

1 **MultiGWAS: An integrative tool for Genome**
2 **Wide Association Studies (GWAS) in tetraploid**
3 **organisms**

4 immediate

5 July 15, 2020

6 **Abstract**

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

Comentatios Ivania

24 MultiGWAS generates five summary results integrating the four tools. (1)
25 Score tables with detailed information on the associations for each tool. (2)
26 Venn diagrams of shared SNPs among the four tools. (3) Heatmaps of signifi-
27 cantive SNP profiles among the four tools. (4) Manhattan and QQ plots for the
28 association found by each tool. And (5) Chord diagrams for the chromosomes
29 vs. SNP by each tool. **Contact:** phreyes@agrosavia.co

Comentatios Luis

31 **Keywords:** GWAS, tetraploids, SNPs, polyploids, software

33 **1 Introduction**

34 The **Genome-wide association studies (GWAS)** are used to identify which variants
35 through the whole genome of a large number of individuals are associated with

36 a specific trait ((?; ?)). This methodology started with humans and several model
 37 plants, such as rice, maize, and *Arabidopsis* ((?; ?; ?; ?; ?)). Because of the advances
 38 in the next-gen sequencing technology and the decline of the sequencing cost in re-
 39 cent years, there is an increase in the availability of genome sequences of different
 40 organisms at a faster rate ((?; ?)). Thus, the GWAS is becoming the standard tool
 41 to understand the genetic bases of either ecological or economic phenotypic varia-
 42 tion for both model and non-model organisms. This increment in GWAS includes
 43 complex species such as polyploids (Fig. ??) ((?; ?)).

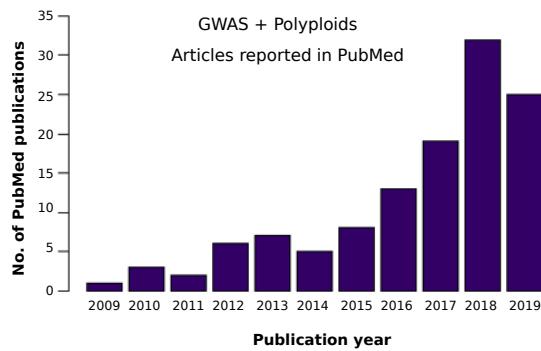


Figure 1: Timeline for articles reported for GWAS studies on polyploid species in PubMed. We present data for completed years.

44 The GWAS for polyploid species has fourth related challenges. First, as all
 45 GWAS, we should replicate the study as a reliable method to validate the results
 46 and recognize real associations. This replication involves finding the same associa-
 47 tions either in several replicates from the study population using the same software
 48 or testing different GWAS tools among the same study population. This approach
 49 involved the use of different parameters, models, or conditions, to test how consis-
 50 tent the results are ((?; ?)). However, the performance of different GWAS software
 51 could affect the results. For example, the threshold *pvalue* for SNP significance
 52 change through four GWAS software (i.e., PLINK, TASSEL, GAPIT, and FaST-LMM)
 53 when sample size varies ((?)). It means that well-ranked SNPs from one package
 54 can be ranked differently in another.

55 Second, although there are many GWAS software available to repeat the analy-
 56 sis under different conditions ((?)), most of them are designed exclusively for the
 57 diploid data matrix ((?)). Therefore, it is often necessary to "diploidizing" the poly-
 58 ploid genomic data in order to replicate the analysis.

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

63 The iPAT allows running in a graphic interface three well-known command-line
 64 GWAS software such as GAPIT, PLINK, and FarmCPU ((?)) . However, the output
 65 from each package is separated. On the other hand, the easyGWAS allows running a

Me confunde al hablar de
software y toois

"The latter approach..."

Software or tools ??

66 GWAS analysis on the web using different algorithms. This analysis could run inde-
67 pendently of both the computer capacity and operating system. However, it needs
68 either several datasets available or a dataset with a large number of individuals to
69 make replicates in order to compare among algorithms. Moreover, the output from
70 different algorithms is separated ((?)). Thus, for both software iPAT and easyG-
71 WAS, the integrative and comparative outputs among software or algorithms are
72 missing.

73 Fourth, the GWAS on polyploids generates a new level of complexity to un-
74 derstand how allele dosage affect the phenotype expression on quantitative traits.
75 Therefore, any tool that compares among software but also models with different
76 allele dosage will contribute to gain a better understanding in how redundancy or
77 complex interaction among alleles affect the phenotype expression and the evolu-
78 tion of new phenotypes among polyploid species .

79 To contribute to sort out all the above fourth challenges, we developed the Multi-
80 GWAS tool that performs GWAS analyses for tetraploid species using four software
81 in parallel. Our tool include GWASPoly ((?)) and the SHEsis tool ((?)) that ac-
82 cept polyploid genomic data, and PLINK ((?)) and TASSEL ((?)) with the use of a
83 "diploidized" genomic matrix. The tool deals with [input file formats](#), [data prepro-](#)
84 [cessing](#), [search for associations by running](#) four GWAS tools in parallel, and [creation](#)
85 [of](#) comparative reports from the output of each software to help the user to decide
86 more intuitively the true or false associations.

Aqui falta la cita del paper
de blueberry y el paper de
Rosyara 2016

87 2 Method

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

96 2.1 Adjustment stage

97 MultiGWAS takes as input a configuration file where the user specifies the genomics
98 data along with the parameters that will be used by the four tools. Once the config-
99 uration file is read and processed, the genomic data files (genotype and phenotype)
100 are preprocessed by cleaning, filtering, and checking data quality. The output of this
101 stage corresponds to the inputs for the four programs at the Multi Analysis stage.

102 2.1.1 Reading configuration file

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

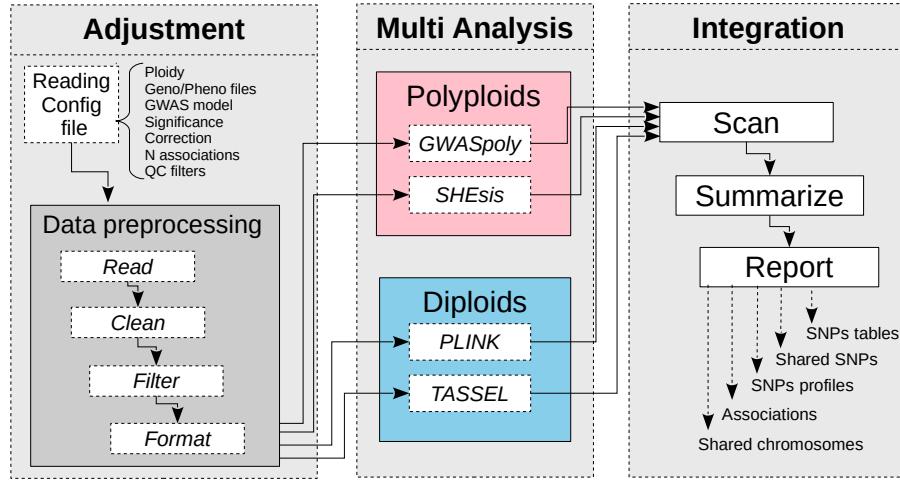


Figure 2: MultiGWAS flowchart has three steps: adjustment, multi analysis, and integration. The first step deals with input data management, reading the configuration file, **reading and preprocessing the input genomic data** (genotype and phenotype). The second step deals with GWAS analysis, configuring and running the four packages in parallel. And the third step deals with summarizing and reporting results using different tabular and graphical visualizations.

104 **Ploidy:** Numerical value for the ploidy level of the genotype, currently MultiGWAS
 105 supports diploids and tetraploids genotypes (2: for diploids, 4: for tetraploids).

106 **Genotype and phenotype input files:** MultiGWAS uses two input files, one for
 107 genotype and one for phenotype. Genotypes files can be either in GWASpoly format
 108 ((?)) using SNP markers in rows and samples in columns (Fig. ???.a) or Variant Call
 109 Format (VCF) (Fig. ???.b) which is transformed into GWASpoly format using NGSEP
 110 4.0.2 ((?)). The phenotype file contains only one trait and uses a matrix format
 111 with the first column for the sample names and the second column for the trait
 112 values (Fig. ???.c).

a.	<pre>Marker, Chrom, Pos, sample01, sample02, sample03, ... c2_41437, 0, 805179, AAAG, AAGG, AAGG, ... c2_24258, 0, 1252430, AAGG, AGGG, GGGG, ... c2_21332, 0, 3499519, TTCC, TTCC, TTCC, ...</pre>	c.	<pre>Individual,Trait sample01, 3.59 sample02, 4.07 sample03, 1.05</pre>
b.	<pre>##fileformat=VCFv4.2 ##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype" #CHROM POS ID REF ALT QUAL FILTER INFO FORMAT sample01 sample02 sample03 0 805179 c2_41437 A G . . PR GT 0/1/1/0 0/1/1/0 0/1/0/0 0 1252430 c2_24258 G A . . PR GT 0/1/0/0 0/1/1/0 0/0/1/0 0 3499519 c2_21332 T C . . PR GT 0/1/1/0 0/1/1/1 0/1/1/0</pre>		

Figure 3: Examples of MultiGWAS input file formats. Figures a and b show genotype files in GWASpoly and VCF formats, respectively, while figure c shows a phenotype file in matrix format. a. Genotype file in GWASpoly format containing column headers and with the first three columns for markers names, chromosomes and positions. The following columns correspond to the marker data of the samples in "ACGT" format (e.g. AAGG, CCTT for tetraploids, AG, CT for diploids). b. Genotype file in VCF format with metadata (first two lines) and header line. The following lines contain genotype information of the samples for each position. VCF marker data can be encoded as simple genotype calls (GT format field, e.g. 0/0/1/1 for tetraploids or 0/1 for diploids) or using the NGSEP custom format fields (?): ACN, ADP or BSDP. c. Phenotype file in matrix format with column headers and sample names followed by their trait values. Both GWASpoly genotype and phenotype files are in CSV (Comma Separated Values) format.).

113 **GWAS model:** MultiGWAS is designed to work with quantitative phenotypes and
 114 can run GWAS analysis using two types of statistical models that we have called *full*
 115 and *naive* models. The *full model* is known in the literature as the Q+K model (?) and
 116 includes control for structure (Q) and relatedness between samples (K), whereas
 117 the *naive model* does not include any type of correction. Both models are based
 118 on linear regression approaches and variations of them are implemented by the
 119 four GWAS packages used by MultiGWAS. The *naive* is modeled with Generalized
 120 Linear Models (GLMs, Phenotype + Genotype), and the *full* is modeled with Mixed
 121 Linear Models (MLMs, Phenotype + Genotype + Structure + Kinship). The default
 122 model used by MultiGWAS is the *full model* (Q+K) (?), which is expressed with the
 123 following equation:

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

124 where y is the vector of observed phenotypes; β is a vector of fixed effects other
 125 than SNP or population group effects; α is a vector of SNP effects (Quantitative
 126 Trait Nucleotides); ν is a vector of population effects; μ is a vector of polygene
 127 background effects; e is a vector of residual effects; Q , modeled as a fixed effect,
 128 refers to the incidence matrix for subpopulation covariates relating y to ν ; and X ,
 129 S and Z are incidence matrices of 1s and 0s relating y to β , α and μ , respectively.

130 **Genome-wide significance:** GWAS searches SNPs associated with the phenotype
 131 in a statistically significant manner. A threshold or significance level α is specified
 132 and compared with the p -value derived for each association score. Standard sig-
 133 nificance levels are 0.01 or 0.05 (?; ?), and MultiGWAS uses an α of 0.05 for the
 134 four GWAS packages. But the threshold is adjusted according to each package, as

135 some packages as GWASpoly and TASSEL calculates the SNP effect for each geno-
136 typic class using different gene action models (see “Multi analysis stage”). So, the
137 number of tested markers may be different in each model (see below) that results
138 in different *p-value* thresholds.

139 **Multiple testing correction:** Due to the massive number of statistical tests per-
140 formed by GWAS, it is necessary to perform a correction method for multiple hy-
141 pothesis testing and adjusting the *p-value* threshold accordingly. Two common
142 methods for multiple hypothesis testing are the false discovery rate (FDR) and the
143 Bonferroni correction. The latter is the default method used by MultiGWAS, which
144 is one of the most stringent methods. However, instead of adjusting the *p-values*,
145 MultiGWAS adjust the threshold below which a *p-value* is considered significant,
146 that is α/m , where α is the significance level and m is the number of tested markers
147 from the genotype matrix.

148 **Number of reported associations:** Criticism has arisen in considering only sta-
149 tistically significant associations as the only possible correct associations (?; ?).
150 Many of low *p-value* associations, closer to being significant, are discarded due to
151 the stringent significance levels, and consequently increasing the number of false
152 negatives. To help to analyze both significant and non-significant associations,
153 MultiGWAS provides the option to specify the number of best-ranked associations
154 (lower *p-values*), adding the corresponding *p-value* to each association found. In this
155 way, it is possible to enlarge the number of results, and we can observe replicabil-
156 ity in the results for different programs. Nevertheless, we present each association
157 with the corresponding *p-value*.

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

- 164 • **MAF of x:** filters out SNPs with minor allele frequency below x (default 0.01);
- 165 • **MIND of x:** filters out all individuals with missing genotypes exceeding $x*100\%$
166 (default 0.1);
- 167 • **GENO of x:** filters out SNPs with missing values exceeding $x*100\%$ (default
168 0.1);
- 169 • **HWE of x:** filters out SNPs which have Hardy-Weinberg equilibrium exact test
170 *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 **GWAS tools:** List of names of GWAS packages to run and integrate into MultiG-
174 WAS analysis. Currently four packages, two for tetraploid organisms: GWASpoly
175 and SHEsis, and two for diploids: PLINK and TASSEL.

176 **2.1.2 Data preprocessing**

177 Once the configuration file is processed, the genomic data is read and cleaned by se-
178 lecting individuals present in both genotype and phenotype. Then, individuals and
179 SNPs with poor quality are removed by considering the previous selected quality-
180 control filters and their thresholds,

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

187 After this process, the genomic data, genotype and phenotype, are converted to
188 the specific formats required for each of the four GWAS packages.

189 **2.2 Multi analysis stage**

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

199 **2.2.1 GWASpoly**

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

209 MultiGWAS is using GWASpoly version 1.3 with all gene action models available
210 to find associations. The MultiGWAS reports the top *N* best-ranked (the SNPs with
211 lowest *p-values*) that the user specified in the *N* input configuration file. The *full*
212 model used by GWASpoly includes the population structure and relatedness, which

213 are estimated using the first five principal components and the kinship matrix, re-
214 spectively, both calculated with the GWASpoly built-in algorithms.

215 **2.2.2 SHEsis**

216 SHEsis is a program based on a linear regression model that includes single-locus
217 association analysis, among others. The software design includes polyploid species.
218 However, their use is mainly in diploids animals and humans ((?; ?)).

219 MultiGWAS is using version 1.0, which does not take account for population
220 structure or relatedness. Despite, MultiGWAS externally estimates relatedness for
221 SHEsis by excluding individuals with cryptic first-degree relatedness using the al-
222 gorithm 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 univari-
226 ate GWAS using two-sample tests and linear regression models.

227 MultiGWAS is using two versions of PLINK: 1.9 and 2.0. Linear regression from
228 PLINK 1.9 performs both naive and full model. For the full model, the software
229 calculates the population structure using the first five principal components calcu-
230 lated with a built-in algorithm integrated into version 1.9. Moreover, version 2.0
231 calculates the kinship coefficients across individuals using a built-in algorithm that
232 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 developed
235 initially for maize but currently used in several species ((?; ?)). For the association
236 analysis, TASSEL includes the general linear model (GLM) and mixed linear model
237 (MLM) that accounts for population structure and relatedness. Moreover, as GWAS-
238 Poly, TASSEL provides three-gene action models to calculate the SNP effect of each
239 genotypic class: general, additive, and dominant, and so the significance threshold
240 depends on each action model.

241 MultiGWAS is using TASSEL 5.0, with all gene action models used to find the
242 N best-ranked associations and reporting the top N best-ranked associations (SNPs
243 with lowest p -values). Naive GWAS uses the GLM, and full GWAS uses the MLM
244 with two parameters: population structure that uses the first five principal compo-
245 nents, and relatedness that uses the kinship matrix with centered IBS method, both
246 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 to
249 identify both significant and best-ranked associations with p -values lower than and
250 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 below
254 a predefined significance threshold. However, the four GWAS packages compute
255 differently *p-values* with the consequence to compute them too high or too low. If
256 *p-values* is too high, it would lead to false negatives or SNPs with real associations
257 with the phenotype, but that does not reach the significance threshold. Conversely,
258 if *p-values* are too low, then it would lead to false positives or SNPs with false asso-
259 ciations with the phenotype, but that reaches the significance threshold.

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

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

275 **2.3.2 Selection of significant and best-ranked associations**

276 MultiGWAS selects two groups of associations from the results of each GWAS pack-
277 age: statistically significant and best-ranked. The latter equally important to the
278 former as they are associations with lowest *p-values* not reaching the significance
279 threshold but representing interesting associations for further analysis (possible
280 false negatives).

281 The significant associations are selected from the ones with *p-values* falling be-
282 low a significant threshold, calculated for each GWAS package; and the best-ranked
283 associations are selected as the closest N to being statistically significant, with N de-
284 fined by the user in the configuration file.

285 The selection of these groups takes into account whether the GWAS package
286 uses only one gene action model, as PLINK and SHEsis do, or uses several ones,
287 as GWASpoly and TASSEL do. In the first case, there is only one resulting set of
288 associations and the selection is straightforward, as described above. However, in
289 the second case, there are several resulting sets of associations, one for each model,
290 and MultiGWAS selects both groups of associations by choosing the gene action
291 model with the highest number of shared SNPs and with the inflation factor closest
292 to one, according to the following equation:

$$score(M_i) = \frac{\sum_{j=1}^k sharedSNPs(M_i, M_j)}{k * N} + 1 - |1 - \lambda(M_i)|$$

293 where $score(M_i)$ is the score for the gene action model M_i , with i from $1..k$, for
 294 a GWAS package with k gene action models; $sharedSNPs(M_i, M_j)$ is the number of
 295 shared SNPs between models M_i and M_j ; N is the number of closest SNPs to being
 296 statistically significant, as it was described above; and $\lambda(M_i)$ is the inflation factor
 297 for the model M_i .

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

303 2.3.3 Integration of results

304 At this stage, MultiGWAS integrates the results to evaluate reproducible results
 305 among tools (Fig ??). However, it still reports a summary of the results of each
 306 tool:

- 307 • A Quantile-Quantile (QQ) plots for the resultant *p-values* of each tool and
 308 the corresponding inflation factor λ to assess the degree of the test statistic
 309 inflation.
- 310 • A Manhattan plot of each tool with two lower thresholds, one for the best-
 311 ranked SNPs, and another for the significant SNPs.

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

317 For each SNP identified more than once, we provide what we call the SNP pro-
 318 file. That is a heat diagram for a specific SNP, where each column is a genotype
 319 state AAAA, AAAB, AABB, ABBC, and BBBB. Moreover, each row corresponds to a
 320 sample. Samples with close genotypes form together clusters. Thus to generate
 321 the clusters, we do not use the phenotype information. However, we present the
 322 phenotype information in the figure as the color. This figure visually provides in-
 323 formation regarding genotype and phenotype information simultaneously for the
 324 whole population. We present colors as tones between white and black for color
 325 blind people.

326 MultiGWAS generates a report, one document with the content previously de-
 327 scribed. Besides, there is a folder with the individual figures just in case the user
 328 needs one.

329 MultiGWAS generates a report, one document with the content previously de-
 330 scribed. Besides, there is a folder with the individual figures just in case the user

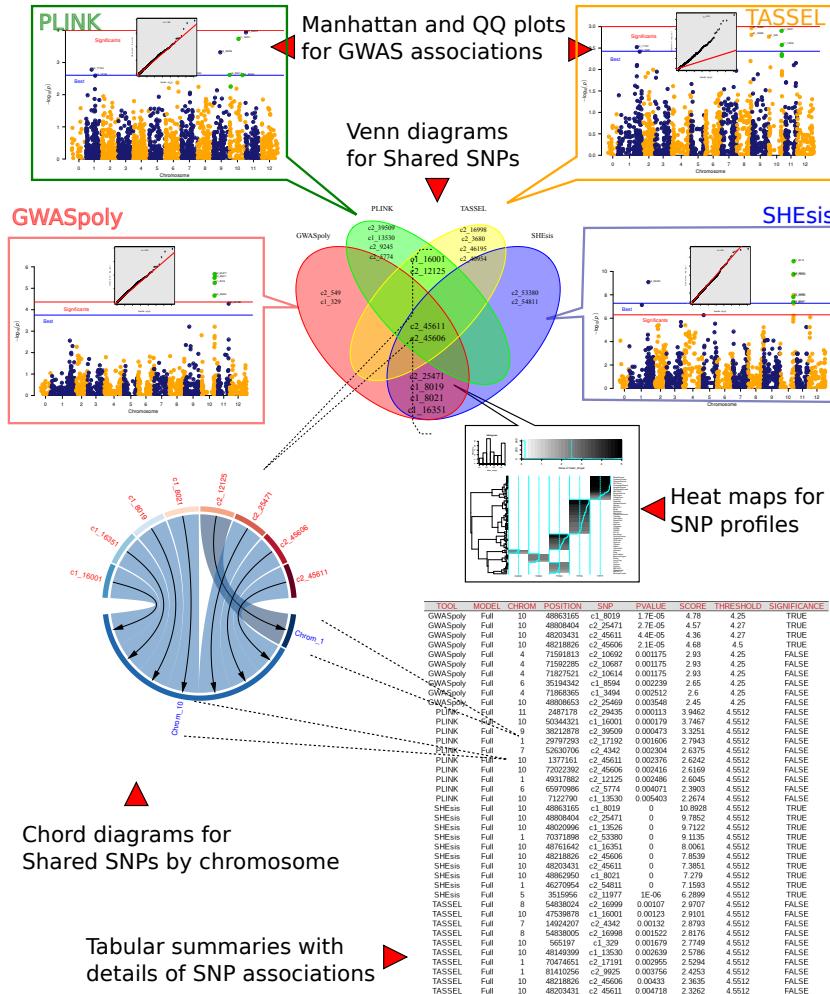


Figure 4: Reports presented by MultiGWAS. For each tool, first a QQ plot that assesses the resultant p-values. Second, a Manhattan plot for each tool with two lines, blue and red, respectively, is the lower limit for the best ranked and significative SNPs. We present two Venn diagrams, one for the significative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

331 needs one. In the supplementary information, we include a report and a description
332 of the report content (Supplementary Material 1)

333 In the following section, we present the results of the functionality of the tool
334 applied on a open dataset of a diversity panel of a tetraploid potato, genotyped and
335 phenotyped as part of the USDA-NIFA Solanaceae Coordinated Agricultural Project
336 (SolCAP) (?).

AquÃ deberÃamos decir
especÃficamente quÃ© datos
son y de donde salieron

337 3 Results

338 All four GWAS packages adopted by MultiGWAS use linear regression approaches.
339 However, they often produce different association results for the same input. Com-
340 puted *p-values* for the same set of SNPs are different between packages. Therefore,
341 SNPs with significant *p-values* for one package maybe not significant for the oth-
342 ers. Alternatively, well-ranked SNPs in one package may be ranked differently in
343 another.

344 To highlight these differences in the results across the four packages, MultiGWAS
345 produces five types of results combining graphics and tables to compare, select, and
346 interpret the set of possible SNPs associated with a trait of interest. The outputs
347 include:

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

354 The complete reports generated by MultiGWAS for both types of analysis: Full and
355 Naive, in the diversity panel of tetraploid potato described above are shown in the
356 supplementary information at <https://github.com/agrosavia-bioinformatics/multiGWAS-Supplementary>.

358 3.1 Manhattan and QQ plots for GWAS associations

359 MultiGWAS uses classical Manhattan and Quantile–Quantile plots (QQ plots) to
360 visualize the results of each package. In both plots, the points are the SNPs and
361 their *p-values* are transformed into scores like $-\log_{10}(p\text{-values})$ (see Fig. ??). The
362 Manhattan plot shows the strength of association of the SNPs (y-axis) distributed at
363 their genomic location (x-axis), so the higher the score, the stronger the association.
364 While the QQ plot compares the expected distribution of *p-values* (y-axis) with the
365 observed distribution (x-axis)..

366 MultiGWAS adds distinctive marks to both plots to identify different types of
367 SNPs: (a) In the Manhattan plots, the significant SNPs are above a red line and the

Corregido: "among packages"
y cita de polygenic traits

Revisar si esto del color
verde es correcto "shared
among all packages". Tam-
biÃ©n revisar si estÃ¡ bien
dicho lo de la comparaciÃ³n
entre pocos genes y muchos
genes. No se si sea necesario
cuantificar la desviaciÃ³n.
Originalmente se cuantificaba
pero era confuso

368 best-ranked SNPs are above a blue line. In addition, **SNPs shared between packages**
 369 **are colored green** (See Fig. ??-b). (b) In the QQ plots, a red diagonal line indicates
 370 the expected distribution under the null hypothesis of no association of SNPs with
 371 the phenotype, both distributions should coincide, and most SNPs should lie on the
 372 diagonal line. **Deviations for a large number of SNPs may reflect inflated *p*-values**
 373 **due to population structure or cryptic relatedness.** But, it is also expected that few
 374 **SNPs deviate from the diagonal for a truly polygenic trait (??).**

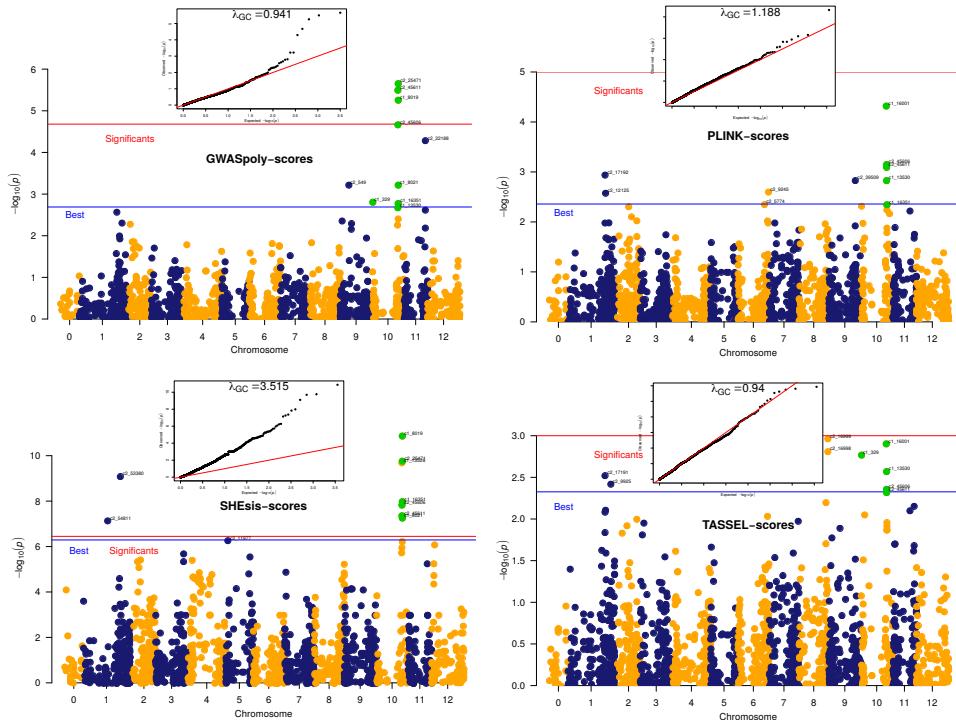


Figure 5: Associations in the tetraploid potato dataset. MultiGWAS shows the associations identified by the four GWAS packages using Manhattan and QQ plots. In the case of the tetraploid potato, several SNPs are observed to be shared between the four packages (green dots). The best-ranked SNPs are above the blue line, but only GWASpoly and SHEsis identified significant associations (SNPs above the red line). However, the inflation factor given by SHEsis is too high ($\lambda = 3.5$, at the top of the QQ plot), which is observed by the high number of SNPs deviating from the red diagonal of the QQ plot.

375 3.2 Tables and Venn diagrams for single and shared SNPs

376 MultiGWAS provides tabular and graphic views to report the best-ranked and signif-
 377 icant SNPs identified by the four GWAS packages in an integrative way (see Figure
 378 ??). Both *p*-values and significance levels have been scaled as $-\log_{10}(p\text{-value})$ to
 379 give high scores to the best statistically evaluated SNPs.

380 First, best-ranked SNPs correspond to the top-scored N SNPs, whether they were

381 assessed significant or not by its package, and with N defined by the user in the
 382 configuration file. These SNPs appears in both a SNPs table (Figure ??-a), and in a
 383 Venn diagram (Figure ??-b). The table lists them by package and sorts by decreasing
 384 score, whereas the Venn diagram emphasizes if they were best-ranked either in a
 385 single package or in several at once (shared).

386 Second, the significant SNPs correspond to the ones valued statistically signif-
 387 icant by each package. They appear in a Venn diagram (Figure ??-c), and in the
 388 SNPs table, marked with significance TRUE (T) in the table of the Figure ??-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	0.000589	3.23	4.50	F
GWASPoly	additive	0.96	c1_8021	10	48862	0.000589	3.23	4.50	F
GWASPoly	additive	0.96	c1_329	10	56519	0.001514	2.82	4.50	F
GWASPoly	additive	0.96	c1_16351	10	48761	0.001622	2.79	4.50	F
PLINK	additive	1.19	c1_16001	10	47539	0.000047	4.33	4.55	F
PLINK	additive	1.19	c2_45606	10	48218	0.000688	3.16	4.55	F
PLINK	additive	1.19	c2_45611	10	48203	0.000786	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

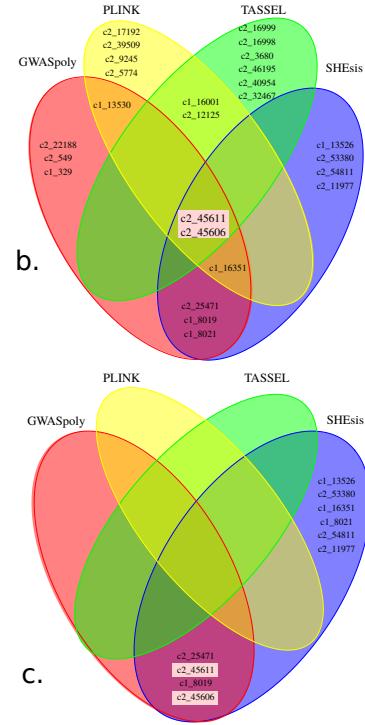


Figure 6: Shared SNPs Views. Tabular and graphical views of SNP associations identified by one or more GWAS packages (shared SNPs). SNPs identified by all packages are marker with red background in all figures. **(a)** Table with details of the N=9 best-ranked SNPs from each GWAS package. Each row corresponds to a single SNP, and the nine columns are tool name, the 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 in the central intersection. Other SNPs identified by more than one packages are in both upper central and lower central intersections. **(c)** Venn diagram of the significant SNPs (score $>$ threshold).

389 3.3 Heat diagrams for the structure of shared SNPs

390 MultiGWAS creates a two-dimensional representation, called the SNP profile, to vi-
 391 sualize each trait by individuals and genotypes as rows and columns, respectively
 392 (Figure ??). At the left, the individuals are grouped in a dendrogram by their geno-
 393 type. At the right, there is the name or ID of each individual. At the bottom, the
 394 genotypes are ordered from left to right, starting from the major to the minor allele
 395 (i.e., AAAA, AAAB, AABB, ABBB, BBBB). At the top, there is a description of the
 396 trait based on a histogram of frequency (top left) and by an assigned color for each
 397 numerical phenotype value using a grayscale (top right). Thus, each individual ap-
 398 pears as a colored line by its phenotype value on its genotype column. For each
 399 column, there is a solid cyan line with the mean of each column and a broken cyan
 400 line that indicates how far the cell deviates from the mean.

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

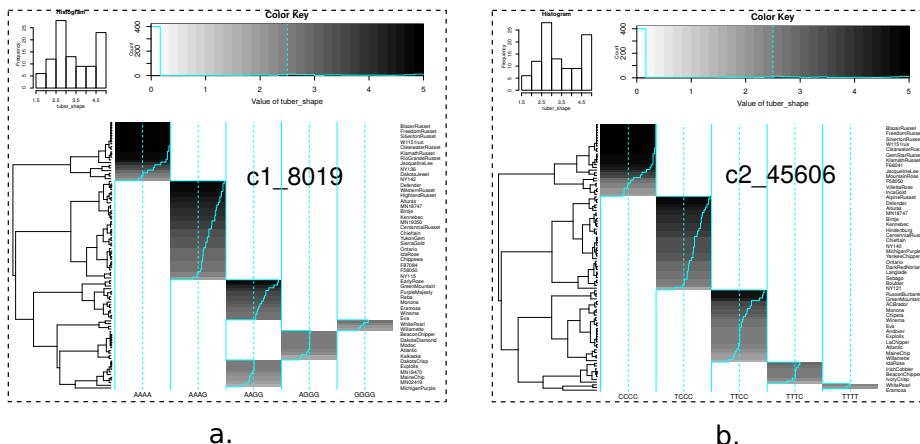


Figure 7: SNP profiles. SNP profiles for two of the best-ranked significant SNPs shown in the figure ??-b. (a) SNP c2_45606 best-ranked by the four packages (central intersection of the Venn diagram Figure ??-b.) (b) SNP c1_8019 best-ranked by the two tetraploid packages (Figure ??-b), and also identified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

403

404 3.4 Chord diagrams for SNPs by chromosome

405 The chord diagrams visualize the location across chromosomes of the best-ranked
 406 associated SNPs shared among the four packages and described in the table ??-a.
 407 Thus, **in the case of the tetraploid potato**, we found that they are located mostly in
 408 chromosome 10 (Figure ??-a). This visualization complements the manhattan plots
 409 from each GWAS package (Figure ??-b).

En esta sección fui bastante radical y terminé eliminando dos párrafos que no me parecen relevantes. Sin embargo los dejo señalados por si ustedes consideran que es importante dejarlo. Si se vuelen a incorporar al texto sugiero reorganizarlos mejor porque el mensaje es confuso en mi opinión. Párrafo 1. 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. Leyenda

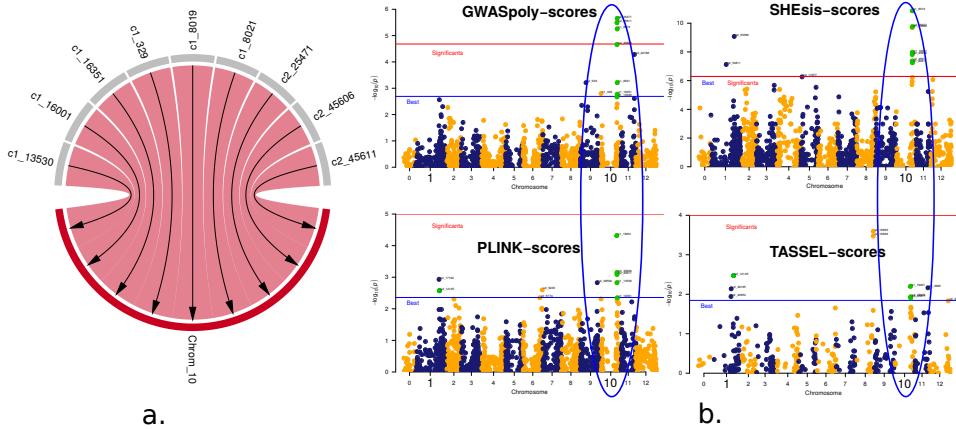


Figure 8: SNPs by chromosome. The position of best-ranked SNPs across chromosomes using two different visualizations. (a) Chord diagram showing that best-ranked SNPs located in chromosome 10. The SNPs are at the top and the chromosomes at the bottom. The arrows connect the best-ranked SNPs with their position in the chromosomes. (b) Manhattan plots from each GWAS packages showing two important locations of associations, chromosome 1 and chromosome 10, marked with blue and red ellipsis, respectively.

4 Availability and Implementation

The core of the MultiGWAS tool runs under 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 ??). Source code, examples, documentation, and installation instructions are available at <https://github.com/agrosavia-bioinformatics/multiGWAS>.

4.1 Input parameters

MutiGWAS uses as the only input a simple configuration text file with the values for the main parameters that drive the analysis. To create the configuration text file, users can choose either a text editor or the MultiGWAS GUI application. If users prefer a text file, it must have the parameter names and values separated by a colon, filenames enclosed in quotation marks, and TRUE or FALSE values to indicate if filters are applied. If the users prefer the GUI applications, they can create the configuration file using the input parameter view. In any case, this file must have the structure shown in the Figure(Figure ??.a.)

```

default:
    ploidy      : 4
    genotypeFile : "example-genotype-tetra.csv"
    phenotypeFile: "example-phenotype.csv"
    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 ploidy level of the organism (2: for diploids, 4: for tetraploids). The input genotype/phenotype filenames. The genome-wide significance threshold. The method for multiple testing correction. The GWAS model. The number of associations to report. The quality control filters choosing TRUE or FALSE. The filters are minor allele frequency, individual missing rate, SNP missing rate, and Hardy-Weinberg threshold. Finally, the GWAS packages selected for the analysis.

425

426 4.2 Using the command line interface

427 The execution of the CLI tool is simple. It only needs to open a Linux console,
 428 change to the folder where the configuration file was created, and type the name of
 429 the executable tool followed by the the filename of the configuration file, like this:

430 multiGWAS Test01.config

431 Then, the tool starts the execution, showing information on the process in the
 432 console window. When it finishes, the results are in a new subfolder called “out-
 433 Test01. The results include a complete HTML report containing the different views
 434 described in the results section, the source graphics and tables supporting the re-
 435 port, and the preprocessed tables from the results generated by the four GWAS
 436 packages used by MultiGWAS.

437 4.3 Using the graphical user interface

438 The interface allows users to save, load or specify the different input parameters
 439 for MultiGWAS in a friendly way (Fig. ??). The input parameters correspond to
 440 the settings included in the configuration file described in the subsection ??.
 441 The interface can be executed by calling the following command from a linux console:

442 jmultiGWAS

Arreglado, no hay folder de salida en este archivo. El orden si coincide

Yo no veo en la figura 9 el nombre del folder de salida que dice en la leyenda que aparece en el archivo de configuración. TambiÃ©n aquÃ hay que revisar el orden de cada descripción, que coincide con la figura. Yo no puede revisarlo porque si agrando la imagen, se me pierde

CambiÃ© “origina graphics” por “source graphics”

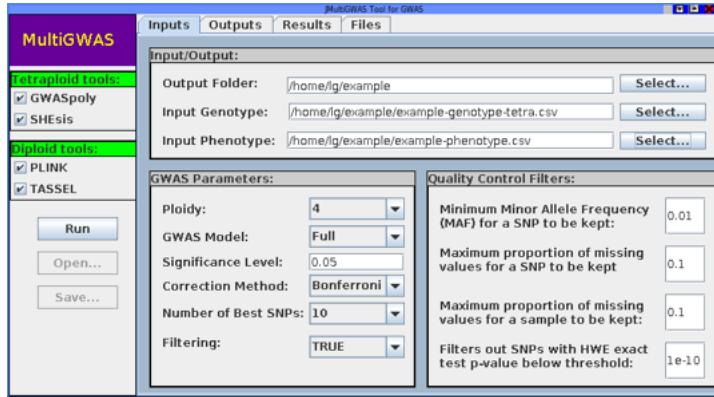


Figure 10: Main view of the MultiGWAS graphical user interface. The interface shows a main view on the center, a toolbar on the left and four tabs on top. From the main view, users can specify the input parameters for the analysis. From the toolbar, users can select the GWAS packages to be used in the analysis (two for tetraploids and two for diploids), and start the analysis with the current parameters (or load a previously saved configuration). And, from the tabs, in addition to specifying the input parameters, users can view the outputs of the process, the results of the analysis as an html report, and browse the source files that support the report.

443 5 Discussion

444 XXXXXXXXXXXXXXXXXXXXXXXXX

445 Challenges studying polyploid organisms are related to the complexity of (1) the
 446 data, and (2) inheritance mechanisms that are under study (?). The difficulties
 447 regarding the data complexity are the uncertainty in the allele dosage and null alle-
 448 les, but these problems are opportunities to improve software at the variant calling
 449 stage. Moreover, there is an on-going understanding of the inheritance mechanisms
 450 for polyploids. For autopolyploids, most loci have a polysomic inheritance. How-
 451 ever, sections of the genome that did not duplicate lead to disomic inheritance for
 452 some loci ((?; ?; ?)).

453 MultiGWAS does not face the problem of the allele dosage uncertainty. It is nec-
 454 essary to measure the impact of the allele dosage uncertainty at the GWAS stage to
 455 understand the effects arisen from this problem at the association stage. However,
 456 MultiGWAS addresses the second challenge variation of inheritance mechanisms by
 457 using different for existing software: two designed for polysomic inheritance ((?;
 458 ?)) together with two for disomic inheritance ((?; ?)). Thus it is a useful tool for
 459 researchers because it looks for significative associations that involve both types of
 460 inheritance.

461 Moreover, GWASpoly ((?)) offers models for different types of polyploid gene
 462 action additive, diploidized additive, duplex dominant, simplex dominant, and gen-
 463 eral. On the other hand, TASSEL ((?)) also models different types of gene action
 464 for diploids general, additive and dominant. We propose an automatic selection of

Luis esta parte revisala por
fa.

465 the gene action model for both tools based on a balance between the factor of inflation
466 and the replicability of the identified SNPs. We inform the user of the selected
467 model based on the automatic strategy; we consider this information helps to under-
468 stand the gene action model for the trait of interest. Even though the main focus
469 is on the resultant SNPs, the model has assumptions that reflect the gene actions
470 for a specific phenotype.

471 Replicability is a strategy to assess results derived from different methods. Multi-
472 GWAS integrates results to check for replicability in the results from different soft-
473 ware. By combining results, we can compare the outputs and look for coincidences.
474 Some software are more sensible, and others are more specific, but the integration
475 allows MultiGWAS to balance specificity and sensitivity.

476 We provide the user with different graphic outputs to help to interpret the re-
477 sults. We report the SNP profile for the SNPs identified by more than one software.
478 The SNP profile gives the researcher visual feedback from the SNP. We previously
479 check for significative SNPs based on the p-value; however, it is essential to go back
480 to the data and check if the SNP is a real association between the genotype and
481 phenotype. For this purpose, we designed the SNP profile.

482 6 Acknowledgements

483 This research was possible thanks to AGROSAVIA project *Investigación en conser-
484 vación, caracterización y uso de los recursos genéticos vegetales*. COLCIENCIAS, to-
485 day Minister of Science, Technology and Innovation of the republic of Colombia,
486 for supporting the postdoctoral researcher L. Garreta at AGROSAVIA during 2019-
487 2020 under the supervision of ICS and PHRH, (Grant number 811-2019). Colom-
488 bian Corporation for Agricultural Research's editorial fund thanks for financing this
489 publication. We thank Andres J. Cortes for helpful discussion.

Este agradecimiento es mÁis actualizado

Luis revisa si tenemos que agrecer a colciencias en otro modo