# Journal of Visualized Experiments

## A pathway association study tool for GWAS analyses of metabolic pathway information

### --Manuscript Draft--

To the Editors of JoVE;

Attached, please find our manuscript entitled "How to run the Pathway Association Study Tool and interpret the results of GWAS analyses in light of metabolic pathway information", which we submit for publication as peer reviewed Scientific Video. We have been invited to submit this paper by Lyndsay Troyer, with whom we have been working on the submission.

Genome wide association studies (GWAS) have become very common in maize, and in some cases can uncover useful variants for the study of the genetic architecture or practical improvement of a trait. In many cases, however, particularly in the case of quantitative traits, we believe that GWAS is not finding all useful information contained in the analyses. In particular, the need to set such a stringent significance threshold on results ensures that many of the results are discarded. Therefore, we created a new tool called PAST (the Pathway Association Study Tool) that assigns SNPs from a GWAS study to genes, and genes to metabolic pathways. The tool then recalculates the probability that the entire pathway is associated with the trait of interest. We have found that since we began working on this tool, a few other groups have released tools, especially for human genetics studies, but we believe no other tool is as versatile, robust, and user friendly than PAST. Lynsday suggested we could present how to use PAST in a JoVE video, and we present here the manuscript that will accompany the video, when it is finished.

We hope you agree that this manuscript is descriptive and useful. My coauthor Adam will present PAST in the video, and we hope this will allow us to reach a much wider audience with our new and potentially highly impactful tool.  Thank you very much in advance for your consideration.


Sincerely,

Marilyn Warburton

1 **TITLE:**
2 **A Pathway Association Study Tool for GWAS Analyses of Metabolic Pathway Information**
3

4 **AUTHORS AND AFFILIATIONS:**
5 Adam Thrash[1], Marilyn L. Warburton[2]
6
7 [1]Institute for Genomics, Biocomputing & Biotechnology, Mississippi State University, Mississippi
8 State, MS, USA
9 [2]USDA-ARS Corn Host Plant Resistance Research Unit, Mississippi State, MS, USA
10
11 **Corresponding Author:**
12 Adam Thrash                        (thrash@igbb.msstate.edu)
13
14 **Email Addresses of Co-Authors:**
15 Marilyn L. Warburton         (marilyn.warburton@usda.gov)
16

21 **SUMMARY:**
22 By running the Pathway Association Study Tool (PAST), either through the Shiny application or
23 through the R console, researchers can gain a deeper understanding of the biological meaning of
24 their genome-wide association study (GWAS) results by investigating the metabolic pathways
25 involved.
26

27 **ABSTRACT:**
28 Recently, a new implementation of a previously described method for interpreting genome-wide
29 association study (GWAS) data using metabolic pathway analysis has been developed and
30 released. The Pathway Association Study Tool (PAST) was developed to address concerns with
31 user-friendliness and slow-running analyses. This new user-friendly tool has been released on
32 Bioconductor and Github. In testing, PAST ran analyses in less than one hour that previously
33 required twenty-four or more hours. In this article, we present the protocol for using either the
34 Shiny application or the R console to run PAST.
35

36 **INTRODUCTION:**
37 Genome-wide association studies (GWAS) are a popular method of studying complex traits and
38 the genomic regions associated with them[1-3]. In this type of study, hundreds of thousands of
39 single nucleotide polymorphism (SNP) markers are tested for their association with the trait, and
40 the significance of the associations is assessed. Marker-trait associations that meet the false
41 discovery rate (FDR) threshold (or some other type of significance threshold) are retained for the
42 study, but true associations may be filtered out. For complex, polygenic traits, the effect of each
43 gene might be small (and thus filtered out), and some alleles are only expressed in specific
44 conditions that might not be present in the study[3]. Thus, while many SNPs may be retained as

45  associated with the trait, each may have a very small effect. Too many SNP calls will be missing,
46  and an interpretation of the biological meaning and genetic architecture of the trait may be
47  incomplete and confusing. Metabolic pathway analysis can help to address some of these issues
48  by focusing on the combined effects of genes grouped according to their biological function[4-6].
49
50  Several studies were completed using a previous implementation of the method described in this
51  article. Aflatoxin accumulation[7], corn earworm resistance[8], and oil biosynthesis[9] were all studied
52  with the previous implementation. While these analyses were successful, the process of analysis
53  was complicated, time-consuming, and cumbersome, because the analysis tools were written in
54  a combination of R, Perl, and Bash, and the pipeline was not automated. Because of the
55  specialized knowledge required to modify this method for each analysis, a new method has now
56  been developed that can be shared with other researchers.
57
58  The Pathway Association Study Tool (PAST)[10] was designed to address the shortcomings of the
59  previous method by requiring less knowledge of programming languages and by running analyses
60  in a shorter period. While the method was tested with maize, PAST makes no species-specific
61  assumptions. PAST can be run through the R console, as a Shiny app, and an online version is
62  expected to soon be available on MaizeGDB.
63
64  **PROTOCOL:**
65
66  **1. Setup**
67
68  1.1. Install R, if it is not already installed.
69
70  NOTE: PAST is written in R and, therefore, requires that its users have R installed. At the time of
71  this writing, installing PAST directly from Bioconductor requires R3.6, but PAST can be installed
72  from Github for users with R3.5. R installation instructions can be downloaded from the following
73  link: https://www.r-project.org/.
74
75  1.2. Install the latest version of RStudio Desktop or update RStudio (optional).
76
77  NOTE: RStudio is a helpful environment for working with the R language. Its installation is
78  recommended, especially for those who choose to run PAST in the command line rather than
79  through the Shiny GUI application. RStudio and its installation instructions can be found at the
80  following link: https://rstudio.com/products/rstudio/.
81
82  1.3. Install PAST from Bioconductor[11] by following the instructions on Bioconductor.
83
84  NOTE: Installation through Bioconductor should handle the installation of PAST's dependencies.
85  Additionally, PAST can be installed from Github[12], but installing from Github will not install
86  dependencies automatically.
87
88  1.4. Install PAST Shiny (optional). Download the file "app.R" from the Releases page of the Github

89 repository: https://github.com/IGBB/PAST/releases/, and remember where the downloaded file
90 is located.

91

92 NOTE: PAST can be used by calling its methods directly with R, but users who are less familiar
93 with R can run the PAST Shiny application, which provides a guided user interface. PAST Shiny is
94 an R script available in the shiny_app branch of the PAST Github repository. PAST Shiny will
95 attempt to install its dependencies during the first run.

96

97 1.5. Begin analysis by starting the application in one of the three ways described below.

98

99 1.5.1. PAST Shiny with RStudio

100

101 1.5.1.1. Using RStudio, create a new project in the folder where app.R is located. Click **File | New**
102 **Project** and select that folder.

103

104 1.5.1.2. Once a new project has been created, open the app.R file downloaded earlier. RStudio
105 recognizes that app.R is a Shiny app and creates a **Run App** button on the bar above the displayed
106 source code. Click **Run App**. RStudio will then launch a window that displays the PAST Shiny
107 application.

108

109 1.5.2. PAST Shiny with R Console

110

111 1.5.2.1. Launch R and run the following code to start the PAST Shiny application:
112 *shiny::runApp('path/to/folder/with/shiny/app.R*'. Replace the text in quotes with the folder to
113 which app.R was downloaded, and keep the quotes.

114

115 1.5.3. PAST without R Shiny

116

117 1.5.3.1. Run library(PAST) in an R Console to load PAST.

118

119 **2. Customize Shiny analysis (optional)**

120

121 2.1. Change the analysis title from "New Analysis" to something that better reflects the type of
122 analysis being run which helps to keep track of multiple analyses (see **Figure 1**).

123

124 [Place **Figure 1** here]

125

126 2.2. Modify the number of cores and the mode. Set the number of cores to any number between
127 1 and the total number on the machine but be aware that devoting more resources to PAST may
128 slow down other operations on the machine. Set the mode based on the description in section
129 6.

130

131 **3. Load GWAS data**

132

133 NOTE: Verify that the GWAS data is tab delimited. Ensure that the association file contains the
134 following columns: trait, marker name, locus or chromosome, position on the chromosome, p-
135 value, and $R^2$ value for the marker. Ensure that the effects file contains the following columns:
136 trait, marker name, locus or chromosome, and position on the chromosome. The order of these
137 columns is not important, as the user can specify the names of the columns when loading the
138 data. Any additional columns are ignored. TASSEL[13] can be used to produce these files.
139
140 3.1. Load GWAS data with PAST Shiny.
141
142 3.1.1. Select an association file and an effects file by using the **Association File** and **Effects File**
143 selection boxes. Change the column names in the **Association Column Name** and **Effects**
144 **Columns Name** input boxes below the file selection boxes to reflect the column names in the
145 data.
146
147 [Place **Figure 2** here]
148
149 3.2. Load GWAS data with PAST in the R Console.
150
151 3.2.1. Modify and run the following code:
152
153 *gwas_data = load_GWAS_data("path/to/association_file.tsv", "path/to/effects_file.tsv",*
154 *association_columns = c("Trait", "Marker", "Locus", "Site", "p", "marker_R2"), effects_columns =*
155 *c("Trait", "Marker", "Locus", "Site", "Effect")*
156
157 NOTE: Change the paths to the actual location of the GWAS files. The values provided for
158 association_columns and effects_columns are the default values. If the names do not match the
159 default values, specify the column names. Otherwise, these can be omitted.
160
161 **4. Load linkage disequilibrium (LD) data**
162
163 NOTE: Verify that the linkage disequilibrium (LD) data is tab delimited and contains the following
164 types of data: Locus, Position1, Site1, Position2, Site2, Distance in base pairs between Position1
165 and Position2, and $R^2$ value.
166
167 4.1. Load LD data with PAST Shiny.
168
169 4.1.1. Select the file containing LD data. Change the column names in the **LD Column Names** input
170 boxes below the file selection box to match the column names in the LD data if necessary.
171
172 [Place **Figure 3** here]
173
174 4.2. Load LD Data with PAST in the R Console.
175
176 4.2.1. Modify and run the following code to load LD data:

177

178  *LD = load_LD(“path/to/LD.tsv”, LD_columns = c("Locus1", "Position1", "Site1", "Position2",*
179  *"Site2", "Dist_bp", "R.2")*

180

181  NOTE: Change the path to the actual location of the LD file. The values provided for LD_columns
182  are the default values. If the names do not match these defaults, specify the correct names of
183  the columns; otherwise, these can be omitted.

184

185  **5. Assign SNPs to genes**

186

187  NOTE: Download or otherwise locate annotations in GFF format. These annotations can often be
188  found in online databases for specific organisms. Be cautious about low quality annotations, as
189  the quality of the annotations data will affect the quality of the pathway analysis. Confirm that
190  the first column of these annotations (the chromosome) matches the format of the
191  locus/chromosome in the association, effects, and LD data. For example, the annotations should
192  not call the first chromosome "chr1" if the GWAS and LD data files call the first chromosome "1".

193

194  5.1. Assign SNPs to genes with PAST Shiny.

195

196  NOTE: More information about determining an appropriate $R^2$ cutoff can be found in Tang et al.[6],
197  in the section called "SNP to gene algorithm for the pathway analysis".

198

199  5.1.1. Select the file containing GFF annotations. Consider what window size and $R^2$ cutoff are
200  most suitable for the species being considered and modify if the defaults do not suit the uploaded
201  data.

202

203  NOTE: Default values in PAST primarily reflect values appropriate for maize. The number of cores
204  set at the beginning of the PAST Shiny analysis (Step 2.2) is used in this step.

205

206  [Place **Figure 4** here]

207

208  5.2. Assign SNPs to genes with PAST in the R Console.

209

210  5.2.1. Modify and run the following code to assign SNPs to genes:

211

212  *genes = assign_SNPs_to_genes(gwas_data, LD, "path/to/annotations.gff", c("gene"), 1000, 0.8,*
213  *2)*

214

215  NOTE: In this sample code, several default suggestions are provided: 1000 is the size of the
216  window around the SNP to search for genes; 0.8 is the cutoff value for $R^2$; 2 is the number of
217  cores used for parallel processing. The path to the annotations should also be changed to the
218  actual location of the annotations file.

219

220  **6. Discover significant pathways**

221

222 NOTE: Verify that the pathways file contains the following data in tab delimited format, with one
223 line for every gene in each pathway: pathway ID – an identifier such as "PWY-6475-1"; pathway
224 description – a lengthier description of what the pathways do such as "trans-lycopene
225 biosynthesis"; gene – a gene in the pathway, which should match the names provided in the
226 annotations. Pathway information can likely be found in online databases for specific organisms,
227 such as MaizeGDB. The second user-specified option is the mode. "Increasing" refers to
228 phenotypes that reflect when an increasing value of the measured trait is desirable, such as yield,
229 while "decreasing" refers to a trait where a decrease in the measured values is beneficial, such
230 as insect damage ratings. The significance of pathways is tested using previously described
231 methods[4,6,14].

232

233 6.1. Discover significant pathways with PAST Shiny.

234

235 6.1.1. Select the file containing pathways data and be sure that the mode is selected in the
236 analysis options. If necessary, change the number of genes that must be in a pathway to retain it
237 for the analysis and the number of permutations used to create the null distribution to test
238 significance of effect.

239

240 [Place **Figure 5** here]

241

242 NOTE: The number of cores and the mode set at the beginning of the PAST Shiny analysis (Step
243 2.2) is used in this step. The default number of genes is currently set at 5 genes, so pathways with
244 fewer known genes will be removed. The user can lower this value to 4 or 3, to include shorter
245 pathways, but doing so will risk false positive results. Increasing this value can increase the power
246 of the analysis but will remove more pathways from the analysis. Changing the number of
247 permutations used, increases and decreases the power of the test.

248

249 6.2. Discover significant pathways with PAST in the R Console.

250

251 6.2.1. Modify and run the following code to discover significant pathways:

252

253 *rugplots_data <- find_pathway_significance(genes, "path/to/pathways.tsv", 5, "increasing",*
254 *1000, 2)*

255

256 NOTE: In this sample code, several suggested defaults are provided. 5 is the minimum number of
257 genes that must be in a pathway in order to keep the pathway in the analysis, increasing refers
258 to an increasing amount of the measured trait (it is recommended that the user run both
259 increasing and decreasing, regardless of trait; data interpretation will differ for the two,
260 however), 1000 is the number of times to sample the effects to determine the null distribution,
261 and 2 is the number of cores used for parallel processing. Change the path to the actual location
262 of the pathways file.

263

264 **7. View Rugplots**

266  7.1. View Rugplots with PAST Shiny.

267

268  7.1.1. Once all inputs are uploaded and set, click **Begin Analysis**. A progress bar will appear and
269  indicate which step of the analysis was last completed. When the analysis completes, PAST Shiny
270  will switch to the **Results** tab. A table of results will be displayed in the left column (labeled
271  "pathways") and the Rugplots will be displayed in the right column (labeled "plots").

272

273  [Place **Figure 6** here]

274

275  7.1.2. Use the slider to control the filtering parameters. When the filtering level is satisfactory,
276  click the **Download Results** button at the bottom left to download all images and tables
277  individually to a ZIP file that is named with the Analysis title. This ZIP file contains the filtered
278  table, the unfiltered table, and one image per pathway in the filtered table.

279

280  [Place **Figure 7** here]

281

282  7.2. View Rugplots with PAST in the R Console

283

284  7.2.1. Modify and run the following code to save the results:

285

286  *plot_pathways(rugplots_data, "pvalue", 0.02, "increasing", "output_folder")*

287

288  NOTE: In this sample code, several suggested defaults are provided. pvalue provides the data
289  that can be used for filtering insignificant pathways after a significance threshold is chosen by the
290  user; 0.02 is the default value used in filtering, and increasing refers to an increasing amount of
291  the measured trait (it is recommended that the user run both increasing and decreasing,
292  regardless of trait; data interpretation will differ for the two, however); output_folder is the
293  folder where the images and tables will be written (this folder must exist prior to running the
294  function). A table of filtered results, the unfiltered results, and individual images for every
295  pathway in the filtered results are written to this folder.

296

297  **REPRESENTATIVE RESULTS:**
298  If results are not produced following a run of the PAST software tool, check to be sure that all
299  input files are correctly formatted. A successful run using the example data in the PAST package,
300  which are based on a maize GWAS of grain color, is shown in **Figure 8**. This table and the resulting
301  image can be downloaded using the Download Results button. An example of the downloaded
302  image is shown in **Figure 2**[10]. Incorrect settings might lead to results that do not make biological
303  sense, but determining incorrectness must be up to the researcher, who should double check the
304  validity of the chosen settings and consider all known evidence regarding the trait of interest.

305

306  **Figure 9**[10] shows the rugplot produced from the pathway analysis of GWAS results created with
307  a maize panel of 288 inbred lines that had been phenotyped for grain color. This simplistic
308  example, where the phenotypes were either "white" or "yellow", was used because the pathway

309 responsible for creating the bright yellow carotenoid pigments is known and should be
310 responsible for most of the phenotype. Thus, we expected to see the trans-lycopene biosynthesis
311 pathway (which produce carotenoids) to be significantly associated with grain color, which it is.
312 Pathway ID and name are listed at the top of the graph. The horizontal axis of the graph ranks all
313 genes that were included in the analysis, arranged from left to right in order of largest effect on
314 the trait to smallest. However, only the genes in the trans-lycopene biosynthesis pathway are
315 marked (at the top of the graph, as hatch marks, appearing in the gene rank of their effect as
316 compared to all other genes in the analysis). There are 7 genes in this pathway. The running
317 enrichment score (ES) is plotted along the vertical axis. The ES for each gene is added into the
318 running total in order of effect and the total is adjusted to the number of genes analyzed. Thus,
319 the score changes as one moves right along the horizontal axis and tends to increase as the larger
320 effect genes are included, but at some point, the increase in the effect is smaller than the
321 adjustment for having added another gene, and the entire score begins to decrease. The apex of
322 the running ES line is marked with a dotted vertical line; this is the ES for the entire pathway and
323 is used by the program to determine if the pathway is chosen and presented as a rugplot.
324
325 **FIGURE LEGENDS:**
326 **Figure 1: Step 2.1.**
327
328 **Figure 2: Step 3.1.**
329
330 **Figure 3: Step 4.1.**
331
332 **Figure 4: Step 5.1.**
333
334 **Figure 5: Step 6.1.**
335
336 **Figure 6: Step 7.1.1.**
337
338 **Figure 7: Step 7.1.2.**
339
340 **Figure 8: Completed run of PAST Shiny.**
341
342 **Figure 9: Pathway image from completed run of PAST (or downloaded from Shiny).** This figure
343 has been cited from Thrash et al.[10].
344
345 **DISCUSSION:**
346 A primary goal of PAST is to bring metabolic pathway analyses of GWAS data to a wider audience,
347 especially for non-human and non-animal organisms. Alternative methods to PAST are often
348 command-line programs that focus on humans or animals. User-friendliness was a primary goal
349 in the development of PAST, both in choosing to develop a Shiny application and in choosing to
350 use R and Bioconductor to release the application. Users do not need to learn how to compile
351 programs in order to use PAST.
352

353 As with most types of analysis software, the results of PAST are only as good as the input data; if
354 the input data has errors or is incorrectly formatted, PAST will fail to run or produce
355 uninformative results. Ensuring that the GWAS data, LD data, annotations, and pathways files are
356 correctly formatted is critical to receiving correct output from PAST. PAST only analyzes bi-allelic
357 markers and can run only one trait for each set of input data. In addition, GWAS data produced
358 by poor genotyping or incorrect or imprecise phenotyping is not likely to produce clear or
359 repeatable results either. PAST can aid in the biological interpretation of GWAS results but is
360 unlikely to clarify chaotic data sets if environmental variation, experimental error, or population
361 structure was not properly accounted for.
362
363 Users can choose to change some parameters of the analysis, both in the Shiny application and
364 by passing those parameters to PAST's functions in the R console. These parameters can change
365 the results reported by PAST, and users should take care when modifying these from the defaults.
366 Because LD is measured by the users, typically using the same marker data set that was also used
367 in the GWAS, the LD measurements are specific to the population. For all studies, especially for
368 species other than maize, (particularly self-pollinating, polyploid, or highly heterogenous
369 species), changes in the defaults may be warranted.
370

377 **REFERENCES:**
378 1.      Rafalski, J. Association genetics in crop improvement. *Current Opinion in Plant Biology*. **13**
379 (2), 174–180 (2010).
380 2.      Yan, J., Warburton, M., Crouch, J. Association Mapping for Enhancing Maize (Zea *mays* L.)
381 Genetic Improvement. *Crop Science*. **51** (2), 433–449 (2011).
382 3.      Xiao, Y., Liu, H., Wu, L., Warburton, M., Yan, J. Genome-wide Association Studies in Maize:
383 Praise and Stargaze. *Molecular Plant*. **10** (3), 359–374 (2017).
384 4.      Wang, K., Li, M., Bucan, M. Pathway-Based Approaches for Analysis of Genomewide
385 Association Studies. *The American Journal of Human Genetics*. **81** (6), 1278–1283 (2007).
386 5.      Weng, L. et al. SNP-based pathway enrichment analysis for genome-wide association
387 studies. *BMC Bioinformatics*. **12** (1), 99 (2011).
388 6.      Tang, J., Perkins, A., Williams, W., Warburton, M. Using genome-wide associations to
389 identify metabolic pathways involved in maize aflatoxin accumulation resistance. *BMC Genomics*.
390 **16** (1), 673 (2015).
391 7.      Warburton, M. et al. Genome-Wide Association Mapping of Aspergillus flavus and
392 Aflatoxin Accumulation Resistance in Maize. *Crop Science*. **55** (5), 1857–1867 (2015).
393 8.      Warburton, M. et al. Genome-Wide Association and Metabolic Pathway Analysis of Corn
394 Earworm Resistance in Maize. *The Plant Genome*. **11** (1), 170069 (2018).
395 9.      Li, H., Thrash, A., Tang, J., He, L., Yan, J., Warburton, M. Leveraging GWAS data to identify
396 metabolic pathways and networks involved in maize lipid biosynthesis. *The Plant Journal*. **98** (5),
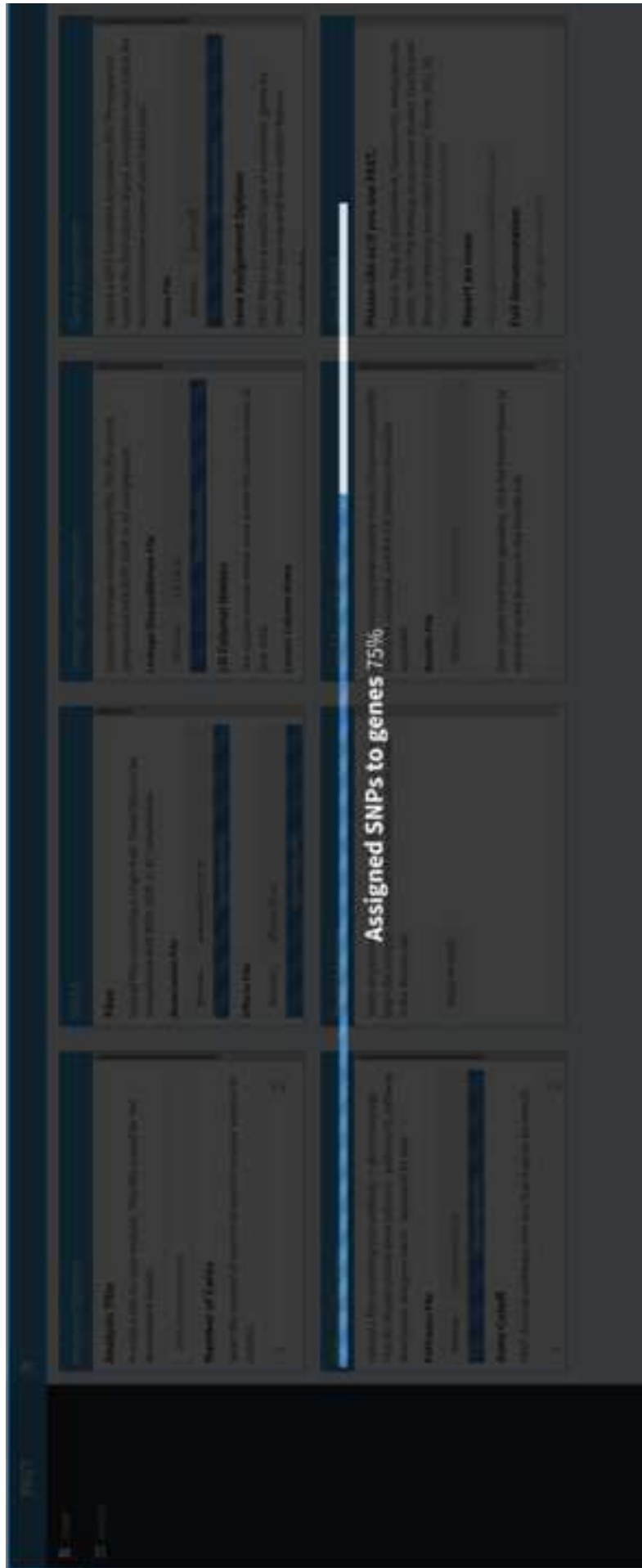
397  853–863 (2019).

398  10.  Thrash, A., Tang, J., DeOrnellis, M., Peterson, D., Warburton, M. PAST: The Pathway

399  Association Studies Tool to Infer Biological Meaning from GWAS Datasets. *Plants*. **9** (1), 58 (2020).
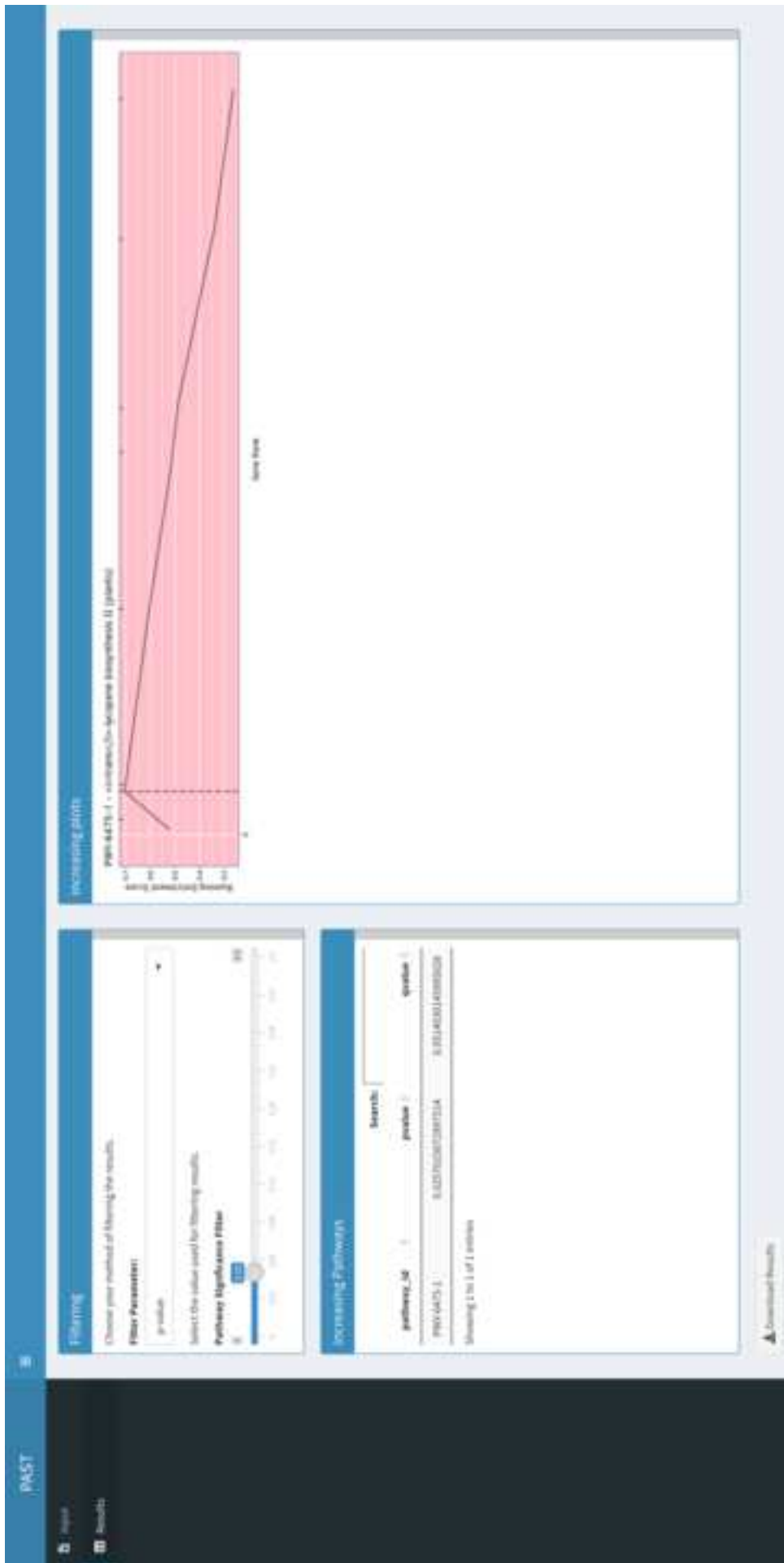
400  11.  Adam, T., Mason, D. PAST: Pathway Association Study Tool (PAST). doi:

401  10.18129/B9.bioc.PAST. Bioconductor version: Release (3.10). (2020).
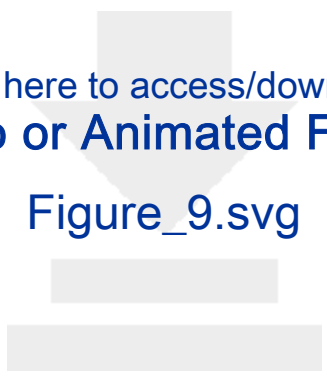
402  12.  Thrash, A., DeOrnellis, M. *IGBB/PAST*. at <https://github.com/IGBB/PAST>. IGBB. (2019).

403  13.  Bradbury, P. et al. TASSEL: software for association mapping of complex traits in diverse

404  samples. *Bioinformatics*. **23** (19), 2633–2635 (2007).

405  14.  Subramanian, A. et al. Gene set enrichment analysis: A knowledge-based approach for

406  interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences*

407  *U.S.A.* 102. 15545–15550 (**2005**).

Assigned SNPs to genes 75%

| Name of Material/ Equipment | Company | Catalog Number | Comments/Description |
|---|---|---|---|
| Computer | NA | NA | Any computer with 8GB RAM should be sufficien |
| R | R Project | NA | R 3.6 or greater is required to install from Bioco |

nt
nductor

To Whom It May Concern:

Please find my remarks addressing editorial and reviewer concerns below. Editorial and reviewer remarks are **bold**, while my remarks are not.

# Editorial Concerns

**1. Please take this opportunity to thoroughly proofread the manuscript to ensure that there are no spelling or grammar issues. The JoVE editor will not copy-edit your manuscript and any errors in the submitted revision may be present in the published version.**

Some errors were corrected in this revision.

**2. Please submit each figure as a vector image file to ensure high resolution throughout production: (.psd, ai, .eps., .svg).**

An SVG for Figure 2 was uploaded. Figure 1 is a screenshot, and vectorizing a raster image wouldn't do much, to my knowledge. Instead, a higher resolution PNG has been attached for Figure 1.

**3. Please ensure that the references appear as the following: [Lastname, F.I., LastName, F.I., LastName, F.I. Article Title. Source. Volume (Issue), FirstPage – LastPage (YEAR).] For more than 6 authors, list only the first author then et al.**

I edited the author names and order of the references.

**4. Please include volume and issue numbers for all references.**

More feedback would be useful on this remark. I checked the references, and the only ones without volume and issue numbers are websites.

**5. Please define all abbreviations before use.**

I have made edits to ensure that all abbreviations are defined before first use.

**6. Please revise the table of the essential supplies, reagents, and equipment. The table should include the name, company, and catalog number of all relevant materials in separate columns in an xls/xlsx file. Please sort the Materials Table alphabetically by the name of the material.**

This paper describes a software tool. There are no supplies or reagents.

**7. Please revise the title for conciseness: A pathway association study tool for GWAS analyses of metabolic pathway information**

The title has been changed as suggested.

**8. Please ensure that all text in the protocol section is written in the imperative tense as if telling someone how to do the technique (e.g., "Do this," "Ensure that," etc.). The actions should be described in the imperative tense in complete sentences wherever possible. Avoid usage of phrases such as "could be," "should be," and "would be" throughout the Protocol. Any text that cannot be written in the imperative tense may be added as a "Note." However, notes should be concise and used sparingly. Please include all safety procedures and use of hoods, etc.**

The protocol has been edited to reflect these suggestions. Because PAST is a software tool, there are no safety procedures specific to PAST.


**9. The Protocol should contain only action items that direct the reader to do something. Please move the discussion about the protocol to the Discussion.**

Remarks on the protocol have either been edited to be imperative or moved.

# Reviewer Concerns

**--It would be good to mention that there are other tools for visualization of omics data (e.g. Pathway Collage through Pathway Tools, MetaMapp, etc.; can be useful in the context of examining differential expression), but that this one is specifically tailored towards GWAS results (or any other distinguishing features that you would like to point out).**

I've added a note about PAST's specific design as a GWAS tool.


**--Line 40: would suggest changing 'FDR threshold' to 'FDR or another type of significance threshold', given that not all GWAS-involving analyses use FDR.**

I have edited as suggested.


**--Is it correct that the metabolic pathway annotations/nomenclature being used are based solely on the input file from the user? It seems so given that the tool would otherwise not be free of species-specific assumptions (as mentioned in line 59), given that it would have needed to assume a certain master set of pathways/model of plant metabolism if not accepting external pathway identifiers (as do other mapping/visualization tools). If solely based on user input, it would likely be good to mention common sources of this data (GFF annotation data and pathway information) and to encourage caution with annotations (particularly if electronically inferred). These would be in addition to the warnings already included, which state that the output is only as good as the input data (so that those warnings are not interpreted to be referring only to the genotypic and phenotypic data). It may also be good to highlight (if true) that the tool will only analyze pathways that are labeled with exactly the same identifier (e.g., if one is upstream of the other but is labeled differently, that they will be analyzed as separate pathways).**

Yes, that's correct. I've added some notes about being wary of the quality of the annotations and that the gene names in the annotations should match those in the pathways file.

**--Perhaps this is already planned, but it would be great in the video to show the two workflows (Shiny and R console) separately, as it is otherwise a bit tricky in written form to see the complete workflow in either case.**

I think this would be the best approach as well.

**--Lines 143-146: can background LD level also be specified, if examining genes in the vicinity of a significant SNP (or is this not a view that will be enabled)? Can a reference be mentioned that describes which values are standard in maize?**

No, PAST doesn't yet have a way to specify background LD. That may be worth investigating as we continue to improve PAST. Tang *et al* (2015) describes how the default $R^2$ was decided upon.

**--Line 172: it seems that there optimally would not be a default for search space window, rather would need to be the result of some analysis relevant to the data set being analyzed. That said, 1 kb is a very conservatively small window, so seems OK (just may have lots of false negatives, but seems better to aim on the too-small side than too-large).**

We do expect that those who use PAST will change this parameter to reflect their knowledge of their species.

**--Line 174: could the default be set to be half of the cores that the user has available? It seems that this would greatly improve the processing time (which was mentioned to be long even in this improved version), if parallel processing is already enabled.**

It's technically possible, but letting the user set this on their own is better. Users should explicitly grant tools permission to use more resources; PAST has no way of knowing whether a user is running resource-intensive processes using the rest of their cores. Based on this feedback, the default has been reduced to 1 in PAST Shiny. The code examples still use 2 cores, since the user is explicitly specifying the number of cores.

**--Lines 187-192: Does significance refer solely to significance of effect estimate for a SNP near a gene? And 'significant pathways' (as mentioned in lines 189-190) would then be the pathways having more than the cutoff number of genes with significant effect estimates? Or has another test been run to determine which pathways are 'significant'? Please be very specific and explicit re: which 'significance' tests are being run within this tool.**

Another test has been run, but the details of this test are described in other manuscripts. I've added a line to inform readers where they can find details on how pathways are tested for significance.

**--Line 218: 'associated' may not have been intended to be taken literally, but \*SNPs\* are 'associated' in a GWAS context, not genes or pathways (unless some of the involved significance tests explicitly test for significance of association of a pathway with marker-trait associations? Pathways are not themselves being tested in GWAS.) Perhaps another word such as 'corresponding to', to avoid confusion given that this paper somewhat operates in a GWAS context?**

I changed this line to match line 201, since the parameters are the same and should be described in the same way.

**--Line 243: are enrichment scores (statistically) of such a nature that they can be cleanly added? One example: PVEs are not quite additive. A bit more detail re: gene set enrichment analysis (as a method) would be useful, if that is what has been used here. (Gene set enrichment analysis was what came back when looking up 'enrichment scores' online.)**

Because this method was described in the release paper for PAST, we didn't discuss it here. However, one of the citations added to address concerns with lines 187-192 does link to Subramanian et al's GSEA paper, where the details of the calculation are described.

**--Lines 272-273: some mention of within-species variation might be worthwhile. For example, it has been found that LD patterns vary substantially even between the various subpopulations of maize. The same has been seen in wild vs. cultivated barley.**

I've expanded our description to indicate that the user generates LD and thus the LD measurements are specific to the population they're studying.

**--Figure 1: Would all 'significant' pathways be depicted within this view? Could the names of the genes and/or number of SNPs in their vicinity exhibiting significant association with the trait under analysis be somehow depicted?**

All significant pathways are depicted in this view, but the size of the image changes depending on how many pathways there are. Individual pathway images are written to the results folder or downloaded from PAST Shiny. The names of the genes are included in the table output that is written to the results folder or downloaded from PAST Shiny. Depicting more information becomes complicated for pathways with many genes, as the image gets quite crowded.