda$  either low (close to 0) or high ( $\lambda > 2$ ). In any other case, the score is  
 307 balanced between the number of shared SNPs and the inflation factor.

### 308 2.3.3 Integration of results

309 At this stage, MultiGWAS integrates the results to evaluate reproducible results  
 310 among tools (Fig 4). However, it still reports a summary of the results of each  
 311 tool:

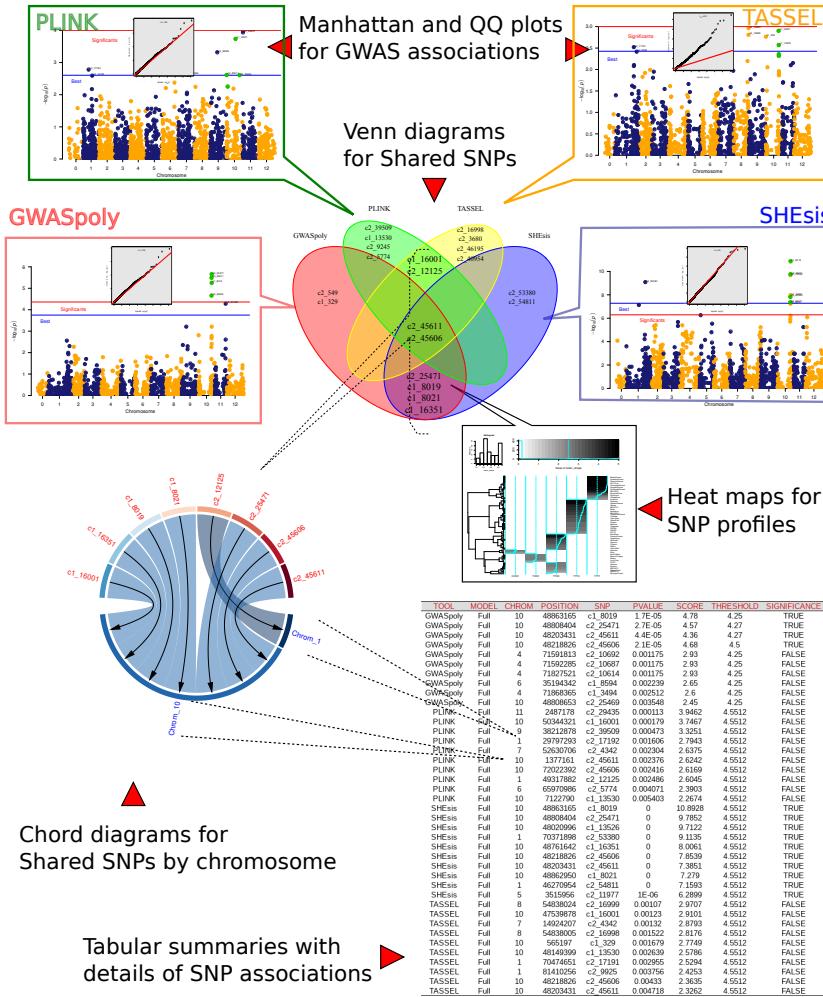
- 312 • A Quantile-Quantile (QQ) plots for the resultant  $p$ -values of each tool and  
 313 the corresponding inflation factor  $\lambda$  to assess the degree of the test statistic  
 314 inflation.
- 315 • A Manhattan plot of each tool with two lower thresholds, one for the best-  
 316 ranked SNPs, and another for the significant SNPs.

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

322 For each SNP identified more than once, we provide what we call the SNP  
 323 profile. That is a heat diagram for a specific SNP, where each column is a geno-  
 324 type state AAAA, AAAB, AABB, ABBB, and BBBB. Moreover, each row corre-  
 325 sponds to a sample. Samples with close genotypes form together clusters. Thus  
 326 to generate the clusters, we do not use the phenotype information. However,  
 327 we present the phenotype information in the figure as the color. This figure  
 328 visually provides information regarding genotype and phenotype information  
 329 simultaneously for the whole population. We present colors as tones between  
 330 white and black for color blind people.

331 MultiGWAS generates a report, one document with the content previously  
 332 described. Besides, there is a folder with the individual figures just in case the  
 333 user needs one. In the supplementary information, we include a report and a  
 334 description of the report content (supplementary information XXX)

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



**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 GWApoly, 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.

**336    3   Results**

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

**343** To alleviate these difficulties, MultiGWAS produces five types of outputs us-  
**344** ing different graphics and tabular views, these outputs are intended to help  
**345** users to compare, select, and interpret the set of possible SNPs associated with  
**346** a trait of interest. The outputs include:

- 347**     • Manhattan and Q-Q plots to show GWAS associations.
- 348**     • Venn diagrams to show associations identified by single or several tools.
- 349**     • Heat diagrams to show the genotypic structure of shared SNPs.
- 350**     • Chord diagrams to show shared SNPs by chromosomes.
- 351**     • Score tables to show detailed information of associations for both sum-  
**352**        mary results from MultiGWAS and particular results from each GWAS  
**353**        package

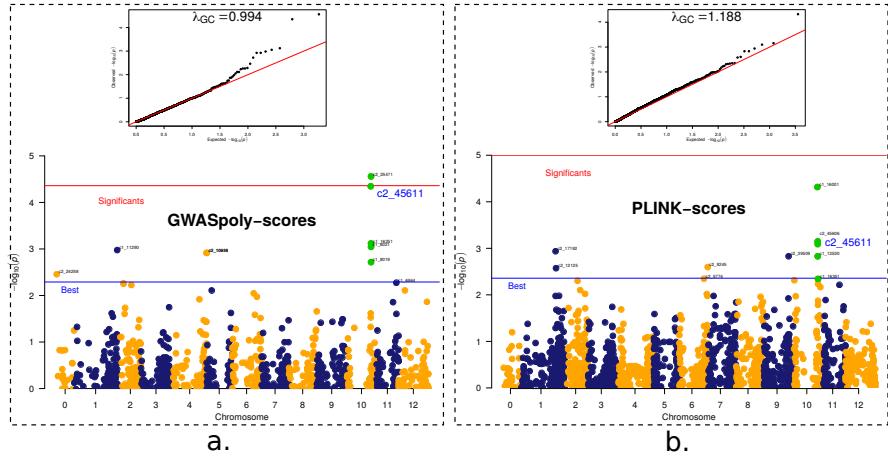
**354** As an example of the functionality of the tool, here we show the outputs re-  
**355** ported by MultiGWAS in the tetraploid potato diversity panel, genotyped and  
**356** phenotyped as part of the USDA-NIFA Solanaceae Coordinated Agricultural Project  
**357** (SolCAP) [16]. The complete report from MultiGWAS for the naive and full  
**358** model is in the Supplementary information (<https://github.com/agrosavia-bioinformatics/multiGWAS>)

**360    3.1 Manhattan and QQ plots for GWAS associations**

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

**371** MultiGWAS adds special marks to the Manhattan and QQ plots to help iden-  
**372** tify different types of SNPs: (a) In Manhattan plots, significant SNPs are above  
**373** a red line, best-ranked SNPs are above a blue line, and shared SNPs (See Figure

374 6.b) are colored in green (b) In QQ plots, a red diagonal line indicates the ex-  
 375 pectation, so potential associations can be observed when the number of SNPs  
 376 deviating from the diagonal is small, as in the case of monogenic traits, or when  
 377 this number is somewhat higher, as in the case of truly polygenic traits. How-  
 378 ever, deviations for a high number of SNPs could reflect inflated  $p$ -values owing  
 379 to population structure or cryptic relatedness.



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

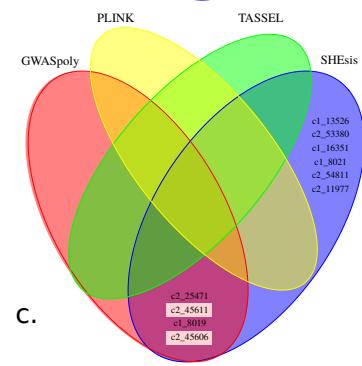
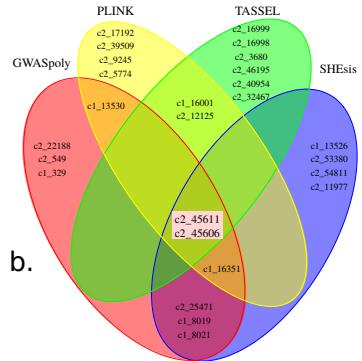
### 380 3.2 Tables and Venn diagrams for single and shared SNPs

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

385 First, best-ranked SNPs correspond to the top-scored  $N$  SNPs, wheter they  
 386 were assesed significant or not by its package, and with  $N$  defined by the user  
 387 in the configuration file. These SNPs are shown both in a SNPs table (Figure  
 388 6.a) and in a Venn diagram (Figure 6.b). The table lists them by package and  
 389 sorts by decreasing score, whereas the Venn diagram shows them emphasizing if  
 390 they were best-ranked either in a single package or in several at once (shared).  
 391 And second, the significant SNPs correspond to the ones assesed statistically  
 392 significant by each package, they are shown in a Venn diagram (Figure 6.c),  
 393 and they are also shown in the SNPs table, marked with significance TRUE (T)  
 394 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	F
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	F
GWASPoly	additive	0.96	c2_22188	11	40777	0.000050	4.30	4.50	F
GWASPoly	additive	0.96	c2_549	9	16527	0.000580	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	0.000786	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	F
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	F
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	F
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



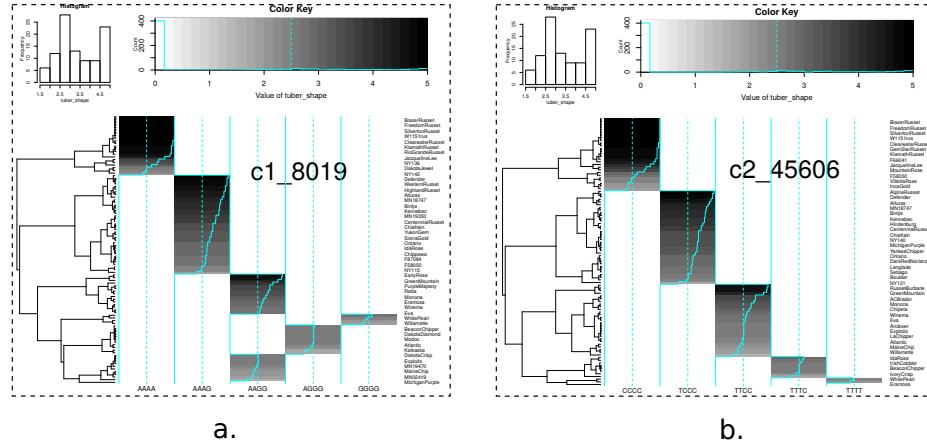
**Figure 6: Shared SNPs Views.** Tabular and graphical views of SNP associations identified by one or more GWAS packages (shared SNPs). SNPs identified by all packages are marker with red background in all figures **(a)** Table with details of the N=9 best-ranked SNPs from each GWAS package. Each row 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  $\alpha$  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).

### 3.3 Heat diagrams for structure of shared SNPs

396 MultiGWAS creates a two-dimensional representation, called SNP profile, to  
 397 visualize each trait by individuals and genotypes as rows and columns, respec-  
 398 tively (Figure 7). At the left, the individuals are grouped in a dendrogram by  
 399 their genotype. At the right, there is the name or ID of each individual. At the  
 400 bottom, the genotypes are ordered from left to right, starting from the major  
 401 to the minor allele (i.e., AAAA, AAAB, AABB, ABAA, BBBB). At the top, there  
 402 is a description of the trait based on a histogram of frequency (top left) and  
 403 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.

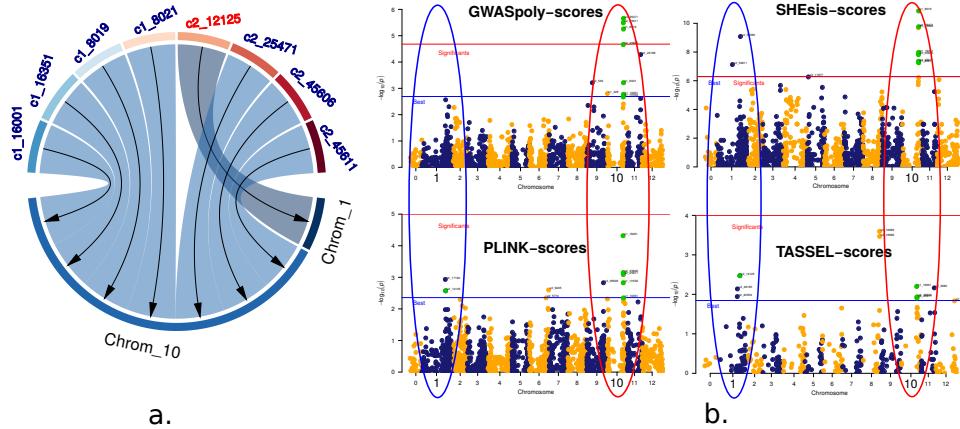


**Figure 7: SNP profiles.** SNP profiles for two of the best-ranked significant SNPs shown in the figure 6.b. (a) SNP c2\_45606 best-ranked by the four packages (central intersection of the Venn diagram Figure 6.b) (b) SNP c1\_8019 best-ranked by the two tetraploid packages (Figure 6.b), and also identified as significant by the same packages (at the bottom of the Figure 6.a).

### 410 3.4 Chord diagrams for SNPs by chromosome

411 Generally, in a typical GWAS analysis the strongest associations are signaled by  
 412 several nearby-correlated SNPs located in the same chromosome, as in manhat-  
 413 tan plots, where these associations form neat peaks with several SNPs showing  
 414 the same signal. Conversely, no peaks are shown when few SNPs correlate with  
 415 a trait.

416 However, when the analysis is performed by several GWAS packages, as  
 417 MultiGWAS does, it can identify correlated SNPs between packages that show  
 418 the same signal, what is presented by MultiGWAS through chord diagrams. For  
 419 example, the Figure 8.a shows the chord diagram for the shared SNPs from  
 420 the best-ranked associations previously described in the Figure 6.b. It can be  
 421 observed that most SNPs relate to chromosome 10 and only one to chromosome  
 422 1, which is also observed in the manhattan plots from each GWAS package  
 423 (Figure 8.b).



**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 packages showing two important locations of associations: chromosome 1 and chromosome 10, marked with a blue and red ellipsis, respectively.

## 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 10). Source code, examples, documentation and installation instructions are available at <https://github.com/agrosavia-bioinformatics/multiGWAS>.

### 4.1 Input parameters

MultiGWAS 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 file can be created either using a general text editor or using the MultiGWAS GUI application (see below). In both cases, the file must have the structure shown in the Figure 9.a, where parameter names and values are separated by colon, filenames are enclosed in quotation marks, and TRUE or FALSE indicates whether filters are applied or not. In the second case, the user creates the config file in a simple and straightforward 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.

## 440    4.2 Using the command line interface

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

445    multiGWAS Test01.config

446    Then, the tool starts the execution, showing information of the process in  
 447    the console window, and when it finishes the results are saved to a new sub-  
 448    folder called “out-Test01. Results include a full html report containing the dif-  
 449    ferent views described in the results section, along with the original graphics  
 450    and summary tables created by MultiGWAS and used to create the html report.  
 451    Additionally, results include the preprocessed tables of the main outputs gener-  
 452    ated by the four GWAS packages used by MultiGWAS.

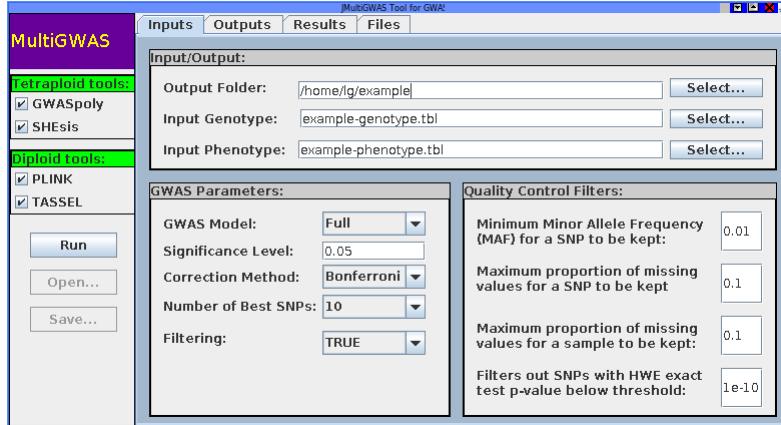
## 453    4.3 Using the graphical user interface

454    The MultiGWAS GUI can be executed by calling from a linux console the follow-  
 455    ing command:

456    jmultiGWAS

457    After it opens, it shows a main frame with a tool bar at left and four tabs  
 458    at the top (Figure 10). From the tool bar, users can select the GWAS packages  
 459    to use in the analysis—two for tetraploids and two for diploids—, and start the  
 460    analysis with the current parameters (or with parameters from a previous con-  
 461    figuration). And, from the tabs, users can input the MultiGWAS parameters,  
 462    and view the process and results of the analysis.

463



**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.

464

## 5 Discussion

465

XXXXXXXXXXXXXXXXXXXXXX

466

## References

467

- [1] María F. Álvarez et al. “Identification of Novel Associations of Candidate Genes with Resistance to Late Blight in Solanum tuberosum Group Phureja”. In: *Frontiers in Plant Science* 8 (2017), p. 1040. ISSN: 1664-462X. DOI: 10.3389/fpls.2017.01040. URL: <http://journal.frontiersin.org/article/10.3389/fpls.2017.01040/full>.
- [2] Ferdouse Begum et al. “Comprehensive literature review and statistical considerations for GWAS meta-analysis”. In: *Nucleic acids research* 40.9 (2012), pp. 3777–3784.
- [3] Jhon Berdugo-Cely et al. “Genetic diversity and association mapping in the colombian central collection of solanum tuberosum L. Andigenum group using SNPs markers”. In: *PLoS ONE* 12.3 (2017). ISSN: 19326203. DOI: 10.1371/journal.pone.0173039.
- [4] Peter M. Bourke et al. “Tools for Genetic Studies in Experimental Populations of Polyploids”. In: *Frontiers in Plant Science* 9 (2018), p. 513. ISSN: 1664-462X. DOI: 10.3389/fpls.2018.00513. URL: <http://journal.frontiersin.org/article/10.3389/fpls.2018.00513/full>.

- 483 [5] Peter J Bradbury et al. “TASSEL: software for association mapping of  
484 complex traits in diverse samples”. In: *Bioinformatics* 23.19 (2007), pp. 2633–  
485 2635. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btm308. URL:  
486 <https://doi.org/10.1093/bioinformatics/btm308>.
- 487 [6] Rita M Cantor, Kenneth Lange, and Janet S Sinsheimer. “Prioritizing GWAS  
488 results: a review of statistical methods and recommendations for their  
489 application”. In: *The American Journal of Human Genetics* 86.1 (2010),  
490 pp. 6–22.
- 491 [7] Jun Cao et al. “Whole-genome sequencing of multiple *Arabidopsis thaliana*  
492 populations”. In: *Nature genetics* 43.10 (2011), p. 956.
- 493 [8] Christopher C. Chang et al. “Second-generation PLINK: Rising to the chal-  
494 lenge of larger and richer datasets”. In: *GigaScience* 4.1 (2015), pp. 1–16.  
495 ISSN: 2047217X. DOI: 10.1186/s13742-015-0047-8. arXiv: 1410.4803.
- 496 [9] Rishika De, William S Bush, and Jason H Moore. “Bioinformatics Chal-  
497 lenges in Genome-Wide Association Studies (GWAS)”. In: *Clinical Bioin-  
498 formatics*. Ed. by Ronald Trent. New York, NY: Springer New York, 2014,  
499 pp. 63–81. ISBN: 978-1-4939-0847-9. DOI: 10.1007/978-1-4939-0847-  
500 9\_5. URL: <https://doi.org/10.1007/978-1-4939-0847-9%7B%5C-%7D5>.
- 501 [10] Robert Eklblom and Juan Galindo. “Applications of next generation se-  
502 quencing in molecular ecology of non-model organisms”. In: *Heredity*  
503 107.1 (2011), pp. 1–15.
- 504 [11] Hans Ellegren. “Genome sequencing and population genomics in non-  
505 model organisms”. In: *Trends in ecology & evolution* 29.1 (2014), pp. 51–  
506 63.
- 507 [12] Luís Felipe V. Ferrão et al. “Insights Into the Genetic Basis of Blueberry  
508 Fruit-Related Traits Using Diploid and Polyploid Models in a GWAS Con-  
509 text”. In: *Frontiers in Ecology and Evolution* 6 (2018), p. 107. ISSN: 2296-  
510 701X. DOI: 10.3389/fevo.2018.00107. URL: <https://www.frontiersin.org/articles/10.3389/fevo.2018.00107/full>.
- 511 [13] Dominik G Grimm et al. “easyGWAS: A Cloud-Based Platform for Com-  
512 paring the Results of Genome-Wide Association Studies”. In: *The Plant  
513 Cell* 29.1 (2017), pp. 5–19. ISSN: 1040-4651. DOI: 10.1105/tpc.16.  
514 00551. URL: <http://www.plantcell.org/content/29/1/5>.
- 515 [14] Anja C Gumpinger et al. *Methods and Tools in Genome-wide Association  
516 Studies*. Vol. 1819. 2018. ISBN: 9781493986187.
- 517 [15] Bin Han and Xuehui Huang. “Sequencing-based genome-wide associa-  
518 tion study in rice”. In: *Current opinion in plant biology* 16.2 (2013),  
519 pp. 133–138.
- 520 [16] Candice N. Hirsch et al. “Retrospective view of North American potato  
521 (*Solanum tuberosum* L.) breeding in the 20th and 21st centuries”. In: *G3:  
522 Genes, Genomes, Genetics* 3.6 (2013), pp. 1003–1013. ISSN: 21601836.  
523 DOI: 10.1534/g3.113.005595.
- 524
- 525

- 526 [17] “How to interpret a genome-wide association study”. In: *JAMA - Journal*  
527 *of the American Medical Association* 299.11 (2008), pp. 1335–1344. ISSN:  
528 00987484. DOI: 10.1001/jama.299.11.1335.
- 529 [18] Avjinder S Kaler and Larry C Purcell. “Estimation of a significance thresh-  
530 old for genome-wide association studies”. In: *BMC Genomics* 20.1 (2019),  
531 p. 618. ISSN: 1471-2164. DOI: 10 . 1186 / s12864 - 019 - 5992 - 7. URL:  
532 <https://doi.org/10.1186/s12864-019-5992-7>.
- 533 [19] Arthur Korte and Ashley Farlow. “The advantages and limitations of trait  
534 analysis with GWAS: a review”. In: *Plant methods* 9.1 (2013), p. 29.
- 535 [20] Gordan Lauc et al. “Genomics meets glycomics—the first GWAS study  
536 of human N-glycome identifies HNF1 $\alpha$  as a master regulator of plasma  
537 protein fucosylation”. In: *PLoS genetics* 6.12 (2010).
- 538 [21] Jie Meng et al. “Genome-wide association analysis of nutrient traits in  
539 the oyster *Crassostrea gigas*: Genetic effect and interaction network”. In:  
540 *BMC Genomics* 20.1 (2019), pp. 1–14. ISSN: 14712164. DOI: 10 . 1186 /  
541 s12864-019-5971-z.
- 542 [22] Robert A. Power, Julian Parkhill, and Tilio De Oliveira. “Microbial genome-  
543 wide association studies: lessons from human GWAS”. In: *Nature Reviews*  
544 *Genetics* 18.1 (2016), pp. 41–50. ISSN: 14710064. DOI: 10 . 1038 / nrg .  
545 2016. 132.
- 546 [23] Shaun Purcell et al. “PLINK: A tool set for whole-genome association and  
547 population-based linkage analyses”. In: *American Journal of Human Genetics* 81.3 (2007), pp. 559–575. ISSN: 00029297. DOI: 10 . 1086 / 519795.
- 549 [24] Hui Ping Qiao et al. “Genetic variants identified by GWAS was associated  
550 with colorectal cancer in the Han Chinese population”. In: *Journal of Can-  
551 cer Research and Therapeutics* 11.2 (2015), pp. 468–470. ISSN: 19984138.  
552 DOI: 10 . 4103 / 0973 - 1482 . 150346.
- 553 [25] Umesh R. Rosyara et al. “Software for Genome-Wide Association Studies  
554 in Autopolyploids and Its Application to Potato”. In: *The Plant Genome* 9.2  
555 (2016), pp. 1–10. ISSN: 1940-3372. DOI: 10 . 3835 / plantgenome2015 .  
556 08 . 0073. URL: <https://dl.sciencesocieties.org/publications/tpg/abstracts/9/2/plantgenome2015.08.0073>.
- 558 [26] Anna W Santure and Dany Garant. “Wild GWAS—association mapping in  
559 natural populations”. In: *Molecular ecology resources* 18.4 (2018), pp. 729–  
560 738.
- 561 [27] Sanjeev Kumar Sharma et al. “Linkage disequilibrium and evaluation of  
562 genome-wide association mapping models in tetraploid potato”. In: *G3:*  
563 *Genes, Genomes, Genetics* 8.10 (2018), pp. 3185–3202. ISSN: 21601836.  
564 DOI: 10 . 1534 / g3 . 118 . 200377.
- 565 [28] Jiawei Shen et al. “SHEsisPlus, a toolset for genetic studies on polyploid  
566 species”. In: *Scientific Reports* 6 (2016), pp. 1–10. ISSN: 20452322. DOI:  
567 10 . 1038 / srep24095. URL: <http://dx.doi.org/10.1038/srep24095>.

- 568 [29] John R Thompson, John Attia, and Cosetta Minelli. “The meta-analysis  
569 of genome-wide association studies”. In: *Briefings in Bioinformatics* 12.3  
570 (2011), pp. 259–269. ISSN: 1467-5463. DOI: 10.1093/bib/bbr020. URL:  
571 <https://doi.org/10.1093/bib/bbr020>.
- 572 [30] Feng Tian et al. “Genome-wide association study of leaf architecture in  
573 the maize nested association mapping population”. In: *Nature genetics*  
574 43.2 (2011), pp. 159–162.
- 575 [31] Yan Y. Yan et al. “Effects of input data quantity on genome-wide asso-  
576 ciation studies (GWAS)”. In: *International Journal of Data Mining and*  
577 *Bioinformatics* 22.1 (2019), pp. 19–43. ISSN: 17485681. DOI: 10.1504/  
578 IJDMB.2019.099286.
- 579 [32] J Yu et al. “A unified mixed-model method for association mapping that  
580 accounts for multiple levels of relatedness.” In: *Nature genetics* 38.2 (2006),  
581 pp. 203–208.
- 582 [33] Jiazheng Yuan et al. “Genome-Wide Association Study of Resistance to  
583 Potato Common Scab”. In: *Potato Research* (2019). ISSN: 18714528. DOI:  
584 10.1007/s11540-019-09437-w.
- 585 [34] Shengkui Zhang et al. “Genome-wide association studies of 11 agronomic  
586 traits in cassava (*Manihot esculenta crantz*)”. In: *Frontiers in Plant Science*  
587 9.April (2018), pp. 1–15. ISSN: 1664462X. DOI: 10.3389/fpls.2018.  
588 00503.