

1 **GenFam: A new web application for gene family-based classification and functional**
2 **enrichment analysis of plant genomes**

3 Renesh Bedre¹ and Kranthi Mandadi^{1,2*}

4 ¹Texas A&M AgriLife Research & Extension Center, 2415 E. Hwy. 83, Weslaco, Texas 78596

5 ²Department of Plant Pathology & Microbiology, Texas A&M University, 2132 TAMU, College
6 Station, Texas 77843

7 *To whom correspondence should be addressed.

8

9

10 **Running Title:** Gene family-based enrichment analysis

11 **ABSTRACT**

12 **Motivation:** Genome-scale studies using next-generation sequencing technologies generate
13 substantial number of differentially-regulated genes. The gene lists need to be further analyzed to
14 identify overrepresented genes and functions in order to guide downstream analyses. Currently
15 available gene enrichment tools rely on functional classifications based on Gene Ontology (GO)
16 terms. A shortcoming of the GO-based classification system is that the GO terms are broad and
17 often redundant, hence necessitating alternate approaches.

18 **Results:** We propose a new functional enrichment approach, GenFam, to classify as well as
19 enrich overrepresented gene functions, based on gene family categories. GenFam offers a unique
20 approach to mine valuable, biologically-relevant information, beyond the conventional GO term
21 based enrichment. GenFam is available as a web-based, graphical-user interface, which allows
22 users to readily input gene lists, and export results in both tabular and graphical formats.
23 Additionally, users can customize analysis parameters, by choosing from the different
24 significance tests to conduct advanced statistics. Currently, GenFam supports gene family
25 classification and enrichment analyses for seventy-eight plant genomes and gene identifiers that
26 are available on Phytozome v12.0 database.

27 **Availability and implementation:** The GenFam application is open-source and accessible
28 through world-wide web at <http://mandadilab.webfactional.com/home/>

29 **Contact:** kkmandadi@tamu.edu

30 **Supplementary information:** Supplementary File 1 and 2

31

32 1 INTRODUCTION

33 In recent years, genome-wide analyses using next-generation sequencing (NGS) technologies,
34 have become indispensable to life science research. Generating large-scale datasets has become
35 relatively straightforward, as opposed to efficiently interpreting the data to gain intuition into
36 biologically-significant mechanisms. Data mining tools that determine, predict, and enrich
37 putative functions among NGS datasets are highly valuable for such genomic analyses (Backes
38 *et al.*, 2007). For instance, RNA-sequencing (RNA-seq) analyses is a high-throughput approach
39 to study transcriptome regulation by determining transcript-level changes in multiple cell- or
40 tissue-types, or among varying experimental conditions (e.g., unstressed vs. stressed). In a
41 typical RNA-seq experiment, the analysis yields hundreds, if not thousands, of genes that are
42 differentially expressed among the experimental conditions. Uncovering enriched biological
43 pathways among these gene lists is a valuable starting step for downstream genetic analyses.

44 The Gene Ontology (GO)-term based enrichment tools (e.g., BinGO, Blast2GO, AgriGO)
45 are commonly used by researchers to infer the enriched pathways in NGS experiments (Bedre *et al.*
46 *et al.*, 2016; Bedre *et al.*, 2015; Chen *et al.*, 2013; Li *et al.*, 2017; Mandadi and Scholthof, 2015;
47 Mandadi and Scholthof, 2012; Schaker *et al.*, 2016). These tools identify overrepresented GO
48 terms associated within a user-defined list of genes by mapping them to the background genome
49 annotations, and calculating statistical probability of enrichment relative to the background. The
50 enrichment tools can classify genes into GO categories or pathways related to biological process,
51 molecular function and cellular locations (Du *et al.*, 2010; Goffard and Weiller, 2007). However,
52 the GO classifications are often broad and provide limited information on specific biological
53 attributes of the gene (Ashburner *et al.*, 2000). For instance, GO terms in molecular function
54 such as nucleic acid binding (GO:0003676) and DNA binding (GO:0003677) do not provide
55 further information on the class of gene that is being enriched. Further, enriched GO terms can
56 be redundant, that need to be manually filtered before interpretation. Given these shortcomings,
57 new methods to analyze and interpret large-scale datasets to gain further insights into
58 biologically-meaningful information are needed.

59 In this study, we present a unique approach to perform classification and enrichment
60 analysis of genes, based on gene family (GenFam). The GenFam offers a meaningful way to
61 determine pertinent gene functions by directly classifying and enriching genes, in a user-defined

62 list, based on the encoded-protein and its associated gene family. We present GenFam as a user-
63 friendly, graphical-user interface application that can be launched on the world-wide web.

64

65 **2 IMPLEMENTATION AND DATA ANALYSIS**

66 **2.1 Background database**

67 GenFam classifies and enriches genes into 128 representative and unique gene families, based on
68 the well-annotated reference plant genome, *Arabidopsis thaliana* (Berardini *et al.*, 2015).

69 GenFam currently supports analysis of genes from seventy-eight plant species. The background
70 gene family database for the genomes was manually curated to remove redundancy among the
71 families. Furthermore, we also determined a common protein domain structure for each gene
72 family based on the protein sequences of the family members. The protein domains were
73 predicted using HMMER (v3.1b2) from protein family database (Pfam release 31.0) (Eddy,
74 2009; Finn *et al.*, 2015). A multi-step annotation approach was used to classify gene sequences
75 to a gene family. First, gene families were assigned based on their sequence similarity to
76 *Arabidopsis* orthologs. Next, remaining sequences were assigned to a gene family based on their
77 Pfam protein domain signature. All the selected 128 gene families, individual gene sequences,
78 and corresponding gene IDs were formatted using the PostgreSQL database to perform
79 classification and enrichment analysis using various statistical methods.

80 **2.2 Statistical enrichment methods**

81 GenFam provides two main functions: i) classification, and ii) enrichment of user-defined gene
82 lists. The enrichment analysis is based on the singular enrichment analysis methods (Huang da *et*
83 *al.*, 2009). In a manner similar to GO term enrichment tools (Backes, *et al.*, 2007; Du, *et al.*,
84 2010; Huang da, *et al.*, 2009), GenFam utilizes the user-defined gene IDs as input to perform
85 statistical enrichment analysis. GenFam accepts different types of gene IDs for the analysis, as
86 defined by the Phytozome database. For example, for rice, it accepts locus (LOC_Os01g06882),
87 transcripts (LOC_Os01g06882.1) and PAC (24120792) IDs. To determine acceptable IDs for all
88 plant species, user can use the “check allowed ID type for each species” function on the GenFam
89 analysis page. Once the appropriate gene IDs are provided, GenFam classifies and identifies
90 specific gene families and members that are overrepresented in the input gene lists. A unique

91 feature of GenFam is that it only utilizes genes categorized to gene family as a reference
92 background, unlike the GO enrichment tools which utilizes the entire genome as a reference
93 background. This feature greatly enhances the sensitivity of the enrichment analysis. GenFam
94 can employ rigorous statistical tests such as the Fisher exact, Chi-Square, Binomial distribution
95 and hypergeometric tests, along with multiple test corrections to control family-wise error rate,
96 in order to report the statistically significant enriched genes.

97 As a default test, GenFam performs the Fisher exact test, which relies on the proportion
98 of observed data, instead of a value of a test statistic to estimate the probability of genes of
99 interest corresponding to a specific category. For instance, suppose we have n differentially
100 expressed genes, and among them, k falls in a particular gene family category, and there are m
101 total genes associated with that gene family in the background reference database among N total
102 genes; then Fisher probability that a given gene family is overrepresented in the input of gene list
103 is calculated as,

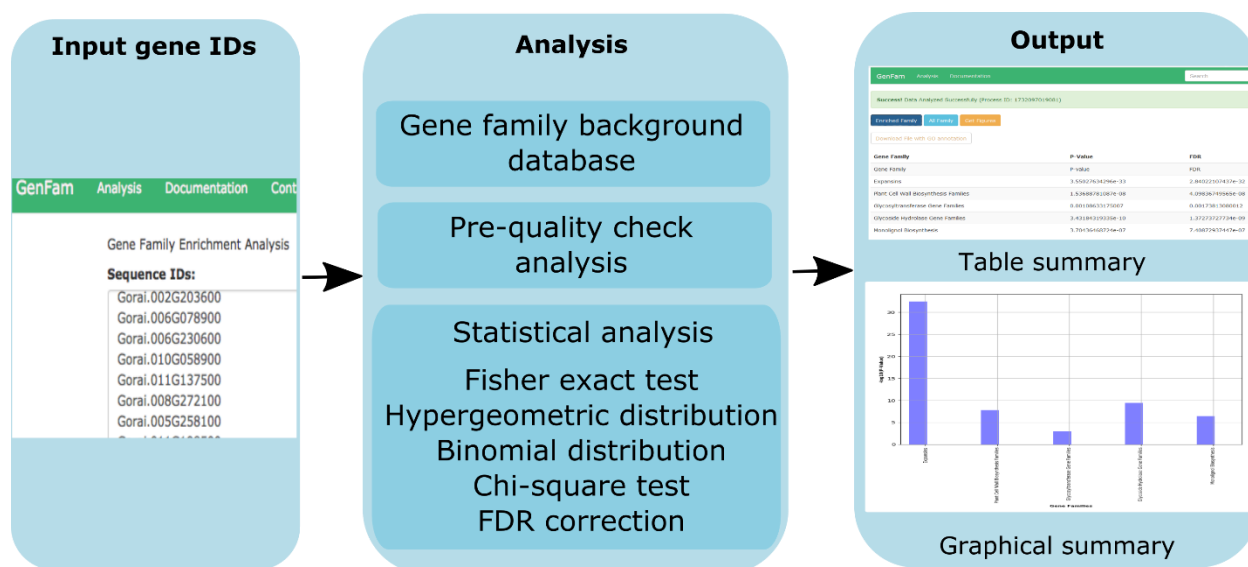
$$104 \quad p = \frac{\binom{n}{k} \binom{N-n}{m-k}}{\binom{N}{m}}$$

105 To address the false positives resulting from multiple comparisons especially when the
106 input gene list is large (>1000), GenFam subsequently employs false discovery rate (FDR)
107 correction methods including the Benjamini-Hochberg (Benjamini and Hochberg, 1995),
108 Bonferroni (Bonferroni, 1936) and Bonferroni-Holm (Holm, 1979). The various statistical tests
109 and FDR methods can be customized by the user as appropriate. Along with enrichment results
110 for the gene families, GenFam also provides information related to GO terms in biological
111 process, molecular function and cellular component categories associated with the enriched gene
112 families. These results can be downloaded as a tabular file (“Enriched Families”) or as a
113 graphical figure of the enriched families (“Get Figures”). If users only want to retrieve the
114 classification of genes, GenFam parses another tabular file containing information of all the
115 annotated gene families (“All Families”).

116 **2.3 Web server implementation**

117 The GenFam web server is implemented using Python3 (<https://www.python.org/>), Django
118 1.11.7 (<https://www.djangoproject.com/>) and PostgreSQL (<https://www.postgresql.org/>)
119 database. All the codes for data formatting and statistical analysis are implemented using Python
120 scripting language. The high-level Python web framework was constructed using Django. The
121 Django web framework was hosted using WebFaction (<https://www.webfaction.com/>). The web-
122 based templates were designed using Bootstrap, HTML, and CSS. GenFam is compatible with
123 all major browsers including Internet Explorer, Microsoft Edge, Google Chrome, Mozilla and
124 Safari. All the precomputed plant gene family background databases were built using advanced
125 PostgreSQL database. The analyzed data was visualized using the matplotlib (Droettboom *et al.*,
126 2016) Python plotting library.

127



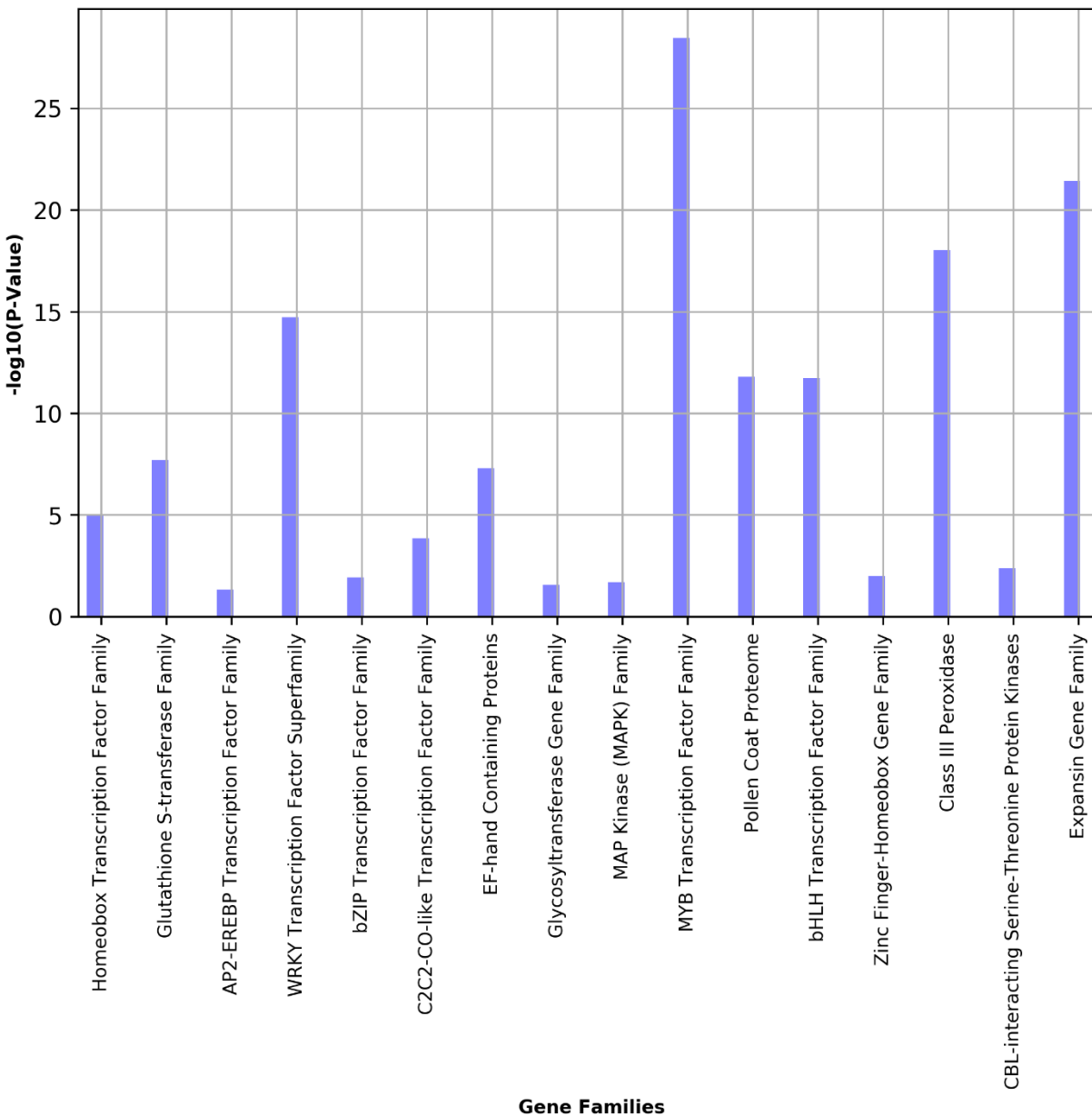
128

129 **Fig. 1.** GenFam workflow. The list of input gene IDs for respective plant species provided by the user are
130 analyzed for enrichment analysis using various statistical tests. The output of the analysis can be viewed
131 and/or downloaded as a table and/or graphical summary. The results page has multiple options to
132 visualize or download data for both enriched and non-enriched categories (all gene families). The detailed
133 output data from a case study are provided in Supplementary Files 1 and 2.

134 2.4 Case study and data analysis

135 To demonstrate the utility of GenFam, we performed two case studies using cotton (a dicot) and
136 rice (a monocot) transcriptome datasets (Bedre, *et al.*, 2015; Dametto *et al.*, 2015). We have

137 previously identified ~662 differentially expressed genes in cotton infected with *Aspergillus*
138 *flavus* (Bedre, *et al.*, 2015). For the first case study, we used GenFam to determine the enriched
139 gene families among these 662 differentially expressed genes, using the options of Fisher exact



140

141

142 **Fig. 2.** Graphical summary of GenFam enrichment analysis of a cotton case study. Results are plotted as
143 bar chart using the $-\log_{10}(\text{P-Value})$ scores. Higher the $-\log_{10}(\text{P-Value})$ value, greater the confidence in
144 enrichment of the gene family.

145 test for statistical enrichment, and the Benjamini-Hochberg (Benjamini and Hochberg, 1995)
146 method to control FDR. The GenFam classification and enrichment analysis revealed
147 overrepresented gene families such as expansins, kinases, peroxidases, and transcription
148 factors—genes that we have hypothesized to mediate cell-wall modifications, antioxidant
149 activity and defense signaling in response to *A. flavus* infection (Bedre, *et al.*, 2015) (**Fig. 1** and
150 **2; Supplementary File 1**). In the second case study, we analyzed ~758 genes which were up-
151 regulated in a cold-tolerant rice genotype (Dametto, *et al.*, 2015). GenFam was able to
152 successfully classify and determine enriched gene families related to aquaporins, peroxidases,
153 glutathione S-transferases, as well as gene families involved in cell wall-related mechanisms
154 (**Supplementary File 2**)—genes that were hypothesized by Dametto *et al.* (2015) to play a role
155 in the rice cold stress response. Together, the information of classified and enriched gene
156 families not only provides understanding of the affected biological processes, but allows the user
157 to readily select favorite gene families for further downstream characterization.

158 A snapshot of the analysis page and workflow is shown in **Fig. 1**. Users have the option
159 to either use the default settings or select desired statistical parameters. The analysis page also
160 guides the users to select gene IDs that are acceptable in GenFam (**Fig. 1**). Users are directed to
161 the results after analysis is completed (**Fig. 1**).

162 **2.5 Output summary**

163 The results are displayed as summary table (HTML) and graphical chart plotted using the -
164 $\log_{10}(\text{P-Value})$ scores. Higher the $-\log_{10}(\text{P-Value})$ value, greater the confidence in enrichment of
165 the gene family (**Fig. 2**). The enriched and non-enriched gene family results can also be
166 downloaded as tabular files, with further details of associated P-value and FDR statistics, and
167 GO terms.

168

169 **3 DISCUSSION**

170 Data mining of big datasets (e.g., NGS data) is a very important step, and approaches that can
171 systematically dissect biologically-relevant information from big data are highly desirable. GO
172 term-based enrichment analyses, although commonly employed, does not provide specific,
173 biologically-relevant, gene family level information. Further, GO classifications can be broad

174 and redundant. We suggest that GenFam is a unique way to extract biologically-relevant, gene
175 family level information among large-scale results. GenFam allows users to readily uncover
176 biologically-relevant functions enriched in large-scale gene datasets by classifying and providing
177 specific information about the enriched gene families— information that could not be inferred by
178 GO enrichment analysis alone. Furthermore, unlike GO enrichment tools, instead of using the
179 whole genome as a background database for enrichment analysis, GenFam uses only genes
180 annotated and classified into a gene family as a reference. This feature ensures decreasing
181 enrichment bias and increasing the accuracy of the analysis (Huang da, *et al.*, 2009). GenFam
182 can be implemented with various statistical enrichment methods such as Fisher exact test,
183 hypergeometric distribution, chi-square test and binomial distribution, thus providing flexibility
184 in the analysis based on the sample size and user preferences. We recommend using Fisher exact
185 test, chi-square test and hypergeometric distribution for smaller datasets (< 1000) (McDonald,
186 2009), and binomial distribution for larger datasets (Khatri and Draghici, 2005; Zheng and
187 Wang, 2008). To control the false positives, GenFam also supports multiple testing corrections
188 (family-wise error rate) algorithms such as Benjamini-Hochberg (Benjamini and Hochberg,
189 1995), Bonferroni (Bonferroni, 1936), and Bonferroni-Holm (Holm, 1979).

190 In conclusion, we suggest that GenFam provides a unique approach to interpret
191 biologically relevant information in big datasets by directly classifying and representing
192 overrepresented genes into gene families. This allows users to readily interpret and identify
193 favorite genes for downstream inquiries.

194

195 **FUNDING**

196 This work was supported by start-up funds from Texas A&M AgriLife Research to K.K.M.

197

198 **ACKNOWLEDGEMENTS**

199 We thank Dr. Sonia Irigoyen (Texas A&M AgriLife Research) for comments and feedback
200 during the preparation of this manuscript.

201

202 CONFLICT OF INTEREST

203 The authors declare no conflict of interest.

204

205 SUPPLEMENTARY DATA

206 **Supplementary File 1:** List of the differentially regulated genes and analysis output of the
207 cotton case study.

208 **Supplementary File 2:** List of the differentially regulated genes and analysis output of the rice
209 case study.

210

211 REFERENCES

- 212 Ashburner, M. et al. (2000) Gene ontology: tool for the unification of biology. The Gene
213 Ontology Consortium. *Nat. Genet.*,**25**,25-29.
- 214 Backes, C. et al. (2007) GeneTrail--advanced gene set enrichment analysis. *Nucleic Acids*
215 *Res.*,**35**,W186-192.
- 216 Bedre, R. et al. (2016) Transcriptome analysis of smooth cordgrass (*Spartina alterniflora*
217 Loisel), a monocot halophyte, reveals candidate genes involved in its adaptation to
218 salinity. *BMC Genomics*,**17**,657.
- 219 Bedre, R. et al. (2015) Genome-wide transcriptome analysis of cotton (*Gossypium hirsutum* L.)
220 identifies candidate gene signatures in response to aflatoxin producing fungus *Aspergillus*
221 *flavus*. *PLoS One*,**10**.
- 222 Benjamini, Y. and Hochberg, Y. (1995) Controlling the false discovery rate - a practical and
223 powerful approach to multiple testing. *J. Roy. Stat. Soc. Ser. B. (Stat. Method.)*,**57**,289-
224 300.
- 225 Berardini, T.Z. et al. (2015) The Arabidopsis information resource: making and mining the “gold
226 standard” annotated reference plant genome. *Genesis*,**53**,474-485.
- 227 Bonferroni, C.E. Teoria statistica delle classi e calcolo delle probabilita. Libreria internazionale
228 Seeber; 1936.
- 229 Chen, E.Y. et al. (2013) Enrichr: interactive and collaborative HTML5 gene list enrichment
230 analysis tool. *BMC Bioinformatics*,**14**,128.
- 231 Dametto, A. et al. (2015) Cold tolerance in rice germinating seeds revealed by deep RNAseq
232 analysis of contrasting indica genotypes. *Plant Sci.*,**238**,1-12.
- 233 Droettboom, M. et al. matplotlib: matplotlib v1. 5.1. In.: doi; 2016.
- 234 Du, Z. et al. (2010) agriGO: a GO analysis toolkit for the agricultural community. *Nucleic Acids*
235 *Res.*,**38**,W64-W70.
- 236 Eddy, S.R. A new generation of homology search tools based on probabilistic inference. In,
237 *Genome Informatics 2009: Genome Informatics Series Vol. 23*. World Scientific; 2009. p.
238 205-211.

- 239 Finn, R.D. et al. (2015) The Pfam protein families database: towards a more sustainable future.
240 *Nucleic Acids Res.*,**44**,D279-D285.
- 241 Goffard, N. and Weiller, G. (2007) PathExpress: a web-based tool to identify relevant pathways
242 in gene expression data. *Nucleic Acids Res.*,**35**,W176-W181.
- 243 Holm, S. (1979) A simple sequentially rejective multiple test procedure. *Scand. J. Stat.*,**6**,65-70.
- 244 Huang da, W., Sherman, B.T. and Lempicki, R.A. (2009) Bioinformatics enrichment tools: paths
245 toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Res.*,**37**,1-
246 13.
- 247 Khatri, P. and Draghici, S. (2005) Ontological analysis of gene expression data: current tools,
248 limitations, and open problems. *Bioinformatics*,**21**,3587-3595.
- 249 Li, Y. et al. (2017) Global identification of alternative splicing via comparative analysis of
250 SMRT- and Illumina-based RNA-seq in strawberry. *Plant J.*,**90**,164-176.
- 251 Mandadi, K.K. and Scholthof, K.-B.G. (2015) Genome-wide analysis of alternative splicing
252 landscapes modulated during plant-virus interactions in *Brachypodium distachyon*. *Plant*
253 *Cell*,**27**,71-85.
- 254 Mandadi, K.K. and Scholthof, K.B. (2012) Characterization of a viral synergism in the monocot
255 *Brachypodium distachyon* reveals distinctly altered host molecular processes associated
256 with disease. *Plant Physiol.*,**160**,1432-1452.
- 257 McDonald, J.H. Handbook of biological statistics. Sparky House Publishing Baltimore, MD;
258 2009.
- 259 Schaker, P.D. et al. (2016) RNAseq transcriptional profiling following whip development in
260 sugarcane smut disease. *PLoS One*,**11**,e0162237.
- 261 Zheng, Q. and Wang, X.J. (2008) GOEAST: a web-based software toolkit for Gene Ontology
262 enrichment analysis. *Nucleic Acids Res.*,**36**,W358-363.

263

264

265

266

267

268

269

Input gene IDs

GenFam Analysis Documentation Cont

Gene Family Enrichment Analysis

Sequence IDs:

Gorai.002G203600
Gorai.006G078900
Gorai.006G230600
Gorai.010G058900
Gorai.011G137500
Gorai.008G272100
Gorai.005G258100

Analysis

Gene family background database

Pre-quality check analysis

Statistical analysis

Fisher exact test
Hypergeometric distribution
Binomial distribution
Chi-square test
FDR correction

Output

GenFam Analysis Documentation

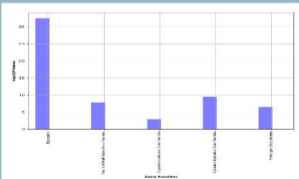
Success! Data Analyzed Successfully (Process ID: 1733397014661)

[Enriched Family](#) [All Families](#) [Cell Figures](#)

[Download File with GO Annotations](#)

Gene Family	P-Value	FDR
Other Family	Protein	706
Enzymes	3.55E-05429E-04	3.84E-021043E-04
HMGB-Like-Wall-Associated-Protein	1.32E-00210E-06	4.00E-021043E-06
Oxyanion-Binding-Domain	6.3E-0007E-0507	6.70E-021043E-07
Protein-Disulfide-Isomerase	2.43E-04103E-10	3.27E-021043E-10
Membrane-Associated-Protein	2.7E-0004E-0404	3.48E-021043E-04

Table summary



Graphical summary

