

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

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

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

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

11 July 15, 2020

12 **Abstract**

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

Comentatios Ivania

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

Comentatios Luis

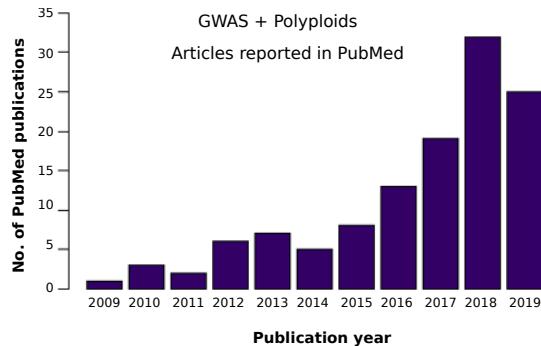
36 vs. SNP by each tool. **Contact:** phreyes@agrosavia.co

37

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

39 1 Introduction

40 The [Genome-wide association studies \(GWAS\)](#) are used to identify which variants
41 through the whole genome of a large number of individuals are associated with a
42 specific trait (Begum et al., 2012; Cantor et al., 2010). This methodology started
43 with humans and several model plants, such as rice, maize, and *Arabidopsis* (Cao
44 et al., 2011; Han and Huang, 2013; Korte and Farlow, 2013; Lauc et al., 2010; Tian
45 et al., 2011). Because of the advances in the next-gen sequencing technology and
46 the decline of the sequencing cost in recent years, there is an increase in the avail-
47 ability of genome sequences of different organisms at a faster rate (Ekblom and
48 Galindo, 2011; Ellegren, 2014). Thus, the GWAS is becoming the standard tool to
49 understand the genetic bases of either ecological or economic phenotypic variation
50 for both model and non-model organisms. This increment in GWAS includes com-
51 plex species such as polyploids (Fig. 1) (Ekblom and Galindo, 2011; Santure and
52 Garant, 2018).



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

55 The GWAS for polyploid species has fourth related challenges. First, as all
56 GWAS, we should replicate the study as a reliable method to validate the results
57 and recognize real associations. This replication involves finding the same associa-
58 tions either in several replicates from the study population using the same software
59 or testing different GWAS tools among the same study population. This approach
60 involved the use of different parameters, models, or conditions, to test how con-
61 sistent the results are (De et al., 2014; Pearson and Manolio, 2008). However, the
62 performance of different GWAS software could affect the results. For example, the
threshold *pvalue* for SNP significance change through four GWAS software (i.e.,
PLINK, TASSEL, GAPIT, and FaST-LMM) when sample size varies (Yan et al., 2019).

Me confunde al hablar de
software y tools

"The latter approach..."

Software or tools ??

63 It means that well-ranked SNPs from one package can be ranked differently in another.

64 Second, although there are many GWAS software available to repeat the analysis under different conditions (Gumpinger et al., 2018), most of them are designed exclusively for the diploid data matrix (Bourke et al., 2018). Therefore, it is often necessary to "diploidizing" the polyploid genomic data in order to replicate the analysis.

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

68 The iPAT allows running in a graphic interface three well-known command-line 69 GWAS software such as GAPIT, PLINK, and FarmCPU (Zhang et al., 2018) . However, the output from each package is separated. On the other hand, the easyGWAS 70 allows running a GWAS analysis on the web using different algorithms. This analysis 71 could run independently of both the computer capacity and operating system. 72 However, it needs either several datasets available or a dataset with a large number 73 of individuals to make replicates in order to compare among algorithms. Moreover, 74 the output from different algorithms is separated (Grimm et al., 2017). Thus, for 75 both software iPAT and easyGWAS, the integrative and comparative outputs among 76 software or algorithms are missing.

77 Fourth, the GWAS on polyploids generates a new level of complexity to understand 78 how allele dosage affect the phenotype expression on quantitative traits. Therefore, any tool that compares among software but also models with different 79 allele dosage will contribute to gain a better understanding in how redundancy or 80 complex interaction among alleles affect the phenotype expression and the evolution 81 of new phenotypes among polyploid species.

82 To contribute to sort out all the above fourth challenges, we developed the Multi- 83 GWAS tool that performs GWAS analyses for tetraploid species using four software 84 in parallel. Our tool include GWASpoly (Rosyara et al., 2016) and the SHEsis tool 85 (Shen et al., 2016) that accept polyploid genomic data, and PLINK (Purcell et al., 86 2007) and TASSEL (Bradbury et al., 2007) with the use of a "diploidized" genomic 87 matrix. The tool deals with [input file formats](#), [data preprocessing](#), [search for associations by running](#) four GWAS tools in parallel, and [creation of comparative reports](#) 88 from the output of each software to help the user to decide more intuitively the true 89 or false associations.

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

99 2 Method

100 The MultiGWAS tool has three main consecutive steps: the adjustment, the multi 101 analysis, and the integration (Fig. 2). In the adjustment step, MultiGWAS processes 102 the configuration file. Then it cleans and filters the genotype and phenotype, and 103 MultiGWAS "diploidize" the genomic data. Next, during the multi analysis, each 104 GWAS tool runs in parallel. Subsequently, in the integration step, the MultiGWAS 105 tool scans the output files from the four packages (i.e., GWASPoly, SHEsis, PLink,

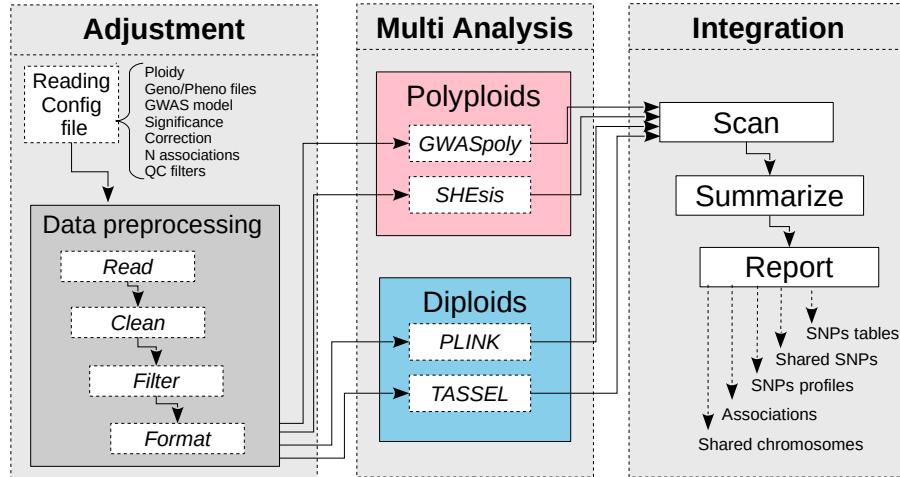


Figure 2: MultiGWAS flowchart has three steps: adjustment, multi analysis, and integration. The first step deals with input data management, reading the configuration file, **reading and preprocessing the input genomic data** (genotype and phenotype). The second step deals with GWAS analysis, configuring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

and TASSEL). Finally, it generates a summary of all results that contains score tables, Venn diagrams, SNP profiles, and Manhattan plots.

2.1 Adjustment stage

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

2.1.1 Reading configuration file

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

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

Genotype and phenotype input files: MultiGWAS uses two input files, one for genotype and one for phenotype. Genotypes files can be either in GWASpoly format (Rosyara et al., 2016) using SNP markers in rows and samples in columns (Fig. 3.a) or Variant Call Format (VCF) (Fig.3.b) which is transformed into GWASpoly format

122 using NGSEP 4.0.2 (Tello et al., 2019). The phenotype file contains only one trait
 123 and uses a matrix format with the first column for the sample names and the second
 124 column for the trait values (Fig. 3.c).

a. GWASpoly genotype file:

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	...

b. VCF genotype file:

```

##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
  
```

c. Phenotype matrix file:

Individual	Trait
sample01	3.59
sample02	4.07
sample03	1.05

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 (Tello et al., 2019): 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.).

125 **GWAS model:** MultiGWAS is designed to work with quantitative phenotypes and
 126 can run GWAS analysis using two types of statistical models that we have called
 127 *full* and *naive* models. The *full model* is known in the literature as the Q+K model Yu
 128 et al., 2006 and includes control for structure (Q) and relatedness between samples
 129 (K), whereas the *naive model* does not include any type of correction. Both models
 130 are based on linear regression approaches and variations of them are implemented
 131 by the four GWAS packages used by MultiGWAS. The *naive* is modeled with Generalized
 132 Linear Models (GLMs, Phenotype + Genotype), and the *full* is modeled with
 133 Mixed Linear Models (MLMs, Phenotype + Genotype + Structure + Kinship). The
 134 default model used by MultiGWAS is the *full model* (Q+K) Yu et al., 2006, which is
 135 expressed with the following equation:

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

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

142 **Genome-wide significance:** GWAS searches SNPs associated with the phenotype
 143 in a statistically significant manner. A threshold or significance level α is specified
 144 and compared with the *p-value* derived for each association score. Standard sig-
 145 nificance levels are 0.01 or 0.05 Gumpinger et al., 2018; Rosyara et al., 2016, and

146 MultiGWAS uses an α of 0.05 for the four GWAS packages. But the threshold is
147 adjusted according to each package, as some packages as GWASpoly and TASSEL
148 calculates the SNP effect for each genotypic class using different gene action models
149 (see “Multi analysis stage”). So, the number of tested markers may be different in
150 each model (see below) that results in different *p-value* thresholds.

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

160 **Number of reported associations:** Criticism has arisen in considering only sta-
161 tistically significant associations as the only possible correct associations Kaler and
162 Purcell, 2019; Thompson et al., 2011. Many of low *p-value* associations, closer to
163 being significant, are discarded due to the stringent significance levels, and conse-
164 quently increasing the number of false negatives. To help to analyze both signif-
165 icant and non-significant associations, MultiGWAS provides the option to specify
166 the number of best-ranked associations (lower *p-values*), adding the corresponding
167 *p-value* to each association found. In this way, it is possible to enlarge the number
168 of results, and we can observe replicability in the results for different programs.
169 Nevertheless, we present each association with the corresponding *p-value*.

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

- 176 • **MAF of x:** filters out SNPs with minor allele frequency below x (default 0.01);
- 177 • **MIND of x:** filters out all individuals with missing genotypes exceeding $x*100\%$
178 (default 0.1);
- 179 • **GENO of x:** filters out SNPs with missing values exceeding $x*100\%$ (default
180 0.1);
- 181 • **HWE of x:** filters out SNPs which have Hardy-Weinberg equilibrium exact test
182 *p-value* below the x threshold.

183 MultiGWAS does the MAF filtering, and uses the PLINK package Gumpinger et al.,
184 2018 for the other three filters: MIND, GENO, and HWE.

185 **GWAS tools:** List of names of GWAS packages to run and integrate into MultiG-
186 WAS analysis. Currently four packages, two for tetraploid organisms: GWASpoly
187 and SHEsis, and two for diploids: PLINK and TASSEL.

188 **2.1.2 Data preprocessing**

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

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

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

201 **2.2 Multi analysis stage**

202 MultiGWAS runs in parallel using two types of statistical models specified in the
203 parameters file, the Full model (Q+K) and Naive (i.e., without any control) where
204 Q refers to population structure and K refers to relatedness, calculated by kinship
205 coefficients across individuals (Sharma et al., 2018). The Full model (Q+K) controls
206 for both population structure and individual relatedness. For population structure,
207 MultiGWAS uses the Principal Component Analysis (PCA) and takes the top five PC
208 as covariates. For relatedness, MultiGWAS uses kinship matrices that TASSEL and
209 GWASpoly calculated separately, and for PLINK and SHEsis, relatedness depends on
210 kinship coefficients calculated with the PLINK 2.0 built-in algorithm (Chang et al.,
211 2015).

212 **2.2.1 GWASpoly**

213 GWASpoly (Rosyara et al., 2016) is an R package designed for GWAS in polyploid
214 species used in several studies in plants (Berdugo-Cely et al., 2017; Ferrão et al.,
215 2018; Sharma et al., 2018; Yuan et al., 2019). GWASpoly uses a Q+K linear mixed
216 model with biallelic SNPs that account for population structure and relatedness.
217 Also, to calculate the SNP effect for each genotypic class, GWASpoly provides eight
218 gene action models: general, additive, simplex dominant alternative, simplex dom-
219 inant reference, duplex dominant alternative, duplex dominant, diplo-general, and
220 diplo-additive. As a consequence, the number of statistical test performed can be
221 different in each action model and so thresholds below which the *p*-values are con-
222 sidered significant.

223 MultiGWAS is using GWASpoly version 1.3 with all gene action models available
224 to find associations. The MultiGWAS reports the top *N* best-ranked (the SNPs with

225 lowest *p*-values) that the user specified in the *N* input configuration file. The *full*
226 model used by GWASpoly includes the population structure and relatedness, which
227 are estimated using the first five principal components and the kinship matrix, re-
228 spectively, both calculated with the GWASpoly built-in algorithms.

229 **2.2.2 SHEsis**

230 SHEsis is a program based on a linear regression model that includes single-locus
231 association analysis, among others. The software design includes polyploid species.
232 However, their use is mainly in diploids animals and humans (Meng et al., 2019;
233 Qiao et al., 2015).

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

238 **2.2.3 PLINK**

239 PLINK is one of the most extensively used programs for GWAS in humans and any
240 diploid species (Power et al., 2016). PLINK includes a range of analyses, including
241 univariate GWAS using two-sample tests and linear regression models.

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

248 **2.2.4 TASSEL**

249 TASSEL is another standard GWAS program based on the Java software developed
250 initially for maize but currently used in several species (Álvarez et al., 2017; Zhang
251 et al., 2018). For the association analysis, TASSEL includes the general linear model
252 (GLM) and mixed linear model (MLM) that accounts for population structure and
253 relatedness. Moreover, as GWASPoly, TASSEL provides three-gene action models to
254 calculate the SNP effect of each genotypic class: general, additive, and dominant,
255 and so the significance threshold depends on each action model.

256 MultiGWAS is using TASSEL 5.0, with all gene action models used to find the
257 *N* best-ranked associations and reporting the top *N* best-ranked associations (SNPs
258 with lowest *p*-values). Naive GWAS uses the GLM, and full GWAS uses the MLM
259 with two parameters: population structure that uses the first five principal compo-
260 nents, and relatedness that uses the kinship matrix with centered IBS method, both
261 calculated with the TASSEL built-in algorithms.

262 **2.3 Integration stage.**

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

266 **2.3.1 Calculation of *p-values* and significance thresholds**

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

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

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

290 **2.3.2 Selection of significant and best-ranked associations**

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

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

300 The selection of these groups takes into account whether the GWAS package
301 uses only one gene action model, as PLINK and SHEsis do, or uses several ones,
302 as GWASpoly and TASSEL do. In the first case, there is only one resulting set of

303 associations and the selection is straightforward, as described above. However, in
 304 the second case, there are several resulting sets of associations, one for each model,
 305 and MultiGWAS selects both groups of associations by choosing the gene action
 306 model with the highest number of shared SNPs and with the inflation factor closest
 307 to one, according to the following equation:

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

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

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

318 2.3.3 Integration of results

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

- 322 • A Quantile-Quantile (QQ) plots for the resultant p -values of each tool and
 323 the corresponding inflation factor λ to assess the degree of the test statistic
 324 inflation.
- 325 • A Manhattan plot of each tool with two lower thresholds, one for the best-
 326 ranked SNPs, and another for the significant SNPs.

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

332 For each SNP identified more than once, we provide what we call the SNP pro-
 333 file. That is a heat diagram for a specific SNP, where each column is a genotype
 334 state AAAA, AAAB, AABB, ABAA, and BBBB. Moreover, each row corresponds to a
 335 sample. Samples with close genotypes form together clusters. Thus to generate
 336 the clusters, we do not use the phenotype information. However, we present the
 337 phenotype information in the figure as the color. This figure visually provides in-
 338 formation regarding genotype and phenotype information simultaneously for the
 339 whole population. We present colors as tones between white and black for color
 340 blind people.

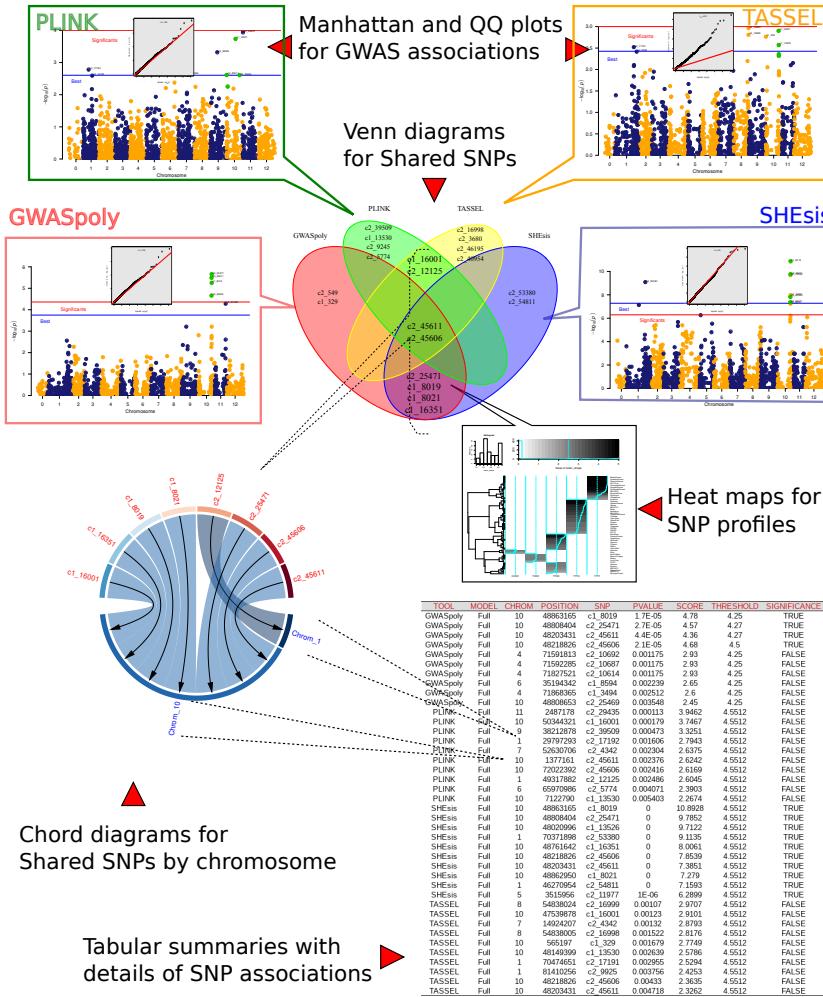


Figure 4: Reports presented by MultiGWAS. For each tool, first a QQ plot that assesses the resultant p-values. Second, a Manhattan plot for each tool with two lines, blue and red, respectively, is the lower limit for the best ranked and significative SNPs. We present two Venn diagrams, one for the significative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

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

344 MultiGWAS generates a report, one document with the content previously de-
345 scribed. Besides, there is a folder with the individual figures just in case the user
346 needs one. In the supplementary information, we include a report and a description
348 of the report content (Supplementary Material 1)

349 In the following section, we present the results of the functionality of the tool
350 applied on a open dataset of a diversity panel of a tetraploid potato, genotyped and
351 phenotyped as part of the USDA-NIFA Solanaceae Coordinated Agricultural Project
352 (SolCAP) (Hirsch et al., 2013).

Aquí deberíamos decir es-
pecificamente qué datos son
y de donde salieron

353 3 Results

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

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

- 364 • Manhattan and Q-Q plots to show GWAS associations.
- 365 • Venn diagrams to show associations identified by single or several tools.
- 366 • Heat diagrams to show the genotypic structure of shared SNPs.
- 367 • Chord diagrams to show shared SNPs by chromosomes.
- 368 • Score tables to show detailed information of associations for both summary
369 results from MultiGWAS and particular results from each GWAS package

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

374 3.1 Manhattan and QQ plots for GWAS associations

375 MultiGWAS uses classical Manhattan and Quantile-Quantile plots (QQ plots) to
376 visualize the results of each package. In both plots, the points are the SNPs and
377 their *p-values* are transformed into scores like $-\log_{10}(p\text{-values})$ (see Fig. 5). The
378 Manhattan plot shows the strength of association of the SNPs (y-axis) distributed at

379 their genomic location (x-axis), so the higher the score, the stronger the association.
 380 While the QQ plot compares the expected distribution of p -values (y-axis) with the
 381 observed distribution (x-axis)..

382 MultiGWAS adds distinctive marks to both plots to identify different types of
 383 SNPs: (a) In the Manhattan plots, the significant SNPs are above a red line and the
 384 best-ranked SNPs are above a blue line. In addition, **SNPs shared between packages**
 385 **are colored green** (See Fig. 6.b). (b) In the QQ plots, a red diagonal line indicates
 386 the expected distribution under the null hypothesis of no association of SNPs with
 387 the phenotype, both distributions should coincide, and most SNPs should lie on the
 388 diagonal line. **Deviations for a large number of SNPs may reflect inflated p -values**
 389 **due to population structure or cryptic relatedness. But, it is also expected that few**
 390 **SNPs deviate from the diagonal for a truly polygenic trait (Power et al., 2016).**

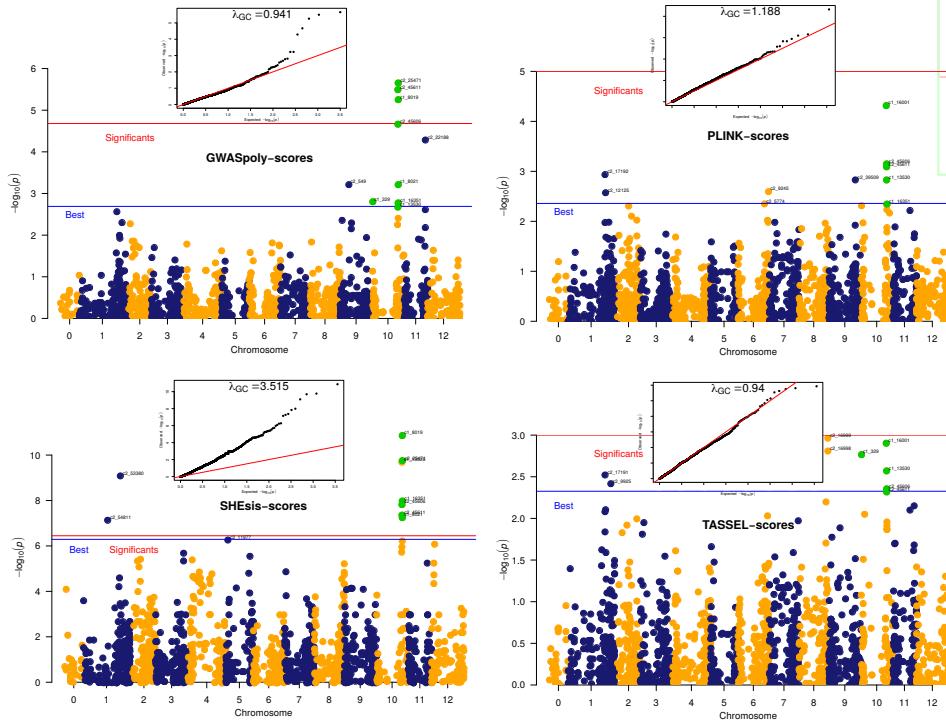


Figure 5: Associations in the tetraploid potato dataset. MultiGWAS shows the associations identified by the four GWAS packages using Manhattan and QQ plots. In the case of the tetraploid potato, several SNPs are observed to be shared between the four packages (green dots). The best-ranked SNPs are above 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.

Corregido: "among packages" y cita de polygenic traits

Revisar si esto del color verde es correcto "shared among all packages". También revisar si está bien dicho lo de la comparación entre pocos genes y muchos genes. No se si sea necesario cuantificar la desviación. Originalmente se cuantificaba pero era confuso

Ya está

La leyenda debe incluir también la descripción de qué significa las líneas rojas, azules y los puntos verdes

391 3.2 Tables and Venn diagrams for single and shared SNPs

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

396 First, best-ranked SNPs correspond to the top-scored N SNPs, whether they were
397 assessed significant or not by its package, and with N defined by the user in the
398 configuration file. These SNPs appears in both a SNPs table (Figure 6.a), and in a
399 Venn diagram (Figure 6.b). The table lists them by package and sorts by decreasing
400 score, whereas the Venn diagram emphasizes if they were best-ranked either in a
401 single package or in several at once (shared).

402 Second, the significant SNPs correspond to the ones valued statistically signifi-
403 cant by each package. They appear in a Venn diagram (Figure 6.c), and in the SNPs
404 table, marked with significance TRUE (T) in the table of the Figure 6.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

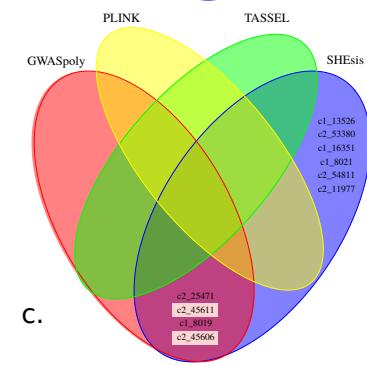
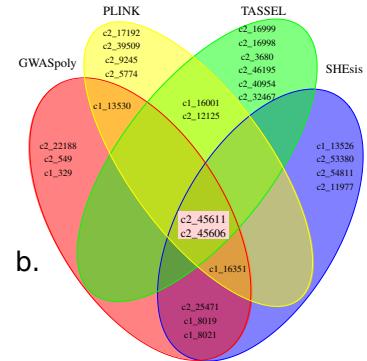


Figure 6: Shared SNPs Views. Tabular and graphical views of SNP associations identified by one or more GWAS packages (shared SNPs). SNPs identified by all packages are marker with red background in all figures. **(a)** Table with details of the N=9 best-ranked SNPs from each GWAS package. Each row corresponds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

3.3 Heat diagrams for the structure of shared SNPs

MultiGWAS creates a two-dimensional representation, called the SNP profile, to visualize each trait by individuals and genotypes as rows and columns, respectively (Figure 7). At the left, the individuals are grouped in a dendrogram by their genotype. At the right, there is the name or ID of each individual. At the bottom, the genotypes are ordered from left to right, starting from the major to the minor allele (i.e., AAAA, AAAB, AABB, ABAA, BBBB). At the top, there is a description of the trait based on a histogram of frequency (top left) and by an assigned color for each numerical phenotype value using a grayscale (top right). Thus, each individual appears as a colored line by its phenotype value on its genotype column. For each

415 column, there is a solid cyan line with the mean of each column and a broken cyan
416 line that indicates how far the cell deviates from the mean.

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

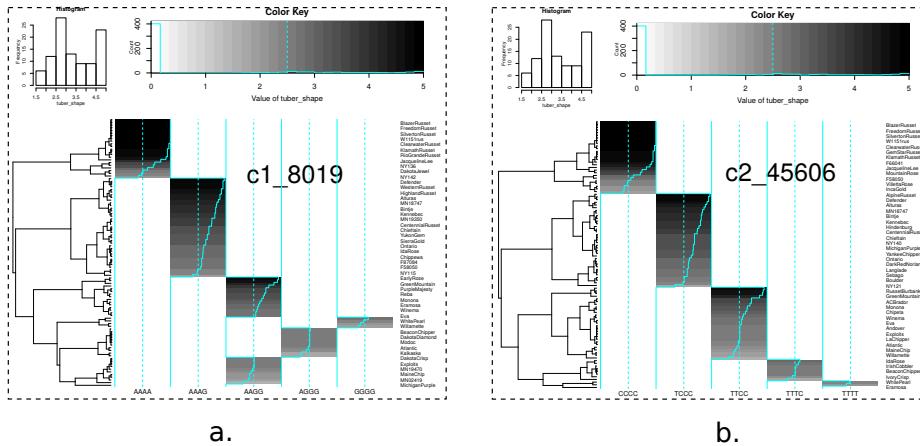


Figure 7: SNP profiles. SNP profiles for two of the best-ranked significant SNPs shown in figure 6.b. (a) SNP c2_45606 best-ranked by the four packages (central intersection of the Venn diagram Figure 6.b) (b) SNP c1_8019 best-ranked by the two tetraploid packages (Figure 6.b), and also identified as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

419

420 3.4 Chord diagrams for SNPs by chromosome

421 The chord diagrams visualize the location across chromosomes of the best-ranked
422 associated SNPs shared among the four packages and described in the table 6.a.
423 Thus, in the case of the tetraploid potato, we found that they are located mostly in
424 chromosome 10 (Figure 8.a). This visualization complements the manhattan plots
425 from each GWAS package (Figure 8.b).

426 4 Availability and Implementation

427 The core of the MultiGWAS tool runs under R and users can interact with the tool by
428 either a command-line interface (CLI) developed in R or a graphical user interface
429 (GUI) developed in Java (Figure 10). Source code, examples, documentation, and
430 installation instructions are available at <https://github.com/agrosavia-bioinformatics/multiGWAS>.

432 4.1 Input parameters

433 MutiGWAS uses as the only input a simple configuration text file with the values
434 for the main parameters that drive the analysis. To create the configuration text

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. Generalmente, en una típica análisis GWAS las asociaciones más fuertes se señalan por varios SNPs cercanos correlacionados ubicados en el mismo cromosoma, como en los plots de Manhattan, donde estas asociaciones forman picos con varios SNPs mostrando la misma señal. Contrairement, no se muestran picos cuando pocos SNPs correlacionados con un trait. Leyenda de la figura 8. Las más asociaciones identificadas en un cromosoma, más amplio es el espacio de su sector.

localización a lo largo del genoma en vez de a lo largo del cromosoma?

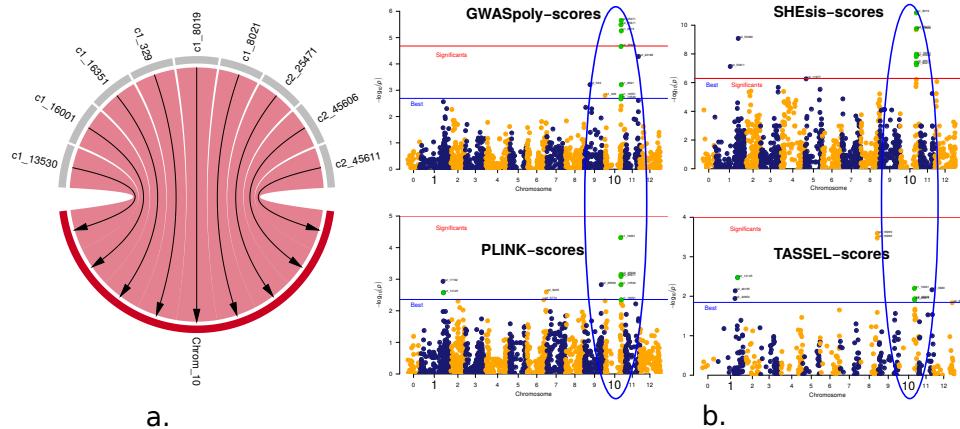


Figure 8: SNPs by chromosome. The position of best-ranked SNPs across chromosomes using two different visualizations. (a) Chord diagram showing that best-ranked SNPs located in chromosome 10. The SNPs are at the top and the chromosomes at the bottom. The arrows connect the best-ranked SNPs with their 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.

435 file, users can choose either a text editor or the MultiGWAS GUI application. If
 436 users prefer a text file, it must have the parameter names and values separated by a
 437 colon, filenames enclosed in quotation marks, and TRUE or FALSE values to indicate
 438 if filters are applied. If the users prefer the GUI applications, they can create the
 439 configuration file using the input parameter view. In any case, this file must have
 440 the structure shown in the Figure(Figure 9.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.

442

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. También aquí hay que revisar el orden de cada descripción, que coincide con la figura. Yo no puedo re-

443 4.2 Using the command line interface

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

447 `multiGWAS Test01.config`

448 Then, the tool starts the execution, showing information on the process in the
449 console window. When it finishes, the results are in a new subfolder called “*out-*
450 *Test01*. [The results include a complete HTML report containing the different views](#)
451 [described in the results section, the source graphics and tables supporting the](#)
452 [report, and the preprocessed tables from the results generated by the four GWAS](#)
453 [packages used by MultiGWAS.](#)

Cambié “origina graphics”
por “source graphics”

454 4.3 Using the graphical user interface

455 The interface allows users to save, load or specify the different input parameters
456 for MultiGWAS in a friendly way (Fig. 10). The input parameters correspond to
457 the settings included in the configuration file described in the subsection 2.1.1. The
458 interface can be executed by calling the following command from a linux console:

459 `jmultiGWAS`

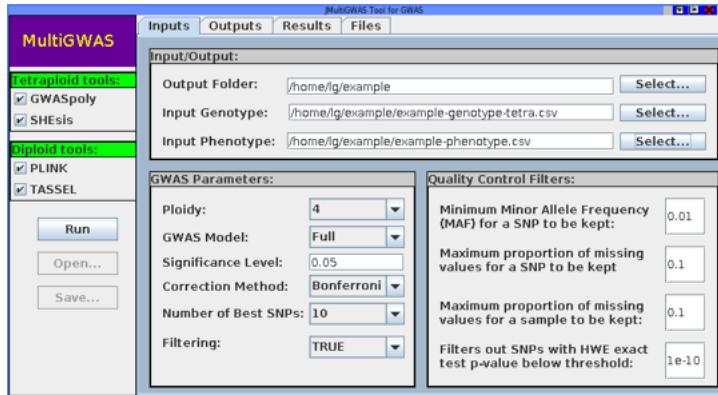


Figure 10: Main view of the MultiGWAS graphical user interface. The interface shows a main view on the center, a toolbar on the left and four tabs on top. From the main view, users can specify the input parameters for the analysis. From the toolbar, users can select the GWAS packages to be used in the 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.

460 **5 Discussion**

461 XXXXXXXXXXXXXXXXXXXXXXXXX

462 Challenges studying polyploid organisms are related to the complexity of (1)
463 the data, and (2) inheritance mechanisms that are under study (Dufresne et al.,
464 2014). The difficulties regarding the data complexity are the uncertainty in the
465 allele dosage and null alleles, but these problems are opportunities to improve soft-
466 ware at the variant calling stage. Moreover, there is an on-going understanding
467 of the inheritance mechanisms for polyploids. For autopolyploids, most loci have
468 a polysomic inheritance. However, sections of the genome that did not duplicate
469 lead to disomic inheritance for some loci (Dufresne et al., 2014; Lynch and Conery,
470 2000; Ohno, 1970).

471 MultiGWAS does not face the problem of the allele dosage uncertainty. It is nec-
472 essary to measure the impact of the allele dosage uncertainty at the GWAS stage to
473 understand the effects arisen from this problem at the association stage. However,
474 MultiGWAS addresses the second challenge variation of inheritance mechanisms
475 by using different for existing software: two designed for polysomic inheritance
476 (Rosyara et al., 2016; Shen et al., 2016) together with two for disomic inheritance
477 (Bradbury et al., 2007; Purcell et al., 2007). Thus it is a useful tool for researchers
478 because it looks for significative associations that involve both types of inheritance.

479 Moreover, GWASpoly (Rosyara et al., 2016) offers models for different types
480 of polyploid gene action additive, diploidized additive, duplex dominant, simplex
481 dominant, and general. On the other hand, TASSEL (Bradbury et al., 2007) also
482 models different types of gene action for diploids general, additive and dominant.
483 We propose an automatic selection of the gene action model for both tools based
484 on a balance between the factor of inflation and the replicability of the identified
485 SNPs. We inform the user of the selected model based on the automatic strategy;
486 we consider this information helps to understand the gene action model for the
487 trait of interest. Even though the main focus is on the resultant SNPs, the model
488 has assumptions that reflect the gene actions for a specific phenotype.

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

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

Luis esta parte revisala por
fa.

501 6 Acknowledgements

502 This research was possible thanks to AGROSAVIA project *Investigación en conser-*
503 *vación, caracterización y uso de los recursos genéticos vegetales.* COLCIENCIAS,
504 today Minister of Science, Technology and Innovation of the republic of Colombia,
505 for supporting the postdoctoral researcher L. Garreta at AGROSAVIA during 2019-
506 2020 under the supervision of ICS and PHRH, (Grant number 811-2019). Colom-
507 bian Corporation for Agricultural Research's editorial fund thanks for financing this
508 publication. We thank Andres J. Cortes for helpful discussion.

Este agradecimiento es más actualizado

Luis revisa si tenemos que agreecer a colciencias en otro modo

509 References

- 510 Álvarez, M. F., Angarita, M., Delgado, M. C., García, C., Jiménez-Gomez, J., Geb-
511 hardt, C., & Mosquera, T. (2017). Identification of Novel Associations of
512 Candidate Genes with Resistance to Late Blight in Solanum tuberosum
513 Group Phureja. *Frontiers in Plant Science*, 8, 1040. <https://doi.org/10.3389/fpls.2017.01040>
- 514 Begum, F., Ghosh, D., Tseng, G. C., & Feingold, E. (2012). Comprehensive literature
515 review and statistical considerations for gwas meta-analysis. *Nucleic acids*
516 *research*, 40(9), 3777–3784.
- 517 Berdugo-Cely, J., Valbuena, R. I., Sánchez-Betancourt, E., Barrero, L. S., & Yock-
518 teng, R. (2017). Genetic diversity and association mapping in the colom-
519 bian central collection of solanum tuberosum L. Andigenum group using
520 SNPs markers. *PLoS ONE*, 12(3). <https://doi.org/10.1371/journal.pone.0173039>
- 521 Bourke, P. M., Voorrips, R. E., Visser, R. G. F., & Maliepaard, C. (2018). Tools for Ge-
522 netic Studies in Experimental Populations of Polyploids. *Frontiers in Plant*
523 *Science*, 9, 513. <https://doi.org/10.3389/fpls.2018.00513>
- 524 Bradbury, P. J., Zhang, Z., Kroon, D. E., Casstevens, T. M., Ramdoss, Y., & Buckler,
525 E. S. (2007). TASSEL: software for association mapping of complex traits
526 in diverse samples. *Bioinformatics*, 23(19), 2633–2635. <https://doi.org/10.1093/bioinformatics/btm308>
- 527 Cantor, R. M., Lange, K., & Sinsheimer, J. S. (2010). Prioritizing gwas results: A
528 review of statistical methods and recommendations for their application.
529 *The American Journal of Human Genetics*, 86(1), 6–22.
- 530 Cao, J., Schneeberger, K., Ossowski, S., Günther, T., Bender, S., Fitz, J., Koenig, D.,
531 Lanz, C., Stegle, O., Lippert, C., Et al. (2011). Whole-genome sequencing
532 of multiple arabidopsis thaliana populations. *Nature genetics*, 43(10), 956.
- 533 Chang, C. C., Chow, C. C., Tellier, L. C., Vattikuti, S., Purcell, S. M., & Lee, J. J.
534 (2015). Second-generation PLINK: Rising to the challenge of larger and
535 richer datasets. *GigaScience*, 4(1), arXiv 1410.4803, 1–16. <https://doi.org/10.1186/s13742-015-0047-8>
- 536 De, R., Bush, W. S., & Moore, J. H. (2014). Bioinformatics Challenges in Genome-
537 Wide Association Studies (GWAS). In R. Trent (Ed.), *Clinical bioinformatics*
- 538
- 539
- 540
- 541

- 542 (pp. 63–81). New York, NY, Springer New York. https://doi.org/10.1007/978-1-4939-0847-9_5
- 543 Dufresne, F., Stift, M., Vergilino, R., & Mable, B. K. (2014). Recent progress and challenges in population genetics of polyploid organisms: An overview of current state-of-the-art molecular and statistical tools. *Molecular Ecology*, 23(1), <https://onlinelibrary.wiley.com/doi/pdf/10.1111/mec.12581>, 40–69. <https://doi.org/10.1111/mec.12581>
- 544 Ekblom, R., & Galindo, J. (2011). Applications of next generation sequencing in molecular ecology of non-model organisms. *Heredity*, 107(1), 1–15.
- 545 Ellegren, H. (2014). Genome sequencing and population genomics in non-model organisms. *Trends in ecology & evolution*, 29(1), 51–63.
- 546 Ferrão, L. F. V., Benevenuto, J., Oliveira, I. d. B., Cellon, C., Olmstead, J., Kirst, M., Resende, M. F. R., & Munoz, P. (2018). Insights Into the Genetic Basis of Blueberry Fruit-Related Traits Using Diploid and Polyploid Models in a GWAS Context. *Frontiers in Ecology and Evolution*, 6, 107. <https://doi.org/10.3389/fevo.2018.00107>
- 547 - Paper for layout. - Many concepts of GWAS, especially structure population.
- 548 Grimm, D. G., Roqueiro, D., Salomé, P. A., Kleeberger, S., Greshake, B., Zhu, W., Liu, C., Lippert, C., Stegle, O., Schölkopf, B., Weigel, D., & Borgwardt, K. M. (2017). easyGWAS: A Cloud-Based Platform for Comparing the Results of Genome-Wide Association Studies. *The Plant Cell*, 29(1), 5–19. <https://doi.org/10.1105/tpc.16.00551>
- 549 Gumpinger, A. C., Roqueiro, D., Grimm, D. G., & Borgwardt, K. M. (2018). *Methods and Tools in Genome-wide Association Studies* (Vol. 1819).
- 550 Han, B., & Huang, X. (2013). Sequencing-based genome-wide association study in rice. *Current opinion in plant biology*, 16(2), 133–138.
- 551 Hirsch, C. N., Hirsch, C. D., Felcher, K., Coombs, J., Zarka, D., Van Deynze, A., De Jong, W., Veilleux, R. E., Jansky, S., Bethke, P., Douches, D. S., & Buell, C. R. (2013). Retrospective view of North American potato (*Solanum tuberosum* L.) breeding in the 20th and 21st centuries. *G3: Genes, Genomes, Genetics*, 3(6), 1003–1013. <https://doi.org/10.1534/g3.113.005595>
- 552 Kaler, A. S., & Purcell, L. C. (2019). Estimation of a significance threshold for genome-wide association studies. *BMC Genomics*, 20(1), 618. <https://doi.org/10.1186/s12864-019-5992-7>
- 553 Korte, A., & Farlow, A. (2013). The advantages and limitations of trait analysis with gwas: A review. *Plant methods*, 9(1), 29.
- 554 Lauc, G., Essafi, A., Huffman, J. E., Hayward, C., Knežević, A., Kattla, J. J., Polašek, O., Gornik, O., Vitart, V., Abrahams, J. L., Et al. (2010). Genomics meets glycomics—the first gwas study of human n-glycome identifies hnf1 α as a master regulator of plasma protein fucosylation. *PLoS genetics*, 6(12).
- 555 Lynch, M., & Conery, J. S. (2000). The evolutionary fate and consequences of duplicate genes. *science*, 290(5494), 1151–1155.
- 556 Meng, J., Song, K., Li, C., Liu, S., Shi, R., Li, B., Wang, T., Li, A., Que, H., Li, L., & Zhang, G. (2019). Genome-wide association analysis of nutrient traits in

- 587 the oyster *Crassostrea gigas*: Genetic effect and interaction network. *BMC*
588 *Genomics*, 20(1), 1–14. <https://doi.org/10.1186/s12864-019-5971-z>
- 589 Ohno, S. (1970). *Evolution by gene duplication*. Springer Berlin Heidelberg. <https://doi.org/10.1007/978-3-642-86659-3>
- 590 Pearson, T. A., & Manolio, T. A. (2008). How to interpret a genome-wide association
591 study. *JAMA - Journal of the American Medical Association*, 299(11), 1335–
592 1344. <https://doi.org/10.1001/jama.299.11.1335>
- 593 Power, R. A., Parkhill, J., & De Oliveira, T. (2016). Microbial genome-wide associa-
594 tion studies: lessons from human GWAS. *Nature Reviews Genetics*, 18(1),
595 41–50. <https://doi.org/10.1038/nrg.2016.132>
- 596 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D., Maller,
597 J., Sklar, P., De Bakker, P. I., Daly, M. J., & Sham, P. C. (2007). PLINK: A tool
598 set for whole-genome association and population-based linkage analyses.
599 *American Journal of Human Genetics*, 81(3), 559–575. <https://doi.org/10.1086/519795>
- 600 Qiao, H. P., Zhang, C. Y., Yu, Z. L., Li, Q. M., Jiao, Y., & Cao, J. P. (2015). Genetic vari-
601 ants identified by GWAS was associated with colorectal cancer in the Han
602 Chinese population. *Journal of Cancer Research and Therapeutics*, 11(2),
603 468–470. <https://doi.org/10.4103/0973-1482.150346>
- 604 Rosyara, U. R., De Jong, W. S., Douches, D. S., & Endelman, J. B. (2016). Software
605 for Genome-Wide Association Studies in Autopolyploids and Its Applica-
606 tion to Potato. *The Plant Genome*, 9(2), 1–10. <https://doi.org/10.3835/plantgenome2015.08.0073>
- 607 Santure, A. W., & Garant, D. (2018). Wild gwas—association mapping in natural
608 populations. *Molecular ecology resources*, 18(4), 729–738.
- 609 Sharma, S. K., MacKenzie, K., McLean, K., Dale, F., Daniels, S., & Bryan, G. J.
610 (2018). Linkage disequilibrium and evaluation of genome-wide associa-
611 tion mapping models in tetraploid potato. *G3: Genes, Genomes, Genetics*,
612 8(10), 3185–3202. <https://doi.org/10.1534/g3.118.200377>
- 613 Shen, J., Li, Z., Chen, J., Song, Z., Zhou, Z., & Shi, Y. (2016). SHEsisPlus, a toolset
614 for genetic studies on polyploid species. *Scientific Reports*, 6, 1–10. <https://doi.org/10.1038/srep24095>
- 615 Tello, D., Gil, J., Loaiza, C. D., Riascos, J. J., Cardozo, N., & Duitama, J. (2019).
616 NGSEP3: accurate variant calling across species and sequencing protocols.
617 *Bioinformatics*, 35(22), 4716–4723. <https://doi.org/10.1093/bioinformatics/btz275>
- 618 Thompson, J. R., Attia, J., & Minelli, C. (2011). The meta-analysis of genome-wide
619 association studies. *Briefings in Bioinformatics*, 12(3), 259–269. <https://doi.org/10.1093/bib/bbr020>
- 620 Tian, F., Bradbury, P. J., Brown, P. J., Hung, H., Sun, Q., Flint-Garcia, S., Rocheford,
621 T. R., McMullen, M. D., Holland, J. B., & Buckler, E. S. (2011). Genome-
622 wide association study of leaf architecture in the maize nested association
623 mapping population. *Nature genetics*, 43(2), 159–162.
- 624 Yan, Y. Y., Burbridge, C., Shi, J., Liu, J., & Kusalik, A. (2019). Effects of input data
625 quantity on genome-wide association studies (GWAS). *International Jour-*

- 632 *nal of Data Mining and Bioinformatics*, 22(1), 19–43. <https://doi.org/10.1504/IJDMB.2019.099286>
- 633
- 634 Yu, J., Pressoir, G., Briggs, W. H., Vroh, I. B., Yamasaki, M., Doebley, J. F., McMullen,
635 M. D., Gaut, B. S., Nielsen, D. M., Holland, J. B., Et al. (2006). A unified
636 mixed-model method for association mapping that accounts for multiple
637 levels of relatedness. *Nature genetics*, 38(2), 203–208.
- 638 Yuan, J., Bizimungu, B., De Koeyer, D., Rosyara, U., Wen, Z., & Lagüe, M. (2019).
639 Genome-Wide Association Study of Resistance to Potato Common Scab.
640 *Potato Research*. <https://doi.org/10.1007/s11540-019-09437-w>
- 641 Zhang, S., Chen, X., Lu, C., Ye, J., Zou, M., Lu, K., Feng, S., Pei, J., Liu, C., Zhou, X.,
642 Ma, P., Li, Z., Liu, C., Liao, Q., Xia, Z., & Wang, W. (2018). Genome-wide
643 association studies of 11 agronomic traits in cassava (*Manihot esculenta*
644 crantz). *Frontiers in Plant Science*, 9(April), 1–15. <https://doi.org/10.3389/fpls.2018.00503>
- 645