

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 August 14, 2020

12 **Abstract**

13 **Summary:** The Genome-Wide Association Studies (GWAS) are essential to
14 determine the genetic bases of either ecological or economic phenotypic variation
15 across individuals within populations of wild and domesticated species. For
16 this research question, current practice is the replication of the GWAS testing
17 different parameters and models to validate the reproducibility of results. How-
18 ever, straightforward methodologies that manage both replication and tetra-
19 ploid data are still missing. To solve this problem, we designed the MultiGWAS,
20 a tool that does GWAS for diploid and tetraploid organisms by executing in par-
21 allel four software, two for polyploid data (GWASPoly and SHEsis) and two for
22 diploids data (PLINK and TASSEL). MultiGWAS has several advantages. It runs
23 either in the command line or in an interface. It manages different genotype
24 formats, including VCF. It executes both the full and naïve models using sev-
25 eral quality filters. Besides, it calculates a score to choose the best gene action
26 model across GWASPoly and TASSEL. Finally, it generates several reports that
27 facilitate the identification of false associations from both the significant and the
28 best-ranked association SNP among the four software. We tested MultiGWAS
29 with tetraploid potato data. The execution demonstrated that the Venn diagram
30 and the other companion reports (i.e., Manhattan and QQ plots, heatmaps for
31 associated SNP profiles, and chord diagrams to trace associated SNP by chromo-
32 somes) were useful to identify associated SNP shared among different models
33 and parameters. Therefore, we confirmed that MultiGWAS is a suitable wrap-
34 ping tool that successfully handles GWAS replication in both diploid and tetra-
35 ploid organisms.

36 **Contact:** phreyes@agrosavia.co

37 Keywords: GWAS on polyploids, GWASPoly, PLINK, SNP, SHEsis, software,
38 TASSEL

39

1 Introduction

40 The Genome-Wide Association Studies (GWAS) comprise statistical tests that iden-
41 tify which variants through the whole genome of a large number of individuals
42 are associated with a specific trait (Begum et al., 2012; Cantor et al., 2010). This
43 methodology started with humans and several model plants, such as rice, maize,
44 and *Arabidopsis* (Cao et al., 2011; Han and Huang, 2013; Korte and Farlow, 2013;
45 Lau et al., 2010; Tian et al., 2011). Because of the advances in the next-gen se-
46 quencing technology and the decline of the sequencing cost in recent years, there
47 is an increase in the availability of genome sequences of different organisms at a
48 faster rate (Ekblom and Galindo, 2011; Ellegren, 2014). Thus, the GWAS is be-
49 coming the standard tool to understand the genetic bases of either ecologically or
50 economically relevant phenotypic variation for both model and non-model organ-
51 isms. This increment includes complex species such as polyploids (Fig. 1) (Ekblom
52 and Galindo, 2011; Santure and Garant, 2018).

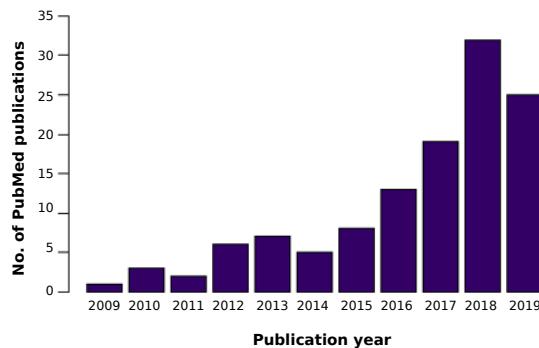


Figure 1: The number of peer-reviewed papers that contains the keywords "GWAS" and "polyploid" in the PubMed database between 2009 and 2019.

53 The GWAS for polyploid species has fourth related challenges. First, replica-
54 tion is critical to validate GWAS results and capture real associations. This approach
55 involved using different parameters, models, or conditions to test how consistent
56 the results are in the same software or different GWAS tools (De et al., 2014; Pear-
57 son and Manolio, 2008). However, the performance of different GWAS software
58 could affect the results. For example, the significance threshold for *pvalue* changes
59 through four GWAS software (i.e., PLINK, TASSEL, GAPIT, and FaST-LMM) when
60 the sample size varies (Yan et al., 2019). It means that well-ranked SNPs from one
61 package can be ranked differently in another.

62 Second, there are very few tools focused on the integration of several GWAS
63 software, to make comparisons under different parameters and conditions across

cambie la figura pero esta en .png

65 them. As far as we are aware, there is only two software with this service in mind,
66 which are iPAT and easyGWAS.

67 The iPAT allows running in a graphic interface three well-known command-line
68 GWAS software such as GAPIT, PLINK, and FarmCPU (Zhang et al., 2018). How-
69 ever, 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 and combin-
71 ing several GWAS results. This analysis runs independently of both the computer
72 capacity and the operating system. Nevertheless, it needs either several datasets
73 to obtain the different GWAS results to make replicates or GWAS results already
74 computed. In either case, the results from different algorithms are also separated
75 (Grimm et al., 2017). Thus, although both software iPAT and easyGWAS integrate
76 with different programs or algorithms, an output that allows them to compare simil-
77 itudes and differences in the association is missing.

78 Third, although there are different GWAS software available to repeat the anal-
79 ysis under different conditions (Gumpinger et al., 2018), most of them are designed
80 exclusively for the diploid data matrix (Bourke et al., 2018). Therefore, it is often
81 necessary to "diploidizing" the polyploid genomic data in order to replicate the anal-
82 ysis. The main consequence of this process is missing the complexity of polyploid
83 data to understand how allele dosage affects the phenotype expression (Ferrão et
84 al., 2018).

85 Finally, for polyploid species, any tool that integrates and compares different
86 gene action among software is key to understanding how redundancy or complex
87 interaction among alleles affects the phenotype expression and the evolution of new
88 phenotypes (Bourke et al., 2018; Ferrão et al., 2018; Rosyara et al., 2016).

89 To overcome these challenges, we developed the MultiGWAS tool that performs
90 GWAS analyses for both diploid and tetraploid species using four software in paral-
91 lel. Our tool includes GWASPoly (Rosyara et al., 2016) and the SHEsis tool (Shen
92 et al., 2016) that accept polyploid genomic data, and PLINK (Purcell et al., 2007)
93 and TASSEL (Bradbury et al., 2007), designed exclusively for diploids, but that in
94 the case of tetraploid data, their use require "diploidizing" genomic matrix. This
95 wrapping tool deals with different input file formats, including VCF. Besides, man-
96 age data preprocessing, search for associations by running four GWAS software in
97 parallel, and create a score to choose between gene action models in GWASPoly and
98 TASSEL. Moreover, create comparative reports from the output of each software to
99 help the user distinguish genuine associations from false positives.

100 2 Method

101 The MultiGWAS tool has three main consecutive steps: the adjustment, the multi
102 analysis, and the integration (Fig. 2). In the adjustment step, MultiGWAS pro-
103 cesses the configuration file. Then it cleans and filters the genotype and phenotype
104 datasets, and in case of tetraploids, MultiGWAS "diploidize" the genomic data. Next,
105 during the multi analysis, each GWAS tool runs in parallel. Subsequently, in the in-
106 tegration step, the MultiGWAS tool scans the output files from the four packages
107 (i.e., GWASPoly, SHEsis, PLINK, and TASSEL). Finally, it generates a summary of

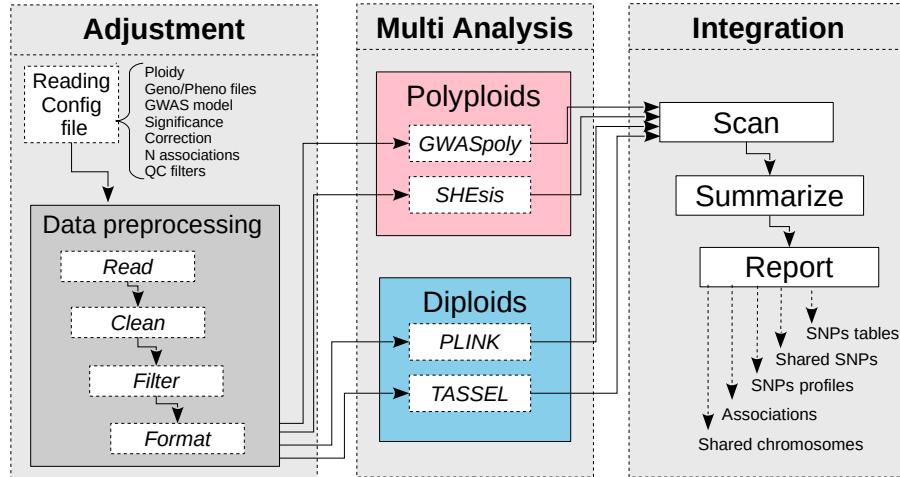


Figure 2: MultiGWAS flowchart has three steps: adjustment, multi analysis, and integration. In the first step, after the input data management upload, MultiGWAS read the configuration file, and preprocess the input data (genotype and phenotype dataset). The second step is the GWAS analysis, where MultiGWAS configure and run the four packages in parallel. Finally, in the third step, MultiGWAS summarize the results and generate a report using different tabular and graphical visualizations.

108 all results that contains score tables, Venn diagrams, SNP profiles, and Manhattan
 109 plots.

110 2.1 Adjustment stage

111 MultiGWAS takes as input a configuration file where the user specifies the genomics
 112 data and the parameters used by the four tools. Once the configuration file is read
 113 and processed, the genomic data files (genotype and phenotype) are then cleaned,
 114 filtered, and checked for data quality. The output of this stage corresponds to the
 115 inputs for the four programs at the Multi Analysis stage.

116 2.1.1 Reading configuration file

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

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

120 **Genotype and phenotype input files:** MultiGWAS uses two input files, one for
 121 the genotype and one for the phenotype. Genotype data can be input in three differ-
 122 ent formats, including a matrix format (Fig. 3.a), a GWASpoly format (Rosyara et
 123 al., 2016) (Fig. 3.b), and Variant Call Format (VCF) (Fig.3.c) which is transformed

124 into GWASpoly format using NGSEP 4.0.2 (Tello et al., 2019). The phenotype file
 125 contains only one trait with the first column for the sample names and the second
 126 column for the trait values (Fig. 3.c).

| a. <pre>Marker, sample01, sample02, sample03, ... c2_41437, AAAG, AAGG, AAAG, ... c2_24258, AAGG, AGGG, GGGG, ... c2_21332, TTCC, TTCC, TTCC, ...</pre> | b. <pre>Marker, Chrom, Pos, sample01, sample02, sample03, ... c2_41437, 0, 805179, AAAG, AAGG, AAAG, ... c2_24258, 0, 1252430, AAGG, AGGG, GGGG, ... c2_21332, 0, 3499519, TTCC, TTCC, TTCC, ...</pre> | c. <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> | d. <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>Individual,Trait</th> </tr> </thead> <tbody> <tr> <td>sample01, 3.59</td> </tr> <tr> <td>sample02, 4.07</td> </tr> <tr> <td>sample03, 1.05</td> </tr> </tbody> </table> | Individual,Trait | sample01, 3.59 | sample02, 4.07 | sample03, 1.05 |
|---|--|---|---|------------------|----------------|----------------|----------------|
| Individual,Trait | | | | | | | |
| sample01, 3.59 | | | | | | | |
| sample02, 4.07 | | | | | | | |
| sample03, 1.05 | | | | | | | |

Figure 3: Examples of MultiGWAS input file formats. Figures a, b and c show examples of genotypes, while figure d shows an example of a phenotype. a. Genotype file in matrix format containing in the first column the marker names and in the following columns the marker data of the samples coded in "ACGT" format (e.g. AAGG, CCTT for tetraploids, AG, CT for diploids). b. Genotype file in GWASpoly format adding the chromosome and marker position to the matrix format. c. 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. d. Phenotype file in a 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).

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

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

138 In this equation, the y is the vector of observed phenotypes. Moreover, the β is a
 139 vector of fixed effects other than SNP or population group effects, the α is a vector of
 140 SNP effects (Quantitative Trait Nucleotides), the ν is a vector of population effects,
 141 the μ is a vector of polygene background effects, and the e is a vector of residual
 142 effects. Besides, Q , modeled as a fixed effect, refers to the incidence matrix for
 143 subpopulation covariates relating y to ν , and X , S and Z are incidence matrices of
 144 1s and 0s relating y to β , α and μ , respectively.

145 **Genome-wide significance:** GWAS searches SNPs associated with the phenotype
 146 in a statistically significant manner. A threshold or significance level α is specified

147 and compared with the *p-value* derived for each association score. Standard significance
148 levels are 0.01 or 0.05 (Gumpinger et al., 2018; Rosyara et al., 2016), and
149 MultiGWAS uses an α of 0.05 for the four GWAS packages. However, in GWASpoly
150 and TASSEL, which calculates the SNP effect for each genotypic class using different
151 gene action models (see “Multi analysis stage”), the threshold is adjusted according
152 to each of those two packages. Therefore, the number of tested markers may be
153 different in each model (see below), impacting the *p-value* thresholds.

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

163 **Number of reported associations:** Criticism has arisen, considering only sta-
164 tistically significant associations as possible correct associations (Kaler and Purcell,
165 2019; Thompson et al., 2011). Many low *p-value* associations, closer to being sig-
166 nificant, are discarded due to the stringent significance levels, which consequently
167 increases the number of false negatives. To avoid this problem, MultiGWAS pro-
168 vides the option to specify the number of best-ranked associations (lower *p-values*),
169 adding the corresponding *p-value* to each association found. In this way, it is pos-
170 sible to enlarge the number of results, and their replicability across the different
171 programs. Nevertheless, the report displays each association with its correspond-
172 ing *p-value*.

173 **Quality control filters:** A control step is necessary to check the input data for
174 the genotype or phenotype errors or poor quality that can lead to spurious GWAS
175 results. MultiGWAS provides the option to select and define thresholds for the fol-
176 lowing filters that control the data quality: Minor Allele Frequency (MAF), individ-
177 ual missing rate (MIND), SNP missing rate (GENO), and Hardy-Weinberg threshold
178 (HWE):

- 179 • **MAF of x:** filters out SNPs with minor allele frequency below x (default 0.01);
- 180 • **MIND of x:** filters out all individuals with missing genotypes exceeding $x*100\%$
181 (default 0.1);
- 182 • **GENO of x:** filters out SNPs with missing values exceeding $x*100\%$ (default
183 0.1);
- 184 • **HWE of x:** filters out SNPs with a *p-value* below the x threshold in the Hardy-
185 Weinberg equilibrium exact test.(for diploids)

186 MultiGWAS does the MAF filtering, and uses the PLINK package (Gumpinger et al.,
187 2018) for the other three filters: MIND, GENO, and HWE.

188 **GWAS tools:** List of names of the four GWAS software to run and integrate into
189 MultiGWAS analysis. They are GWASpoly and SHEsis (designed for polyploid data),
190 and PLINK and TASSEL (designed for diploid data).

191 2.1.2 Data preprocessing

192 Once the configuration file is processed, the genomic data is read and cleaned by
193 selecting individuals present in both genotype and phenotype. Then, MultiGWAS
194 removes individuals and SNPs with poor quality following the previous selected
195 quality-control filters and their thresholds,

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

202 After this process, MultiGWAS converts the genomic data, genotype, and pheno-
203 type datasets to the specific formats required for each of the four GWAS packages.

204 2.2 Multi analysis stage

205 MultiGWAS runs in parallel using two types of statistical models specified in the
206 parameters file, the Full model (Q+K) and Naive (i.e., without any control) where
207 Q refers to population structure, and K refers to relatedness, calculated by kinship
208 coefficients across individuals (Sharma et al., 2018). The Full model (Q+K) controls
209 for both population structure and individual relatedness. For population structure,
210 MultiGWAS uses the Principal Component Analysis (PCA) and takes the top five PC
211 as covariates. For relatedness, MultiGWAS uses kinship matrices that TASSEL and
212 GWASpoly calculated separately, and for PLINK and SHEsis, relatedness depends on
213 kinship coefficients calculated with the PLINK 2.0 built-in algorithm (Chang et al.,
214 2015).

215 2.2.1 GWASpoly

216 GWASpoly (Rosyara et al., 2016) is an R package designed for GWAS in polyploid
217 species used in several studies in plants (Berdugo-Cely et al., 2017; Ferrão et al.,
218 2018; Sharma et al., 2018; Yuan et al., 2019). GWASpoly uses a Q+K linear mixed
219 model with biallelic SNPs that account for population structure and relatedness.
220 Also, to calculate the SNP effect for each genotypic class, GWASpoly provides eight
221 gene action models: general, additive, simplex dominant alternative, simplex dom-
222 inant reference, duplex dominant alternative, duplex dominant, diplo-general, and

223 diplo-additive. Consequently, the number of statistical tests performed can be differ-
224 ent in each action model and so thresholds below which the *p-values* are considered
225 significant.

226 MultiGWAS is using GWASPoly version 1.3 with all gene action models available
227 to find associations. The MultiGWAS reports the top N best-ranked (the SNPs with
228 lowest *p-values*) that the user specified in the N input configuration file. The *full*
229 model used by GWASPoly includes the population structure and relatedness, which
230 are estimated using the first five principal components and the kinship matrix, re-
231 spectively, both calculated with the GWASPoly built-in algorithms.

232 2.2.2 SHEsis

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

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

241 2.2.3 PLINK

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

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

251 2.2.4 TASSEL

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

259 MultiGWAS uses TASSEL 5.0, with all gene action models used to find the N best-
260 ranked associations and reporting the top N best-ranked associations (SNPs with
261 lowest *p-values*). Naive GWAS uses the GLM, and full GWAS uses the MLM with two
262 parameters: population structure that uses the first five principal components, and

263 relatedness that uses the kinship matrix with centered IBS method, both calculated
264 with the TASSEL built-in algorithms.

265 **2.3 Integration stage.**

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

269 **2.3.1 Calculation of *p*-values and significance thresholds**

270 GWAS packages compute *p*-value as a measure of association between each SNP and
271 the trait of interest. The statistically significant associations are those their *p*-value
272 drops below a predefined significance threshold. Since a GWAS analysis performs
273 a large number of tests to look for possible associations, one for each SNP, then
274 some correction in the *p*-values is needed to reduce the possibility of identifying
275 false positives, or SNPs with false associations with the phenotype, but that reach
276 the significance threshold.

277 MultiGWAS provides two methods for adjusting *p*-values and significance thresh-
278 old: the false discovery rate (FDR) that adjust *p*-values, and the Bonferroni cor-
279 rection, that adjusts the threshold. By default, MultiGWAS uses the Bonferroni
280 correction that uses the significance level α/m , with α defined by the user in the
281 configuration file, and m is the number of tested markers to adjust the significance
282 threshold in the GWAS study.

283 However, the significance threshold can be different for each GWAS package as
284 some of them use several action models to calculate the SNP effect of each genotypic
285 class. For both PLINK and SHEsis packages, which use only one model, m is equal
286 to the total number of SNPs. However, for both GWASpoly and TASSEL packages,
287 which use eight and three gene action models, respectively, m is equal to the number
288 of tests performed in each model, which is different between models.

289 Furthermore, most GWAS packages compute both *p*-values and thresholds differ-
290 ently, with the consequence that significant associations identified by one package
291 do not reach the threshold of significance in the others. Thus, it could result in the
292 loss of real associations, the so-called false negatives. To overcome these difficul-
293 ties, MultiGWAS reports two sets of associations: significant and best-ranked (those
294 closest to being statistically significant), as described below.

295 **2.3.2 Selection of significant and best-ranked associations**

296 MultiGWAS reports two groups of associations from the results of the four GWAS
297 packages: the statistically significant associations with *p*-values below a threshold
298 of significance, and the best-ranked associations with the lowest *p*-values, but not
299 reaching the limit to be statistically significant. However, they are representing
300 interesting associations for further analysis (possible false negatives).

301 PLINK and SHEsis have a unique gene action model (see section 2.2.2 and
302 2.2.3). However, in the case of GWASpoly and TASSEL, which have eight and three

303 models respectively, MultiGWAS automatically selects the "best gene action model"
304 from each package and takes the associations from it. This selection within GWAS-
305 Poly and TASSEL has three criteria: the inflation factor (I), the shared SNPs (R),
306 and the significant SNPs (S).

307 Each gene action model is scored using the following equation:

308
$$score(M_i) = I_i + R_i + S_i$$

309 where $score(M_i)$ is the score for the gene action model M_i , with i from 1.. k ,
310 for a GWAS package with k gene action models. I_i is the score for the inflation
311 factor defined as $I_i = 1 - |1 - \lambda(M_i)|$, where $\lambda(M_i)$ is the inflation factor for the
312 M_i model. R_i is the score of the shared SNPs defined as $R_i = \sum_{j=1}^k |M_i \sim M_j|$, where
313 $|M_i \sim M_j|$ is the number of SNPs shared between M_i and M_j models, normalized by
314 the maximum number of SNPs shared between all models. And, S_i is the number of
315 significant SNPs of model M_i normalized by the total number SNPs shared among
316 all models.

317 The score is high when an M_i model has an inflation factor λ close to 1, iden-
318 tifies a high number of shared SNPs, and contains one or more significant SNPs.
319 Conversely, the score is low when the M_i model has an inflation factor λ either low
320 (close to 0) or high ($\lambda > 2$), which identifies a small number of shared SNPs, and
321 contains 0 or few significant SNPs. In any other case, the score results from the
322 balance among the inflation factor, the number of shared SNPs, and the number of
323 significant SNPs.

324 **2.3.3 Integration of results**

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

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

333 To present the replicability, we use two sets: (1) the set of all the significative SNPs
334 provided by each tool and (2) the set of all the best-ranked SNPs. For each set,
335 we present a Venn diagram that displays all SNPs predicted exclusively by one tool
336 and intersections that help identify the SNPs predicted by one, two, three, or all the
337 tools. Also, this information is present on the tables for the two sets.

338 For each SNP identified more than once, MultiGWAS provides its SNP profile.
339 That is a heat diagram for a specific SNP, where each column is a genotype state
340 AAAA, AAAB, AABB, ABAA, and BBBB. Moreover, each row corresponds to a sam-
341 ple. Samples with close genotypes form together clusters. Thus to generate the

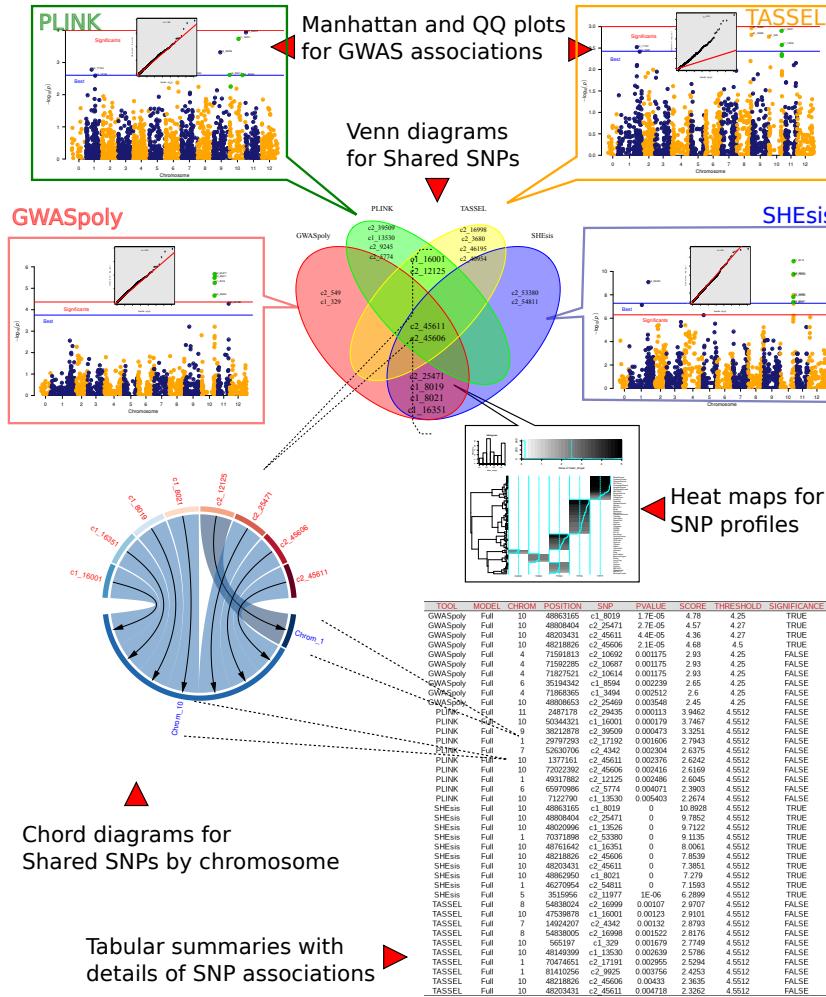


Figure 4: Reports presented by MultiGWAS. For each tool, first, a QQ plot that assesses the resultant p -values. Second, a Manhattan plot for each tool with two lines, blue and red, represents the lower limit for the best ranked and significative SNPs, respectively. We present two Venn diagrams, one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

342 clusters, we do not use the phenotype information. However, we present the phe-
343 notypic information in the figure as the color. This figure visually provides informa-
344 tion regarding genotype and phenotype information simultaneously for the whole
345 population. We present colors as tones between white and black for color blind
346 people.

347 MultiGWAS generates a report, one document with the content previously de-
348 scribed. Besides, there is a folder with the individual figures just in case the user
349 needs one (Supplementary Material 1).

350 In the following section, we present the results of the functionality of the tool,
351 configured with a Full GWAS model using quality filters, and applied on an open
352 dataset of a diversity panel of a tetraploid potato, genotyped and phenotyped as part
353 of the USDA-NIFA Solanaceae Coordinated Agricultural Project (SolCAP) Hirsch et
354 al., 2013. The complete report of this analysis and the report of a second analysis
355 using a naive GWAS model without quality filters are presented in the supplemen-
356 tary materials S1 and S2, respectively.

357 3 Results

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

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

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

374 The complete reports generated by MultiGWAS for both types of analysis, full
375 and naive, applied to the diversity panel of tetraploid potato, are supplementary in-
376 formation at <https://github.com/agrosavia-bioinformatics/MultiGWAS/tree/master/docs/supplements>.

378 3.1 Manhattan and QQ plots for GWAS associations

379 MultiGWAS uses classical Manhattan and Quantile–Quantile plots (QQ plots) to
380 visualize each package’s results. In both plots, the points are the SNPs and their
381 p-values are transformed into scores like $-\log_{10}(p\text{-values})$ (see Fig. 5). The Man-
382 hattan plot shows the strength of association of the SNPs (y-axis) distributed at their
383 genomic location (x-axis), so the higher the score, the stronger the association. At
384 the same time, the QQ plot compares the expected distribution of *p*-values (y-axis)
385 with the observed distribution (x-axis).

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

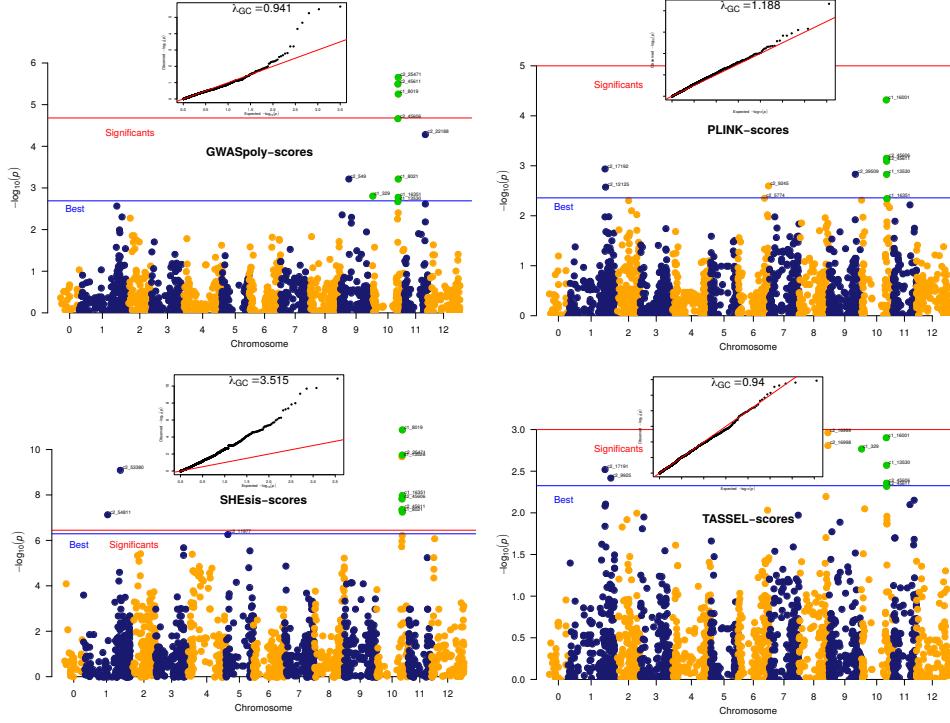


Figure 5: Associations in the tetraploid potato dataset. MultiGWAS shows the associations identified by the four GWAS packages using Manhattan and QQ plots. The tetraploid potato data showed several SNPs shared between the four software (green dots). The best-ranked SNPs are above the blue line, but only GWASpoly and SHEsis identified significant associations (SNPs above the red line) for this dataset. 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.

395 3.2 Tables and Venn diagrams for single and shared SNPs

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

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

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

a.

| TOOL | MDL | IF | SNP | CHR | POS | PVAL | SCR | THR | SGN |
|----------|----------|------|----------|-----|----------|-------|-------|------|-----|
| GWASpoly | additive | 0.99 | c2_25471 | 10 | 48808404 | 0.000 | 5.28 | 4.48 | T |
| GWASpoly | additive | 0.99 | c2_45611 | 10 | 48203431 | 0.000 | 5.07 | 4.48 | T |
| GWASpoly | additive | 0.99 | c1_8019 | 10 | 48863165 | 0.000 | 4.93 | 4.48 | T |
| GWASpoly | additive | 0.99 | c2_45606 | 10 | 48218826 | 0.000 | 4.32 | 4.48 | F |
| GWASpoly | additive | 0.99 | c2_549 | 9 | 16527499 | 0.001 | 3.25 | 4.48 | F |
| GWASpoly | additive | 0.99 | c2_9925 | 1 | 81410256 | 0.002 | 2.77 | 4.48 | F |
| GWASpoly | additive | 0.99 | c1_8021 | 10 | 48862950 | 0.002 | 2.66 | 4.48 | F |
| GWASpoly | additive | 0.99 | c2_12125 | 1 | 71450400 | 0.002 | 2.64 | 4.48 | F |
| PLINK | additive | 1.28 | c1_16001 | 10 | 47539878 | 0.000 | 3.94 | 4.52 | F |
| PLINK | additive | 1.28 | c2_17192 | 1 | 70472766 | 0.001 | 2.86 | 4.52 | F |
| PLINK | additive | 1.28 | c2_12125 | 1 | 71450400 | 0.002 | 2.75 | 4.52 | F |
| PLINK | additive | 1.28 | c2_45606 | 10 | 48218826 | 0.002 | 2.72 | 4.52 | F |
| PLINK | additive | 1.28 | c2_45611 | 10 | 48203431 | 0.002 | 2.64 | 4.52 | F |
| PLINK | additive | 1.28 | c2_14903 | 1 | 87322718 | 0.003 | 2.50 | 4.52 | F |
| PLINK | additive | 1.28 | c1_13530 | 10 | 48149399 | 0.003 | 2.50 | 4.52 | F |
| PLINK | additive | 1.28 | c2_2201 | 1 | 77738822 | 0.003 | 2.49 | 4.52 | F |
| SHEsis | general | 3.56 | c1_8019 | 10 | 48863165 | 0.000 | 10.99 | 4.52 | T |
| SHEsis | general | 3.56 | c1_13526 | 10 | 48020996 | 0.000 | 10.05 | 4.52 | T |
| SHEsis | general | 3.56 | c2_45603 | 10 | 48073593 | 0.000 | 9.89 | 4.52 | T |
| SHEsis | general | 3.56 | c2_25471 | 10 | 48808404 | 0.000 | 9.65 | 4.52 | T |
| SHEsis | general | 3.56 | c2_53380 | 1 | 70371898 | 0.000 | 8.97 | 4.52 | T |
| SHEsis | general | 3.56 | c2_45606 | 10 | 48218826 | 0.000 | 8.17 | 4.52 | T |
| SHEsis | general | 3.56 | c1_16351 | 10 | 48761642 | 0.000 | 8.00 | 4.52 | T |
| SHEsis | general | 3.56 | c2_45611 | 10 | 48203431 | 0.000 | 7.73 | 4.52 | T |
| TASSEL | general | 1.00 | c2_16999 | 8 | 54838024 | 0.001 | 2.96 | 4.52 | F |
| TASSEL | general | 1.00 | c2_4342 | 7 | 14924207 | 0.001 | 2.92 | 4.52 | F |
| TASSEL | general | 1.00 | c2_16998 | 8 | 54838005 | 0.001 | 2.86 | 4.52 | F |
| TASSEL | general | 1.00 | c2_17191 | 1 | 70474651 | 0.002 | 2.67 | 4.52 | F |
| TASSEL | general | 1.00 | c2_9925 | 1 | 81410256 | 0.002 | 2.65 | 4.52 | F |
| TASSEL | general | 1.00 | c1_16001 | 10 | 47539878 | 0.002 | 2.63 | 4.52 | F |
| TASSEL | general | 1.00 | c2_45606 | 10 | 48218826 | 0.005 | 2.34 | 4.52 | F |
| TASSEL | general | 1.00 | c2_45611 | 10 | 48203431 | 0.005 | 2.31 | 4.52 | F |

Column headers: MDL: Model, IF: Inflation factor, SNP: marker name, CHR: Chromosome, PVAL: p-value, SCR: score as -log10 (p-value), THR: significance threshold as -log10 (α / m), where α is the significance level, and m is the number of tested markers, and SGN: significance threshold as true (T) or false (F) whether score > threshold or not.

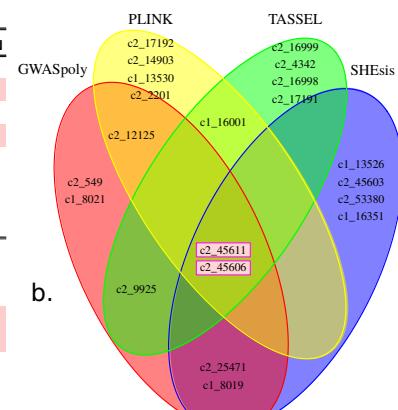
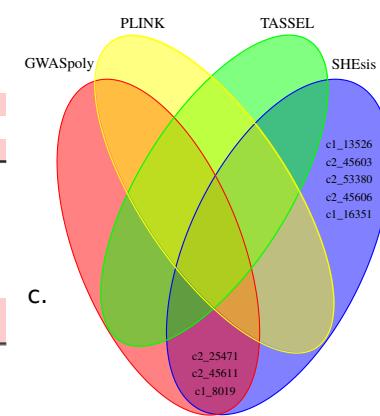
**b.****c.**

Figure 6: Shared SNPs Views. Tabular and graphical views of SNP associations identified by one or more GWAS packages (shared SNPs). SNPs identified by all packages are marker with red background in all figures. (a) Table with details of the N=8 best-ranked SNPs from each GWAS package. Each row corresponds to a single SNP. (b) Venn diagram of the best-ranked SNPs. SNPs identified by all packages are in the central intersection. Other shared SNPs are in both upper central and lower central intersections. (c) Venn diagram of the significant SNPs (score > threshold).

409 In this analysis, both the polyploid packages GWASpoly and SHEsis identified
 410 the SNPs c2_25471, c2_45611, and c1_8019. Of these SNPs, c1_8019 has been
 411 reported in previous studies to be associated with tuber shape and depth of eye
 412 traits (Rosyara et al., 2016; Sharma et al., 2018). Furthermore, in another analy-
 413 sis of MultiGWAS using a naive model without filters (Supplemental Material S2),
 414 the SNP c1_8019 was co-identified by three packages: GWASpoly, SHEsis, and the
 415 diploid PLINK package.

416 3.3 Heat diagrams for the structure of shared SNPs

417 MultiGWAS creates a two-dimensional representation, called the SNP profile, to vi-
 418 sualize each trait by individuals and genotypes as rows and columns, respectively
 419 (Figure 7). At the left, the individuals are grouped in a dendrogram by their geno-
 420 type. At the right, there is the name or ID of each individual. At the bottom, the
 421 genotypes are ordered from left to right, starting from the major to the minor allele
 422 (i.e., AAAA, AAAB, AABB, ABBA, BBBB). At the top, there is a description of the trait
 423 based on a histogram of frequency (top left) and an assigned color for each numer-
 424 ical phenotype value using a grayscale (top right). Thus, each individual appears
 425 as a colored line by its phenotype value on its genotype column. For each column,
 426 there is a solid cyan line with the mean of each column and a broken cyan line that
 427 indicates how far the cell deviates from the mean.

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

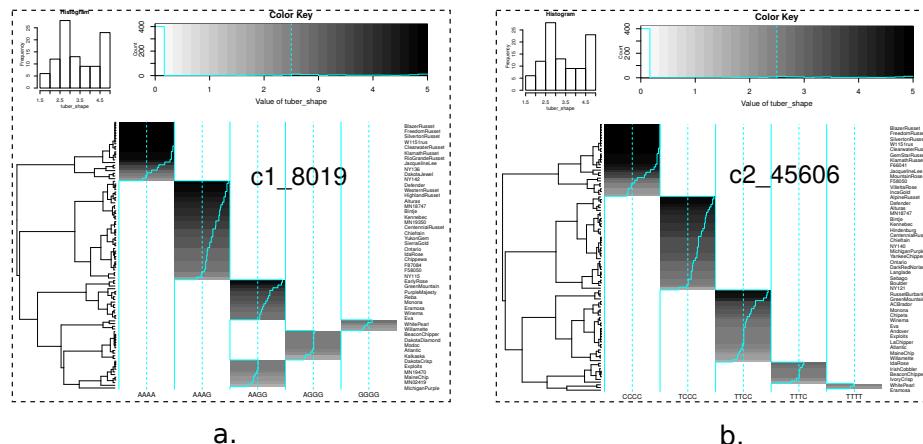


Figure 7: SNP profiles. SNP profiles for two of the best-ranked significant SNPs shown in figure 6.b. (a) SNP c2_45606 best-ranked by the four packages (central intersection of the Venn diagram Figure 6.b) (b) SNP c1_8019 best-ranked by the two tetraploid packages (Figure 6.b), and also identified as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

430 3.4 Chord diagrams for SNPs by chromosome

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

LUIS: Xq no se ven los del
cromosoma 1, en el material
suplementario si estan

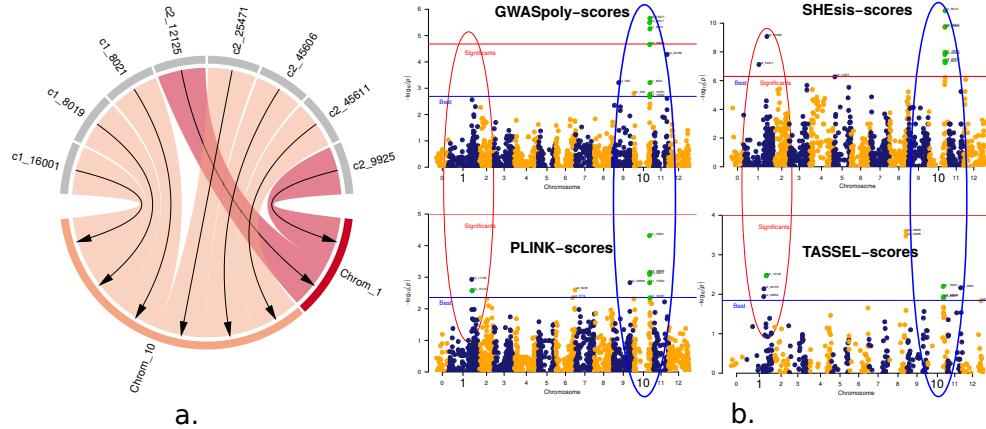


Figure 8: SNPs by chromosome. The position of best-ranked SNPs across chromosomes using two different visualizations. (a) Chord diagram showing that best-ranked SNPs located in chromosome 10. The SNPs are at the top and the chromosomes at the bottom. The arrows connect the best-ranked SNPs with their 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 10). 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 showed in Figure 9.

```

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 organism's ploidy level (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.

452 4.2 Using the command line interface

453 The execution of the CLI tool is easy. It only needs to open a Linux console, change
 454 to the folder where is the configuration file, and type the executable tool's name,
 455 followed by the filename of the configuration file, like this:

456 `multiGWAS Test01.config`

457 Then, the tool starts the execution, showing information on the process in the
 458 console window. When it finishes, the results are in a new subfolder called "*out-Test01*". The results include a complete HTML report containing the different views
 459 described in the results section, the source graphics and tables supporting the re-
 460 port, and the preprocessed tables from the results generated by the four GWAS
 461 packages used by MultiGWAS.
 462

463 4.3 Using the graphical user interface

464 The interface allows users to save, load, or specify the different input parameters
 465 for MultiGWAS in a friendly way (Fig. 10). The input parameters correspond to the
 466 settings included in the configuration file described in subsection 2.1.1. It executes
 467 by calling the following command from a Linux console:

468 `jmultiGWAS`

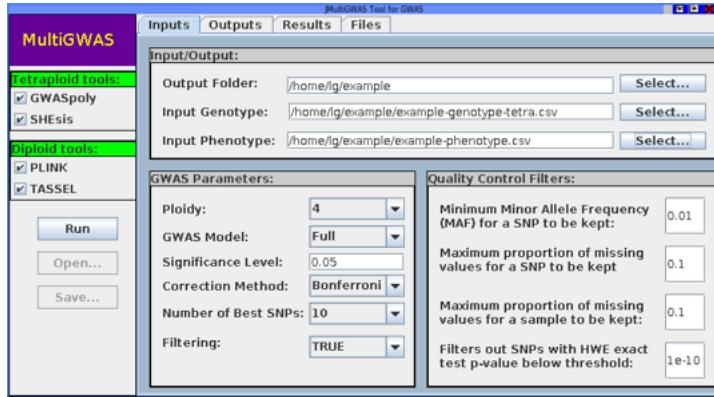


Figure 10: Main view of the MultiGWAS graphical user interface. The interface has a toolbar at the left side and four tabs at the top. In the toolbar, users can select the GWAS packages (Two for tetraploids and two for diploids). The analysis starts with the current parameters or loading a previously saved configuration. In the Input tab, users can set the parameters and quality control filters. The Output tab shows the execution of each process. In the Results tab, users can browse the HTML report of the current analysis generated by the tool. Finally, in the Files tab, users can browse the source files of each software and access the produced data across the analysis.

469 5 Discussion

470 The reanalysis of potato data with MultiGWAS showed that this wrapping tool is
 471 handy to improve the GWAS in tetraploid species. Through MultiGWAS performance,
 472 we could test its effectiveness to answer some of the challenges analyzing
 473 polyploid organisms. They include the integration and replication among parameters
 474 and software, the diploidization of polyploid data, and the incorporation of
 475 different inheritance mechanisms (Dufresne et al., 2014).

476 The main advantage of MultiGWAS is that it replicates the GWAS analysis among
 477 four software and integrates the results obtained across software, models, and parameters.
 478 Depending on the software, users often have to choose between sensitivity or specificity.
 479 However, using MultiGWAS, users do not have to choose between both approaches because they can observe their effect in the analysis within the
 480 same wrapping environment.

lo dejé como diferentes mecanismos de herencia.

482 Another difficulty for replication among software is the variability of structures
 483 for the genomic input data. Currently, the most common format for next-generation
 484 sequencing variant data is the VCF (Variant Call Format) (Danecek et al., 2011;
 485 Ebbert et al., 2014). One of the advantages of VCF is its versatility in summarizing
 486 important genome information for hundreds or thousands of individuals and SNPs,
 487 including information about levels of ploidy. MultiGWAS is different from most of
 488 GWAS software available because it allows the VCF files as an input (but see VarStats
 489 tool in VTC).

490 Moreover, the MultiGWAS is the unique wrapping tool we are aware of that

491 facilitates to understand the effect of diploidizing the tetraploid data in the perfor-
492 mance of the analysis directly. The graphical outputs are a handy approach to find
493 similar results. The SNP profile allows identifying what the significant associations
494 detected by more than one software are. Furthermore, although MultiGWAS checks
495 for significative SNPs based on the *p-value*, it is essential to go back to the data and
496 check if the SNP is a real association between the genotype and phenotype. For this
497 purpose, the SNP profile gives visual feedback for the accuracy of the association.

498 Finally, the MultiGWAS allows comparing among the gene action models that
499 offer GWASPoly and TASSEL. GWASPoly (Rosyara et al., 2016) provides models
500 of different types of polyploid gene action, including additive, diploidized additive,
501 duplex dominant, simplex dominant, and general. On the other hand, TASSEL
502 (Bradbury et al., 2007) also models different types of gene action for general, ad-
503 ditive, and dominant diploids. To choose among models, we propose an automatic
504 selection of the gene action model for both tools based on a balance between three
505 criteria: the inflation factor, the replicability of identified SNPs, and the significance
506 of identified SNPs. . This inflation index is a new tool for comparison that does not
507 offer either GWASPoly or TASSEL. This automatic strategy will help to understand
508 the gene action model for the trait of interest. Although the main focus is on the re-
509 sultant SNPs, the model has assumptions that reflect the gene actions for a specific
510 phenotype.

511 POR HACER.. FUTUROS desarrollos o extensiones y desafios en poliploides.
512 (heterocigosidad, allele dosage, poder de software en poliploides, hitchhiking, HW).
513 Aunque shesis y gwaspoly manejan plodias mayores, hemos visto que hacer gwas
514 para un organismo poliploide requiere tener en cuenta otros factores y por eso aun
515 queremos estudiar esto (HW, Inheritance) antes de extender la herramienta a otras
516 plodias.

517 On the other hand, although MultiGWAS does not solve the uncertainty in the
518 allele dosage and null alleles.

519 MultiGWAS through the active comparison among models addresses the search
520 of the inheritance mechanisms by comparing among two designed for polysomic
521 inheritance software (Rosyara et al., 2016; Shen et al., 2016) with two software for
522 disomic inheritance (Bradbury et al., 2007; Purcell et al., 2007). Understanding the
523 inheritance mechanisms for polyploids organism is an open question. For autopoly-
524 ploids, most loci have a polysomic heritage. However, sections of the genome that
525 did not duplicate lead to disomic inheritance for some loci (Dufresne et al., 2014;
526 Lynch and Conery, 2000; Ohno, 1970). Thus it is a useful tool for researchers be-
527 cause it looks for significative associations that involve both types of inheritance.

534 6 Acknowledgements

535 This research was possible thanks to AGROSAVIA five-years macroproject entitled
536 *Investigación en conservación, caracterización y uso de los recursos genéticos vegetales*.

537 We thanks to the Minister of Science, Technology and Innovation of the re-
538 public of Colombia (previously COLCIENCIAS), for supporting the postdoctoral re-
539 searcher L. Garreta at AGROSAVIA during 2019-2020 under the supervision of ICS

Andrés sugiere aquí discutir brevemente las diferencias en poder de cada software y cómo esta información se integra o se balancea

Andrés sugiere aquí discutir con relación a su comentario anterior, es decir, recalca que ambos criterios no necesariamente pesan lo mismo

la alternativa sería que el usuario diera los pesos a estos parámetros para balancear lo que menciona Andrés?

Andrés sugiere un plot sencillo de heterocigosidad alineado con el de Manhattan para inspeccionar estos casos

estaba encadenado antes con un pedacito que mencionaba que estos eran algunos de los problemas de poliploides y se sugería revisar el efecto de estos en el GWAS pero ese pedacito se perdió y creo que queda fuera de contexto a menos que digamos lo anterior

Comentario de Andrés: Podría haber un caso valioso que solo se captura una vez? (teniendo en cuenta que no todos los métodos tienen el mismo poder, y por ello no necesariamente tendrían el mismo peso)

si pero se muestran todos los que se capturan una vez y los que se capturan varias veces Discusión

Esta es una sugerencia de Andrés: Tener en cuenta a la hora del calcular el score

⁵⁴⁰ and PHRH, (Grant number 811-2019). The editorial of AGROSAVIA gave for fi-
⁵⁴¹ nancial supporting for this publication. Finally to Andres J. Cortes for his valuable
⁵⁴² comments to improve this manuscript.

⁵⁴³ 7 Author Contributions

⁵⁴⁴ LG, ICS, and PHRH conceived the idea. LG developed MultiGWAS. MP tested Multi-
⁵⁴⁵ GWAS. All authors wrote and approved the final version of the manuscript.

⁵⁴⁶ References

- ⁵⁴⁷ Álvarez, M. F., Angarita, M., Delgado, M. C., García, C., Jiménez-Gomez, J., Geb-
⁵⁴⁸ hardt, C., & Mosquera, T. (2017). Identification of Novel Associations of
⁵⁴⁹ Candidate Genes with Resistance to Late Blight in Solanum tuberosum
⁵⁵⁰ Group Phureja. *Frontiers in Plant Science*, 8, 1040. <https://doi.org/10.3389/fpls.2017.01040>
- ⁵⁵¹ Begum, F., Ghosh, D., Tseng, G. C., & Feingold, E. (2012). Comprehensive literature
⁵⁵² review and statistical considerations for gwas meta-analysis. *Nucleic acids*
⁵⁵³ *research*, 40(9), 3777–3784.
- ⁵⁵⁴ Berdugo-Cely, J., Valbuena, R. I., Sánchez-Betancourt, E., Barrero, L. S., & Yock-
⁵⁵⁵ teng, R. (2017). Genetic diversity and association mapping in the colom-
⁵⁵⁶ bian central collection of solanum tuberosum L. Andigenum group using
⁵⁵⁷ SNPs markers. *PLoS ONE*, 12(3). <https://doi.org/10.1371/journal.pone.0173039>
- ⁵⁵⁸ Bourke, P. M., Voorrips, R. E., Visser, R. G. F., & Maliepaard, C. (2018). Tools for Ge-
⁵⁵⁹ netic Studies in Experimental Populations of Polyploids. *Frontiers in Plant*
⁵⁶⁰ *Science*, 9, 513. <https://doi.org/10.3389/fpls.2018.00513>
- ⁵⁶¹ Bradbury, P. J., Zhang, Z., Kroon, D. E., Casstevens, T. M., Ramdoss, Y., & Buckler,
⁵⁶² E. S. (2007). TASSEL: software for association mapping of complex traits
⁵⁶³ in diverse samples. *Bioinformatics*, 23(19), 2633–2635. <https://doi.org/10.1093/bioinformatics/btm308>
- ⁵⁶⁴ Cantor, R. M., Lange, K., & Sinsheimer, J. S. (2010). Prioritizing gwas results: A
⁵⁶⁵ review of statistical methods and recommendations for their application.
⁵⁶⁶ *The American Journal of Human Genetics*, 86(1), 6–22.
- ⁵⁶⁷ Cao, J., Schneeberger, K., Ossowski, S., Günther, T., Bender, S., Fitz, J., Koenig, D.,
⁵⁶⁸ Lanz, C., Stegle, O., Lippert, C., Et al. (2011). Whole-genome sequencing
⁵⁶⁹ of multiple arabidopsis thaliana populations. *Nature genetics*, 43(10), 956.
- ⁵⁷⁰ Chang, C. C., Chow, C. C., Tellier, L. C., Vattikuti, S., Purcell, S. M., & Lee, J. J.
⁵⁷¹ (2015). Second-generation PLINK: Rising to the challenge of larger and
⁵⁷² richer datasets. *GigaScience*, 4(1), arXiv 1410.4803, 1–16. <https://doi.org/10.1186/s13742-015-0047-8>

- 577 Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Hand-
578 saker, R. E., Lunter, G., Marth, G. T., Sherry, S. T., McVean, G., Durbin,
579 R., & Group, 1. G. P A. (2011). The variant call format and VCFtools.
580 *Bioinformatics*, 27(15), <https://academic.oup.com/bioinformatics/article-pdf/27/15/2156/1125001/btr330.pdf>, 2156–2158. <https://doi.org/10.1093/bioinformatics/btr330>
- 583 De, R., Bush, W. S., & Moore, J. H. (2014). Bioinformatics Challenges in Genome-
584 Wide Association Studies (GWAS). In R. Trent (Ed.), *Clinical bioinformatics*
585 (pp. 63–81). New York, NY, Springer New York. https://doi.org/10.1007/978-1-4939-0847-9_5
- 587 Dufresne, F., Stift, M., Vergilino, R., & Mable, B. K. (2014). Recent progress and
588 challenges in population genetics of polyploid organisms: An overview of
589 current state-of-the-art molecular and statistical tools. *Molecular Ecology*,
590 23(1), <https://onlinelibrary.wiley.com/doi/pdf/10.1111/mec.12581>, 40–
591 69. <https://doi.org/10.1111/mec.12581>
- 592 Ebbert, M. T., Wadsworth, M. E., Boehme, K. L., Hoyt, K. L., Sharp, A. R., D O'Fallon,
593 B., Kauwe, J. S., & Ridge, P. G. (2014). Variant tool chest: An improved tool
594 to analyze and manipulate variant call format (vcf) files. *BMC bioinformat-
595 ics*, 15(S7), S12.
- 596 Ekblom, R., & Galindo, J. (2011). Applications of next generation sequencing in
597 molecular ecology of non-model organisms. *Heredity*, 107(1), 1–15.
- 598 Ellegren, H. (2014). Genome sequencing and population genomics in non-model
599 organisms. *Trends in ecology & evolution*, 29(1), 51–63.
- 600 Ferrão, L. F. V., Benevenuto, J., Oliveira, I. d. B., Cellon, C., Olmstead, J., Kirst,
601 M., Resende, M. F. R., & Munoz, P (2018). Insights Into the Genetic Basis
602 of Blueberry Fruit-Related Traits Using Diploid and Polyploid Models in a
603 GWAS Context. *Frontiers in Ecology and Evolution*, 6, 107. <https://doi.org/10.3389/fevo.2018.00107>
605 - Paper for layout. - Many concepts of GWAS, especially structure popula-
606 tion.
- 607 Grimm, D. G., Roqueiro, D., Salomé, P. A., Kleeberger, S., Greshake, B., Zhu, W., Liu,
608 C., Lippert, C., Stegle, O., Schölkopf, B., Weigel, D., & Borgwardt, K. M.
609 (2017). easyGWAS: A Cloud-Based Platform for Comparing the Results of
610 Genome-Wide Association Studies. *The Plant Cell*, 29(1), 5–19. <https://doi.org/10.1105/tpc.16.00551>
- 612 Gumpinger, A. C., Roqueiro, D., Grimm, D. G., & Borgwardt, K. M. (2018). *Methods
613 and Tools in Genome-wide Association Studies* (Vol. 1819).
- 614 Han, B., & Huang, X. (2013). Sequencing-based genome-wide association study in
615 rice. *Current opinion in plant biology*, 16(2), 133–138.
- 616 Hirsch, C. N., Hirsch, C. D., Felcher, K., Coombs, J., Zarka, D., Van Deynze, A., De
617 Jong, W., Veilleux, R. E., Jansky, S., Bethke, P., Douches, D. S., & Buell, C. R.
618 (2013). Retrospective view of North American potato (*Solanum tuberosum*
619 L.) breeding in the 20th and 21st centuries. *G3: Genes, Genomes, Genetics*,
620 3(6), 1003–1013. <https://doi.org/10.1534/g3.113.005595>

- 621 Kaler, A. S., & Purcell, L. C. (2019). Estimation of a significance threshold for genome-
622 wide association studies. *BMC Genomics*, 20(1), 618. <https://doi.org/10.1186/s12864-019-5992-7>
- 623
624 Korte, A., & Farlow, A. (2013). The advantages and limitations of trait analysis with
625 gwas: A review. *Plant methods*, 9(1), 29.
- 626 Lauc, G., Essafi, A., Huffman, J. E., Hayward, C., Knežević, A., Kattla, J. J., Polašek,
627 O., Gornik, O., Vitart, V., Abrahams, J. L., Et al. (2010). Genomics meets
628 glycomics—the first gwas study of human n-glycome identifies hnf1α as a
629 master regulator of plasma protein fucosylation. *PLoS genetics*, 6(12).
- 630 Lynch, M., & Conery, J. S. (2000). The evolutionary fate and consequences of du-
631 plicate genes. *science*, 290(5494), 1151–1155.
- 632 Meng, J., Song, K., Li, C., Liu, S., Shi, R., Li, B., Wang, T., Li, A., Que, H., Li, L., &
633 Zhang, G. (2019). Genome-wide association analysis of nutrient traits in
634 the oyster *Crassostrea gigas*: Genetic effect and interaction network. *BMC*
635 *Genomics*, 20(1), 1–14. <https://doi.org/10.1186/s12864-019-5971-z>
- 636 Ohno, S. (1970). *Evolution by gene duplication*. Springer Berlin Heidelberg. <https://doi.org/10.1007/978-3-642-86659-3>
- 637
638 Pearson, T. A., & Manolio, T. A. (2008). How to interpret a genome-wide association
639 study. *JAMA - Journal of the American Medical Association*, 299(11), 1335–
640 1344. <https://doi.org/10.1001/jama.299.11.1335>
- 641 Power, R. A., Parkhill, J., & De Oliveira, T. (2016). Microbial genome-wide associa-
642 tion studies: lessons from human GWAS. *Nature Reviews Genetics*, 18(1),
643 41–50. <https://doi.org/10.1038/nrg.2016.132>
- 644 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D., Maller,
645 J., Sklar, P., De Bakker, P. I., Daly, M. J., & Sham, P. C. (2007). PLINK: A tool
646 set for whole-genome association and population-based linkage analyses.
647 *American Journal of Human Genetics*, 81(3), 559–575. <https://doi.org/10.1086/519795>
- 648
649 Qiao, H. P., Zhang, C. Y., Yu, Z. L., Li, Q. M., Jiao, Y., & Cao, J. P. (2015). Genetic vari-
650 ants identified by GWAS was associated with colorectal cancer in the Han
651 Chinese population. *Journal of Cancer Research and Therapeutics*, 11(2),
652 468–470. <https://doi.org/10.4103/0973-1482.150346>
- 653 Rosyara, U. R., De Jong, W. S., Douches, D. S., & Endelman, J. B. (2016). Software
654 for Genome-Wide Association Studies in Autopolyploids and Its Applica-
655 tion to Potato. *The Plant Genome*, 9(2), 1–10. <https://doi.org/10.3835/plantgenome2015.08.0073>
- 656
657 Santure, A. W., & Garant, D. (2018). Wild gwas—association mapping in natural
658 populations. *Molecular ecology resources*, 18(4), 729–738.
- 659 Sharma, S. K., MacKenzie, K., McLean, K., Dale, F., Daniels, S., & Bryan, G. J.
660 (2018). Linkage disequilibrium and evaluation of genome-wide associa-
661 tion mapping models in tetraploid potato. *G3: Genes, Genomes, Genetics*,
662 8(10), 3185–3202. <https://doi.org/10.1534/g3.118.200377>
- 663 Shen, J., Li, Z., Chen, J., Song, Z., Zhou, Z., & Shi, Y. (2016). SHEsisPlus, a toolset
664 for genetic studies on polyploid species. *Scientific Reports*, 6, 1–10. <https://doi.org/10.1038/srep24095>
- 665

- 666 Tello, D., Gil, J., Loaiza, C. D., Riascos, J. J., Cardozo, N., & Duitama, J. (2019).
667 NGSEP3: accurate variant calling across species and sequencing protocols.
668 *Bioinformatics*, 35(22), 4716–4723. <https://doi.org/10.1093/bioinformatics/btz275>
- 670 Thompson, J. R., Attia, J., & Minelli, C. (2011). The meta-analysis of genome-wide
671 association studies. *Briefings in Bioinformatics*, 12(3), 259–269. <https://doi.org/10.1093/bib/bbr020>
- 673 Tian, F., Bradbury, P. J., Brown, P. J., Hung, H., Sun, Q., Flint-Garcia, S., Rocheford,
674 T. R., McMullen, M. D., Holland, J. B., & Buckler, E. S. (2011). Genome-
675 wide association study of leaf architecture in the maize nested association
676 mapping population. *Nature genetics*, 43(2), 159–162.
- 677 Yan, Y. Y., Burbridge, C., Shi, J., Liu, J., & Kusalik, A. (2019). Effects of input data
678 quantity on genome-wide association studies (GWAS). *International Journal
679 of Data Mining and Bioinformatics*, 22(1), 19–43. <https://doi.org/10.1504/IJDMB.2019.099286>
- 681 Yu, J., Pressoir, G., Briggs, W. H., Vroh, I. B., Yamasaki, M., Doebley, J. F., McMullen,
682 M. D., Gaut, B. S., Nielsen, D. M., Holland, J. B., Et al. (2006). A unified
683 mixed-model method for association mapping that accounts for multiple
684 levels of relatedness. *Nature genetics*, 38(2), 203–208.
- 685 Yuan, J., Bizimungu, B., De Koeyer, D., Rosyara, U., Wen, Z., & Lagüe, M. (2019).
686 Genome-Wide Association Study of Resistance to Potato Common Scab.
687 *Potato Research*. <https://doi.org/10.1007/s11540-019-09437-w>
- 688 Zhang, S., Chen, X., Lu, C., Ye, J., Zou, M., Lu, K., Feng, S., Pei, J., Liu, C., Zhou, X.,
689 Ma, P., Li, Z., Liu, C., Liao, Q., Xia, Z., & Wang, W. (2018). Genome-wide
690 association studies of 11 agronomic traits in cassava (*Manihot esculenta*
691 crantz). *Frontiers in Plant Science*, 9(April), 1–15. <https://doi.org/10.3389/fpls.2018.00503>