

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

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

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

12 June 11, 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

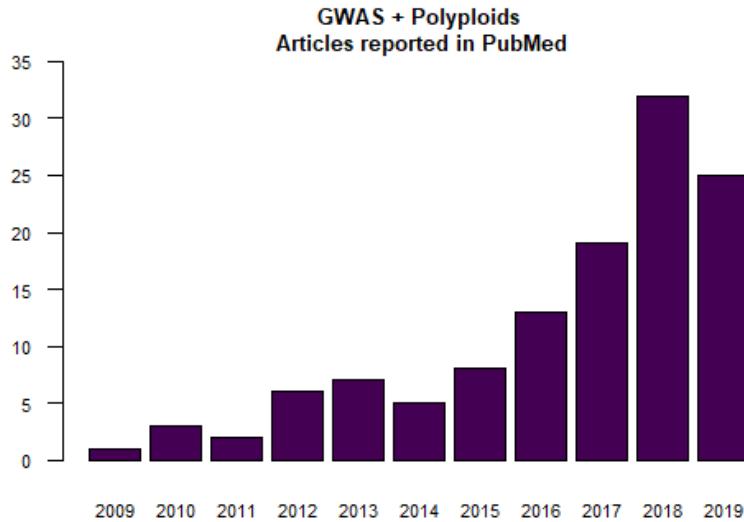


Figure 1: Timeline for articles reported for GWAS studies on polyploid species in PubMed. We present data for completed years.

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 [?, ?]. This methodology started with humans and several model
 45 plants, such as rice, maize, and *Arabidopsis* [?, ?, ?, ?, ?]. Because of the ad-
 46 vances in the next-gen sequencing technology and the decline of the sequencing
 47 cost in recent years, there is an increase in the availability of genome sequences
 48 of different organisms at a faster rate [?, ?]. Thus, the GWAS is becoming the
 49 standard tool to understand the genetic bases of either ecological or economic
 50 phenotypic variation for both model and non-model organisms. This increment
 51 in GWAS includes complex species such as polyploids (Fig 1) [?, ?].

52 The GWAS for polyploid species has three related challenges. First, as all
 53 GWAS, we should replicate the study as a reliable method to validate the results
 54 and recognize real associations. This replication involves finding the same as-

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

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

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

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

81 To solve all the three challenges above, we developed the MultiGWAS tool
82 that performs GWAS analyses for tetraploid species using four software in paral-
83 lel. Our tool include GWASPoly [?] and the SHEsis tool [?] that accept polyploid
84 genomic data, and PLINK [?] and TASSEL [?] with the use of a "diploidized" ge-
85 nomic matrix. The tool deals with preprocessing data, running four GWAS tools
86 in parallel, and create comparative reports from the output of each software to
87 help the user to decide more intuitively the true or false associations.

88 **2 Method**

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

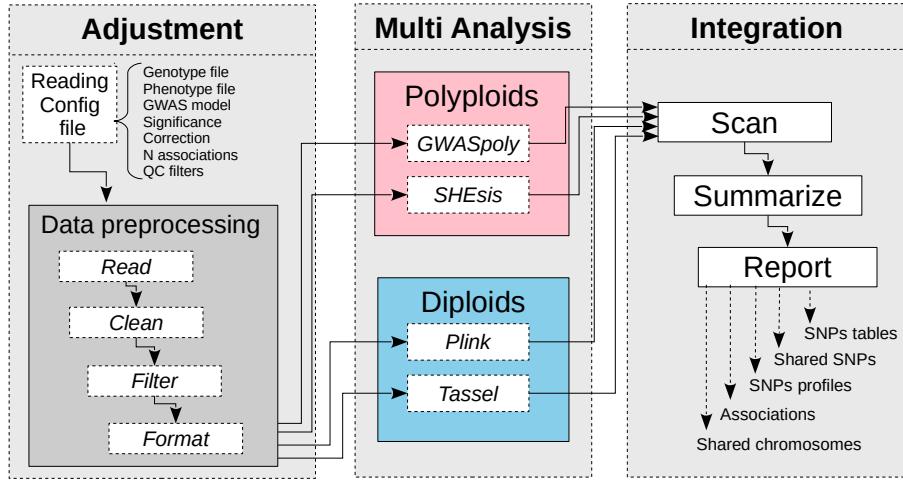


Figure 2: MultiGWAS flowchart has three consecutive steps: adjustment, multi analysis, and integration. The adjustment step manages the input data, reads the configuration file, and preprocesses the input genomic data (genotype and phenotype). The multi analysis step configures and runs the four GWAS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

98 2.1 Adjustment stage

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

104 **2.1.1 Reading configuration file**

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

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

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

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

123 The vector y represents the observed phenotypes depends on the following
 124 factors: the fixed effect vector β , the SNP effects vector α , the population ef-
 125 fect vector ν , the polygene background effect vector μ , and, the residual effect
 126 vector e . The Q , modeled as a fixed effect, refers to the incidence matrix for
 127 subpopulation covariates relating y to ν . Moreover, X , S , and Z are incidence
 128 matrices of ones and zeros relating y to β , α , and μ , respectively.

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

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

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

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

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

- 163 • **MAF of x :** filters out SNPs with minor allele frequency below x (default
164 0.01);
- 165 • **MIND of x :** filters out all individuals with missing genotypes exceeding
166 $x^*100\%$ (default 0.1);
- 167 • **GENO of x :** filters out SNPs with missing values exceeding $x^*100\%$ (de-
168 fault 0.1);
- 169 • **HWE of x :** filters out SNPs which have Hardy-Weinberg equilibrium exact
170 test *p-value* below the x threshold.

171 MultiGWAS does the MAF filtering and uses the PLINK package [?] for the other
172 three filters: MIND, GENO, and HWE.

173 2.1.2 Data preprocessing

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

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

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

186 2.2 Multi analysis stage

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

197 2.2.1 GWASpoly

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

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

214 2.2.2 SHEsis

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

217 ploid species. However, their use is mainly in diploids animals and humans
218 [?, ?].

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

223 2.2.3 PLINK

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

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

233 2.2.4 TASSEL

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

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

247 2.3 Integration stage.

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

251 2.3.1 Calculation of *p*-values and significance thresholds

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

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

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

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

277 2.3.2 Selection of significant and best-ranked associations

278 After corrections, significant associations are selected as the ones with *p*-values
279 falling below a significant threshold on each GWAS package. Nevertheless, as
280 described above, it is equally important to know the best-ranked associations
281 closer to being statistically significant, as they may represent associations to
282 consider for posterior analysis.

283 In the case of GWAS packages with only one gene action model (PLINK and
284 SHESIS), the best-ranked associations are the top N identified by the package.
285 However, in GWAS packages with several gene action models (GWASpoly and
286 TASSEL), the best-ranked associations are selected as the top N from the “best
287 action model”, the one with more shared SNP associations. In other words,
288 from the associations identified in more than one model.

289 2.3.3 Integration of results

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

- 293 • A Quantile-Quantile (QQ) plots for the resultant *p*-values of each tool and
294 the corresponding inflation factor λ to assess the degree of the test statistic
295 inflation.

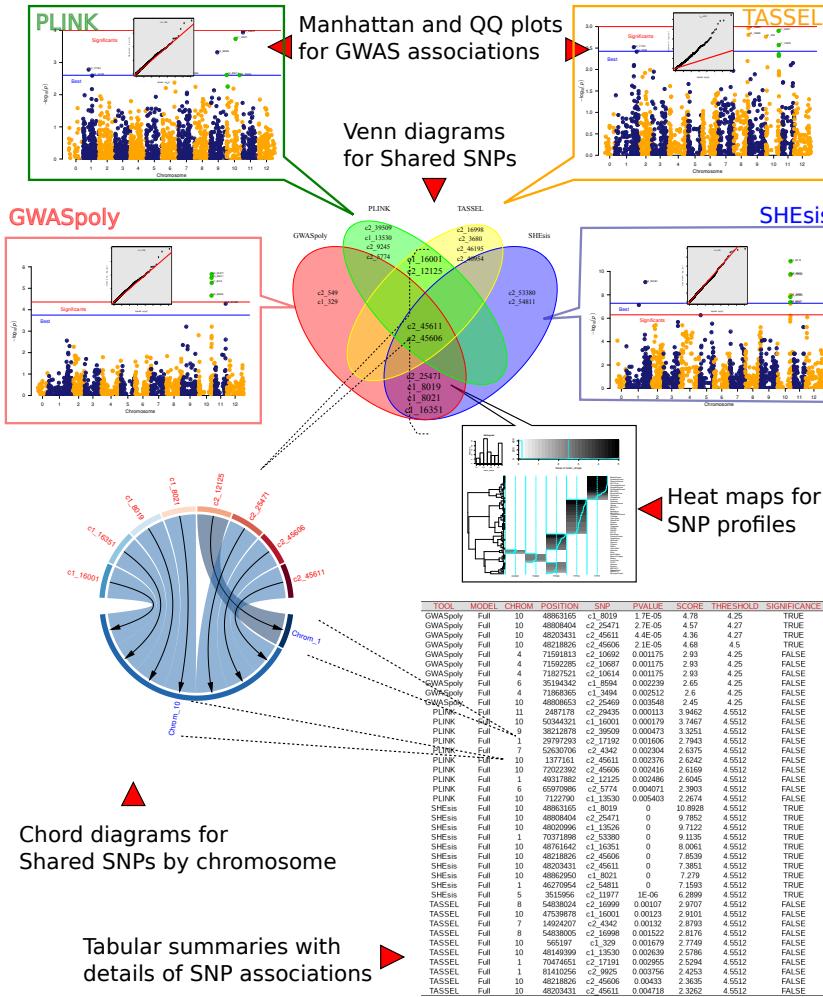


Figure 4: Reports presented by MultiGWAS. For each tool, first a QQ plot that assesses the resultant p-values. Second, a Manhattan plot for each tool with two lines, blue and red, respectively, is the lower limit for the best ranked and significative SNPs. We present two Venn diagrams, one for the significative 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.

- 296 • A Manhattan plot of each tool with two lower thresholds, one for the best-
 297 ranked SNPs, and another for the significant SNPs.
- 298 To present the replicability, we use two sets: (1) the set of all the significative
 299 SNPs provided by each tool and (2) the set of all the best-ranked SNPs. For
 300 each set, we present a Venn diagram that shows SNPs predicted exclusively by
 301 one tool and intersections that help to identify the SNPs predicted by one, two,
 302 three, or all the tools. Also, we provide detailed tables for the two sets.
 303 For each SNP **identified** more than once, we provide what we call the SNP
 304 profile. That is a heat diagram for a specific SNP, where each column is a geno-
 305 type state AAAA, AAAB, AABB, ABBB, **and** BBBB. Moreover, each row corre-
 306 sponds to a sample. Samples with close genotypes form together clusters. Thus
 307 to generate the clusters, we do not use the phenotype information. However,
 308 we present the phenotype information in the figure as the color. This figure
 309 visually provides information regarding genotype and phenotype information
 310 simultaneously for the whole population. We present colors as tones between
 311 white and black for color blind people.
 312 MultiGWAS generates a report, one document with the content previously
 313 described. Besides, there is a folder with the individual figures just in case the
 314 user needs one. In the supplementary information, we include a report and a
 315 description of the report content (**supplementary information XXX**)
 316 In the following section, we present the results applied to a public dataset.

3 Results

- 317 Most of the GWAS packages used by MultiGWAS are based on a linear regression
 318 approaches, but they often produce dissimilar association results for the same
 319 input. For example, computed *p-values* for the same set of SNPs are different
 320 between packages; SNPs with significant *p-values* for one package may be not
 321 significant for the others; or well-ranked SNPs in one package may be ranked
 322 differently in another.
 323 To alleviate these difficulties, MultiGWAS produces five types of outputs us-
 324 ing different graphics and tabular views, these outputs are intended to help
 325 users to compare, select, and interpret the set of possible SNPs associated with
 326 a trait of interest. The outputs include:
 327
 328 • Manhattan and Q-Q plots to show GWAS associations.
 329 • Venn diagrams to show associations identified by single or several tools.
 330 • Heat diagrams to show the genotypic structure of shared SNPs.
 331 • Chord diagrams to show shared SNPs by chromosomes.
 332 • Score tables to show detailed information of associations for both sum-
 333 mary results from MultiGWAS and particular results from each GWAS
 334 package

335 As an example of the functionality of the tool, here we show the outputs re-
336 ported by MultiGWAS in the tetraploid potato diversity panel, genotyped and
337 phenotyped as part of the USDA-NIFA Solanaceae Coordinated Agricultural Project
338 (SolCAP) [?]. The complete report from MultiGWAS for the naive and full
339 model is in the Supplementary information (<https://github.com/agrosavia-bioinformatics/multiGWAS>)
340

341 **3.1 Manhattan and QQ plots for GWAS associations**

342 MultiGWAS uses classical Manhattan and Quantile–Quantile plots (QQ plots)
343 to visualize the results of GWAS analysis from each package. In both plots,
344 SNPs are represented by dots and their *p-values* are transformed to scores as
345 $-\log_{10}(p\text{-values})$ (see Figure 5). The Manhattan plot displays the SNP asso-
346 ciation strength (y-axis) distributed in their genomic location (x-axis), so the
347 higher the score the stronger the association. Whereas the QQ plot is used to
348 visually compare the expected distribution of *p-values* (y-axis) vs. the observed
349 distribution (x-axis), so under the null hypothesis of no association of SNPs with
350 the phenotype, both distributions should coincide, and most SNPs should lie on
351 a diagonal line.

352 MultiGWAS adds special marks to the Manhattan and QQ plots to help iden-
353 tify different types of SNPs: (a) In Manhattan plots, significant SNPs are above
354 a red line, best-ranked SNPs are above a blue line, and shared SNPs (See Figure
355 6.b) are colored in green (b) In QQ plots, a red diagonal line indicates the ex-
356 pectation, so potential associations can be observed when the number of SNPs
357 deviating from the diagonal is small, as in the case of monogenic traits, or when
358 this number is somewhat higher, as in the case of truly polygenic traits. How-
359 ever, deviations for a high number of SNPs could reflect inflated *p-values* owing
360 to population structure or cryptic relatedness.

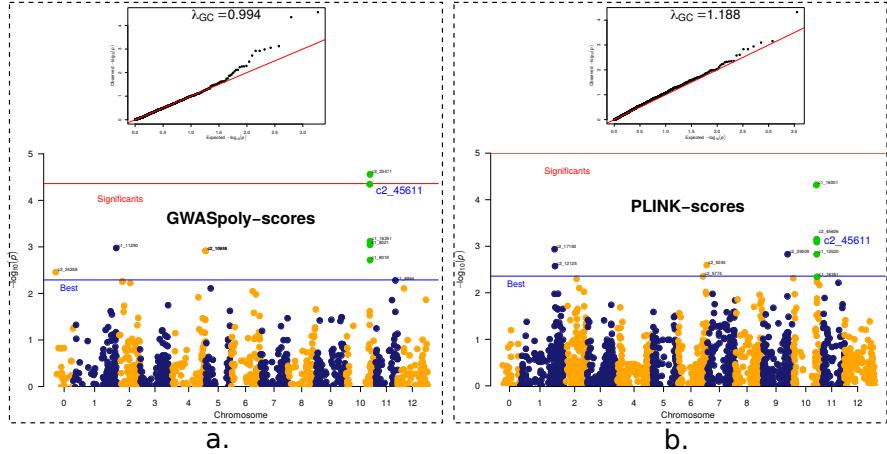


Figure 5: MultiGWAS visualization of associations. MultiGWAS creates Manhattan and QQ plots for GWAS results of each GWAS packages. Here we show the plots for one tetraploid package, GWASpoly (a), and other diploid package, PLINK (b).

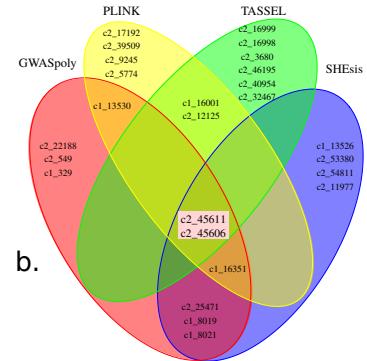
361 3.2 Tables and Venn diagrams for single and shared SNPs

362 MultiGWAS provides tabular and graphic views to report in an integrated way
 363 both the best-ranked and significant SNPs identified by the four GWAS pack-
 364 ages (see Figure 6). Both p -values and significance levels have been scaled as
 365 $-\log_{10}(p\text{-value})$ to give high scores to the best statistically evaluated SNPs.

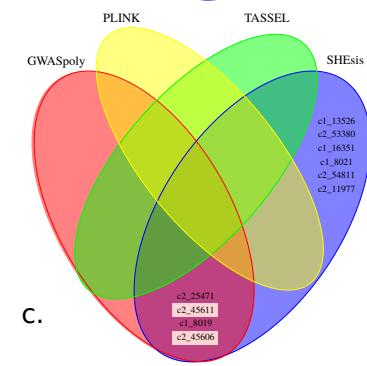
366 First, best-ranked SNPs correspond to the top-scored N SNPs, wheter they
 367 were assesed significant or not by its package, and with N defined by the user
 368 in the configuration file. These SNPs are shown both in a SNPs table (Figure
 369 6.a) and in a Venn diagram (Figure 6.b). The table lists them by package and
 370 sorts by decreasing score, whereas the Venn diagram shows them emphasizing if
 371 they were best-ranked either in a single package or in several at once (shared).
 372 And second, the significant SNPs correspond to the ones assesed statistically
 373 significant by each package, they are shown in a Venn diagram (Figure 6.c),
 374 and they are also shown in the SNPs table, marked with significance TRUE (T)
 375 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



b.



c.

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

3.3 Heat diagrams for structure of shared SNPs

376 MultiGWAS creates a two-dimensional representation, called SNP profile, to
 377 visualize each trait by individuals and genotypes as rows and columns, respec-
 378 tively (Figure 7). At the left, the individuals are grouped in a dendrogram by
 379 their genotype. At the right, there is the name or ID of each individual. At the
 380 bottom, the genotypes are ordered from left to right, starting from the major
 381 to the minor allele (i.e., AAAA, AAAB, AABB, ABAA, BBBB). At the top, there
 382 is a description of the trait based on a histogram of frequency (top left) and
 383 by an assigned color for each numerical phenotype value using a grayscale (top
 384

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.

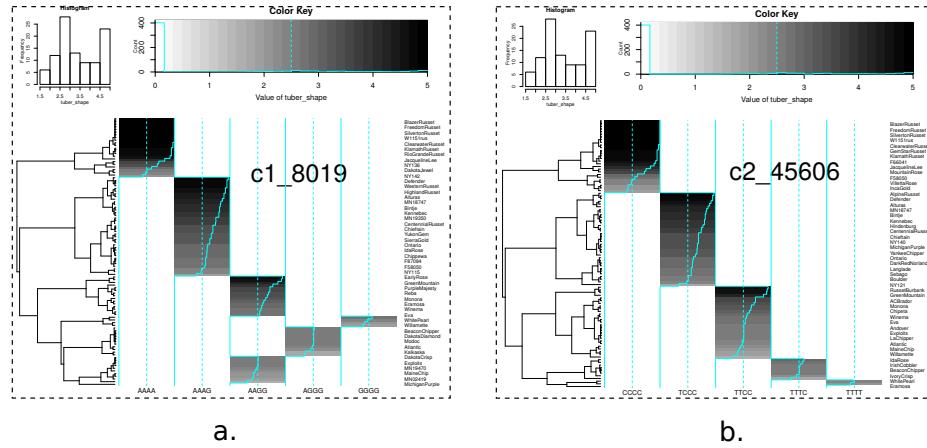


Figure 7: SNP profiles. SNP profiles for two of the best-ranked significant SNPs shown in the figure 6.b. (a) SNP c2_45606 best-ranked by the four packages (central intersection of the Venn diagram Figure 6.b) (b) SNP c1_8019 best-ranked by the two tetraploid packages (Figure 6.b), and also identified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

3.4 Chord diagrams for SNPs by chromosome

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

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

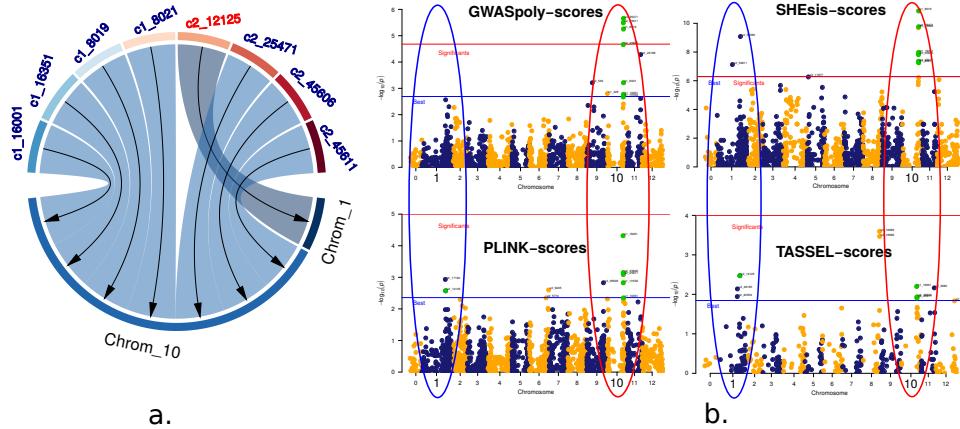


Figure 8: SNPs by chromosome. The figure shows how the best-ranked SNPs relate to chromosomes. (a) Chord diagram showing that most SNPs related to chromosome 10. SNPs are at the top of the diagram, chromosomes at the bottom, and associations are represented by arrows drawn from SNPs to their chromosomes. 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.

421 4.2 Using the command line interface

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

426 multiGWAS Test01.config

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

434 4.3 Using the graphical user interface

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

437 jmultiGWAS

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

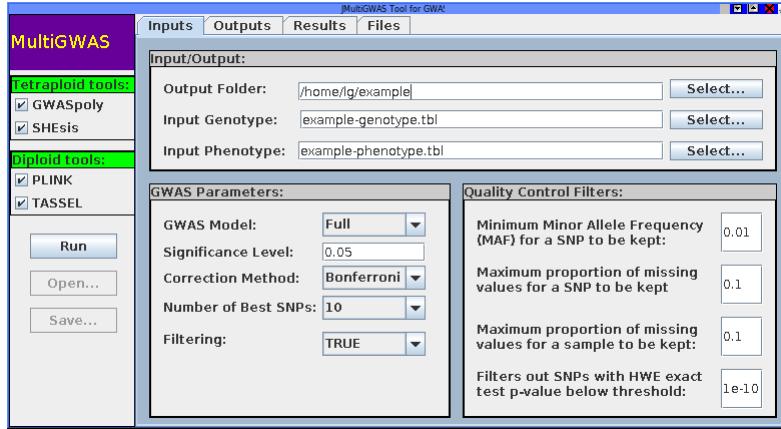


Figure 10: MultiGWAS GUI application. Main view of the MultiGWAS GUI application (“Inputs” view) where users can create the configuration file by setting values for input parameters. The GUI contains other three views: “Outputs” view shows the logs of the running process. “Results” view shows a report 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.

5 Discussion

XXXXXXXXXXXXXXXXXXXXXX