

# 1 Status of Genome Function Annotation in Model Organisms and

## 2 Crops

3 Bo Xue and Seung Y Rhee\*

4 Department of Plant Biology, Carnegie Institution for Science, Stanford, CA 94305, USA

5 Emails: [srhee@carnegiescience.edu](mailto:srhee@carnegiescience.edu) (SYR); [bxue@carnegiescience.edu](mailto:bxue@carnegiescience.edu) (BX)

6 \*Corresponding author

7 **Keywords:** genome annotation, Gene Ontology, model organisms, food crops, bioenergy

8 crops, function annotation, genes of unknown function, Arabidopsis, sorghum, rice

## 9 Abstract

10 Since the entry into genome-enabled biology 20 years ago, much progress has been made in  
11 determining, describing, and disseminating functions of genes and their products. Yet, this  
12 information is still difficult to access by many, especially across genomes. To provide easy  
13 access to the status of genome function annotation for model organisms and bioenergy and  
14 food crop species, we created a web application (<https://genomeannotation.rheelab.org>) to  
15 visualize and download genome annotation data for 27 species. The summary graphics and  
16 data tables will be updated semi-annually and snapshots archived to provide a historical record  
17 of the progress of genome function annotation efforts.

## 18 Background

19 Rapid advances in DNA sequencing technologies made genome sequences widely available  
20 and revealed a plethora of genes encoded within the genomes in the last two decades [1]. The

21 timely invention and wide adoption of the Gene Ontology (GO) system transformed how gene  
22 and protein functions are described, quantified, and compared across many organisms [2,3].  
23 Despite this tremendous progress in genome biology, it is still nontrivial for scientists to get a  
24 snapshot of the status of genome function annotation across species.

25         There are several reasons for the difficulty in obtaining the status of genome function  
26 annotation across species. First, genome sequences and their annotations are hosted across  
27 multiple databases that use different gene/protein/sequence identifier systems. For example,  
28 Phytozome [4] uses its own database identifiers for its genes and does not provide cross-  
29 database identifier (ID) mapping functionalities. Although some databases include cross  
30 database references and provide tools to map IDs, such as UniProt's Retrieve/ID mapping and  
31 BioMart's ID conversion [5], these tools are not available for all sequenced genomes. Second,  
32 gene function information is not generally annotated using the GO system in the literature and  
33 databases. Third, genome function annotation databases generally only include annotated  
34 genes and it is not trivial to retrieve the number and identity of unannotated genes. Importance  
35 of unannotated genes is exemplified by a recent success in identifying the minimal bacterial  
36 genome that included 473 essential genes [6]. Among these were 149 whose molecular  
37 functions remain unknown.

38         To provide scientists and students an easy way to access the status of genome function  
39 annotations of model species and bioenergy and food crops, we created a web application that  
40 displays these data graphically and tabularly. The website retrieves data from multiple  
41 databases, and generates plots that show the percentages of genes with experimental,  
42 computational, or no annotations. The snapshots are updated semi-annually and past  
43 snapshots will be archived.

## 44 Results and Discussion

45 To represent the status of genome function annotation, we selected three groups of organisms:  
46 model organisms, bioenergy model and crop species, and most annotated plant species  
47 (**Figure 1**). Model organisms are important experimental tools for investigating biological  
48 processes and represent key reference points of biological knowledge for other species [7–9].  
49 This panel includes: *Arabidopsis thaliana*, *Caenorhabditis elegans*, *Danio rerio*, *Drosophila*  
50 *melanogaster*, *Mus musculus*, and *Saccharomyces cerevisiae* (**Fig. 1A**). We also included  
51 *Homo sapiens*, a species for which many model organisms are studied. Next, we selected  
52 bioenergy models and crops, which are important in expanding the renewable energy sector  
53 needed to combat the climate crisis and steward a more sustainable environment. Biomass is  
54 currently the biggest source of renewable energy [10] and is projected to become the biggest  
55 source of primary energy by 2050 [11]. The bioenergy models and crops we selected include:  
56 *Brachypodium distachyon*, *Chlamydomonas reinhardtii*, *Glycine max*, *Miscanthus sinensis*,  
57 *Panicum hallii*, *Panicum virgatum*, *Physcomitrium patens*, *Populus trichocarpa*, *Sorghum*  
58 *bicolor*, and *Setaria italica* (**Fig. 1B**). Finally, we selected ten additional plant species that have  
59 the most number of GO annotations in UniProt [12], which include: *Oryza sativa Japonica Group*  
60 (rice), *Gossypium hirsutum* (cotton), *Spinacia oleracea* (spinach), *Zea mays* (corn), *Medicago*  
61 *truncatula*, *Solanum tuberosum* (potato), *Ricinus communis* (castor bean), *Nicotiana tabacum*  
62 (tobacco), *Papaver somniferum* (opium poppy), *Triticum aestivum* (wheat) (**Fig. 1C**). These  
63 include the world's most important cereal crops, such as corn, rice, wheat, and vegetable crops  
64 such as potato [13].

65 There are several ways of accessing the status of genome function annotation for the 27  
66 species. From the front page, visitors can get a quick summary of the state of the genome  
67 function annotation as pie charts for the three groups of species (**Figure 1**). These pie charts  
68 show the percentage of genes that have: 1) annotations with experimental evidence (green); 2)

69 only the annotations that are computationally generated (blue); or 3) no annotations or  
70 annotations as being unknown (**Figure 1**). Of the 7 selected model organisms, *S. cerevisiae* has  
71 the highest percentage of genes with experimental evidence and the least number of genes  
72 unannotated or annotated as having unknown function, followed by *H. sapiens* and *A. thaliana*.  
73 Among the model organisms, *C. elegans* is the least known species with the greatest number of  
74 genes unannotated or annotated as having unknown function. Most of the plant species have  
75 few GO annotations based on experimental support to be even visible in the pie charts. Visitors  
76 can get more detailed information of any of the species by clicking on the species name below  
77 the pie charts. Each species page shows additional information about the annotation status,  
78 including displaying the portion of genes annotated to at least one GO domain (molecular  
79 function, cellular component, and biological process [2,3]) as well as a Venn diagram showing  
80 the overlap of genes annotated to more than one GO domain (**Figure 2**). This page also has  
81 links to source data and a tabular format of the annotation summary for browsing and  
82 downloading.

83 In developing our web application, we came across a few hurdles. First, there was not a  
84 single site where all data were available. To obtain GO annotations from the 27 species, we had  
85 to visit at least three databases. A positive finding was that all sites that had GO annotations  
86 were using the GO Annotation File (GAF) format. Nevertheless, having a single-entry point  
87 where GO annotations of any species can be accessed would be useful. Second, our website  
88 includes genes that are unannotated, which is often missing in gene function annotations and  
89 enrichment analyses [14]. Currently, extracting genes that are not annotated is not trivial and  
90 requires many steps that are different for each species. Including the unannotated genes in a  
91 genome into GAF files would facilitate many downstream applications.

92 To our surprise, some plant species with well-maintained, species-specific databases  
93 seem to have a low number of experimentally supported GO annotated genes in UniProt.

94 Outside of TAIR that provides GO annotations for *A. thaliana* [19], we were not able to find any  
95 database that provides experimental evidence codes to their GO annotations. Apart from  
96 *Nicotiana tabacum* and *Papaver somniferum*, all plants species on our website are included in  
97 the most recent version of Phytozome V13, but their GO terms are assigned computationally [4].  
98 The Sol Genomics Network (SGN) (<https://solgenomics.net> accessed 13 June 2022) [15] hosts  
99 genome annotations of Solanaceae species, including *Nicotiana tabacum* and *Solanum*  
100 *tuberosum*. An annotation file for *Nicotiana tabacum* is available [16] but they are assigned with  
101 computational support coming from InterProScan [17]. SpudDB [18] (<http://spuddb.uga.edu/>  
102 accessed 13 June 2022) provides GO annotation for *Solanum tuberosum* but they are  
103 generated with InterProScan and by best hit to the Arabidopsis proteome (TAIR10) [19].  
104 MaizeGDB [20] (<https://www.maizegdb.org/> accessed 13 June 2022) provides GO annotation  
105 for *Zea mays* that are assigned with GO annotation tools including Argot2.5, FANN-GO, and  
106 PANNZER [21], which are all computational annotations. SpinachBase  
107 (<http://www.spinachbase.org/> accessed 13 June 2022) provides a centralized access to  
108 *Spinacia oleracea*, and their GO annotations are generated computationally with Blast2GO [22].  
109 *Oryza sativa Japonica Group* GO annotations can be found on Rice Genome Annotation Project  
110 [23] and they are assigned with BLASTP searches against Arabidopsis GO-curated proteins  
111 [24]. Gramene [31] (<https://www.gramene.org/> accessed 13 June 2022) hosts genome data for  
112 many species but we could not find GO annotations with evidence codes. We were not able to  
113 find species-specific databases that provide GO annotations for *Triticum aestivum*, *Gossypium*  
114 *hirsutu*, *Medicago truncatula*, *Papaver somniferum* or *Ricinus communis*. In summary, most  
115 plant genome databases stop at computationally generating GO annotations and some  
116 important species do not appear to have dedicated databases. More efforts are needed in both  
117 experimentally validating functional annotations made from computational approaches and  
118 curating experimentally supported function descriptions in the literature into structured  
119 annotations such as GO, which will be crucial for accelerating gene function discovery.

## 120 Conclusions

121 Our website provides a convenient way to obtain the current state of genome function  
122 annotation for model organisms and crops for bioenergy, food, and medicine. Our website  
123 shows how much is annotated and unannotated in the 27 species that represent some of the  
124 most intensely studied and arguably the most valuable organisms for science and society. By  
125 proxy, these charts illustrate how much is known and unknown. These snapshots will be  
126 updated on a semi-annual basis, and comparing the charts across time will reflect how  
127 biological knowledge changes over time. These snapshots can be useful in many contexts  
128 including research projects, grant proposals, review articles, annual reports, and outreach  
129 materials. The data summarized on this website can be linked to their sources, which can be  
130 used for a variety of investigations. Successful examples include exploring why certain proteins  
131 remain unannotated [25], developing pipelines to infer function without relying on sequence  
132 similarity [26], and assessing annotation coverage across bacterial proteomes [27]. As our  
133 society transitions into biology-enabled manufacturing [32], fundamental knowledge of how  
134 genes and their products function at various scales will be crucial in ushering in the era of bio-  
135 economy.

## 136 Methods

### 137 Selecting species and data retrieval

138 For the seven model organisms, gene function annotations were downloaded as GO Annotation  
139 Files (GAF files) from the GO consortium website  
140 (<http://current.geneontology.org/products/pages/downloads.html> accessed 13 June 2022) of the  
141 2022-05-16 release. Genes found in a genome were retrieved from the source indicated on the

142 GO annotation download page as General Feature Format (GFF) files. A detailed description of  
143 the files used to generate charts on our website, including data for the other category of  
144 species, can be found in **Table S1**.

145 Genome annotation and gene list for bioenergy models and crops were downloaded  
146 from Phytozome version V13 (<https://data.jgi.doe.gov/refine-download/phytozome> accessed 13  
147 June 2022). Although some species in this category had GO annotations in the GO consortium  
148 database, the sequence identifiers (IDs) for genes could not easily be mapped to Phytozome  
149 IDs. To maintain consistency within this category, all annotation files were downloaded from  
150 Phytozome. All Phytozome GO annotations are computationally generated [4]. Gene lists were  
151 also retrieved from Phytozome V13.

152 For the last category of plant species, we selected the most annotated plant species  
153 from the UniProt GO annotation database [28] GAF files hosted on the GO consortium website  
154 (<http://current.geneontology.org/products/pages/downloads.html> 2022-05-16 release, accessed  
155 13 June 2022). We downloaded these species reference proteomes from the UniProt release  
156 2022\_02 and retrieved the number of corresponding genes.

157 Using the evidence codes provided by GAF files, we generated the numbers of genes  
158 annotated with GO supported by experimental evidence. If a gene has at least one GO term  
159 annotated using any of the following codes: EXP (Inferred from Experiment), IDA (Inferred from  
160 Direct Assay), IPI (Inferred from Physical Interaction), IMP (Inferred from Mutant Phenotype),  
161 IGI (Inferred from Genetic Interaction), or IEP (Inferred from Expression Pattern), we  
162 categorized the gene as having “Experimental Evidence” for function. Genes that have at least  
163 one annotated GO term, but no terms have the evidence codes described above, are  
164 categorized as “Predicted”. Since Phytozome has only computationally generated GO  
165 annotations, all of their genes are categorized as having their functions “Predicted”. By  
166 subtracting the annotated genes from the total number of genes, we retrieved the number of

167 genes without any GO annotations. These numbers were used to generate pie charts to show  
168 the proportions of genes in each category for every species.

169 All files were processed with scripts written in Python (3.10). All pie charts were  
170 generated using Python Matplotlib version 3.5.2 and Venn diagrams were generated using  
171 Python matplotlib-venn version 0.11.7. The repository of codes can be found at GitHub  
172 (<https://github.com/bxuecarnegie/AnnotationStats>).

## 173 Creating the Website

174 To create a website for hosting our charts, we used Node.js [29] for our server-side  
175 environment, which provides the Application Program Interface (API) for the front end to retrieve  
176 the plots generated by Python. The front end of the website uses AngularJS [30].

## 177 Declarations

## 178 Ethics approval and consent to participate

179 Not applicable

## 180 Consent for publication

181 Not applicable

## 182 Availability of data and materials

183 Data used in this study are all publicly available. GO annotation files were downloaded from  
184 (<http://current.geneontology.org/annotations/index.html> 2022-05-16 release, accessed 13 June  
185 2022) and Phytozome (<https://data.jgi.doe.gov/refine-download/phytozome> V13, accessed 13  
186 June 2022). Gene data were downloaded from sources indicated on the GO  
187 (<http://current.geneontology.org/products/pages/downloads.html> accessed 13 June 2022),  
188 Phytozome, and UniProt (<https://www.uniprot.org/> accessed 13 June 2022). Supplemental  
189 Table S1 provides detailed information on all species annotation and gene source databases,  
190 downloaded versions, and URLs. Graphs and statistics data generated in this study are  
191 available at (<http://genomeannotation.rheelab.org/> accessed 13 June 2022). Scripts used to  
192 process the data and generate the graphs are written in Python 3 and are available at GitHub  
193 (<https://github.com/bxuecarnegie/AnnotationStats> accessed 13 June 2022).

## 194 Competing interests

195 The authors declare no competing interests.

## 196 Funding

197 This work was supported, in part, by the U.S. Department of Energy, Office of Science, Office of  
198 Biological and Environmental Research, Genomic Science Program grant nos. DE-SC0018277,  
199 DE-SC0008769, DE-SC0020366 and DE-SC0021286 and the U.S. National Science  
200 Foundation grants MCB-1617020 and IOS-1546838. This work was done on the ancestral land

201 of the Muwekma Ohlone Tribe, which was and continues to be of great importance to the  
202 Ohlone people.

## 203 Authors' contributions

204 SYR conceived the project and BX implemented the project. BX and SYR wrote and edited the  
205 manuscript.

## 206 Acknowledgements

207 We would like to thank the members of Rhee lab for the discussions and suggestions on the  
208 project.

## 209 References

- 210 1. O'Leary NA, Wright MW, Brister JR, Ciuffo S, Haddad D, McVeigh R, et al. Reference  
211 sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional  
212 annotation. *Nucleic Acids Res.* [academic.oup.com](http://academic.oup.com); 2016;44:D733–45.
- 213 2. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene Ontology: tool  
214 for the unification of biology. *Nat Genet.* Nature Publishing Group; 2000;25:25–9.
- 215 3. Acids research N, 2021. The Gene Ontology resource: enriching a GOLD mine. *Nucleic Acids*  
216 *Res.* Oxford University Press; 2021;49:D325–34.
- 217 4. Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, et al. Phytozome: a  
218 comparative platform for green plant genomics. *Nucleic Acids Res.* 2012;40:D1178–86.
- 219 5. Guberman JM, Ai J, Arnaiz O, Baran J, Blake A, Baldock R, et al. BioMart Central Portal: an  
220 open database network for the biological community. *Database* . [academic.oup.com](http://academic.oup.com);  
221 2011;2011:bar041.
- 222 6. Hutchison CA 3rd, Chuang R-Y, Noskov VN, Assad-Garcia N, Deerinck TJ, Ellisman MH, et  
223 al. Design and synthesis of a minimal bacterial genome. *Science.* 2016;351:aad6253.
- 224 7. Ankeny RA, Leonelli S. *Model Organisms. Elements in the Philosophy of Biology.* Cambridge  
225 University Press; 2020.

- 226 8. Fields S, Johnston M. Cell biology. Whither model organism research? *Science*.  
227 2005;307:1885–6.
- 228 9. Jones AM, Chory J, Dangl JL, Estelle M, Jacobsen SE, Meyerowitz EM, et al. The impact of  
229 *Arabidopsis* on human health: diversifying our portfolio. *Cell*. 2008;133:939–43.
- 230 10. U.S. energy facts explained - consumption and production - U.S. Energy Information  
231 Administration (EIA) [Internet]. [cited 2022 Jun 5]. Available from:  
232 <https://www.eia.gov/energyexplained/us-energy-facts/>
- 233 11. Reid WV, Ali MK, Field CB. The future of bioenergy. *Glob Chang Biol*. 2020;26:274–86.
- 234 12. UniProt Consortium. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res*.  
235 [academic.oup.com](http://academic.oup.com); 2019;47:D506–15.
- 236 13. FIGURE 21: World production of crops, main commodities [Internet]. FAO Statistical  
237 Yearbook 2021 Datasets. Food and Agriculture Organization of the United Nations; 2021.  
238 Available from: <http://www.fao.org/3/cb4477en/StatYearbook2021-fig21.xlsx>
- 239 14. Higgins DP, Weisman CM, Lui DS, D’Agostino FA, Walker AK. Defining characteristics and  
240 conservation of poorly annotated genes in *Caenorhabditis elegans* using WormCat 2.0.  
241 *Genetics* [Internet]. 2022; Available from: <http://dx.doi.org/10.1093/genetics/iyac085>
- 242 15. Fernandez-Pozo N, Menda N, Edwards JD, Saha S, Teclé IY, Strickler SR, et al. The Sol  
243 Genomics Network (SGN)--from genotype to phenotype to breeding. *Nucleic Acids Res*.  
244 2015;43:D1036–41.
- 245 16. Edwards KD, Fernandez-Pozo N, Drake-Stowe K, Humphry M, Evans AD, Bombarely A, et  
246 al. A reference genome for *Nicotiana tabacum* enables map-based cloning of homeologous loci  
247 implicated in nitrogen utilization efficiency. *BMC Genomics*. Springer; 2017;18:448.
- 248 17. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-  
249 scale protein function classification. *Bioinformatics*. 2014;30:1236–40.
- 250 18. Hirsch CD, Hamilton JP, Childs KL, Cepela J, Crisovan E, Vaillancourt B, et al. Spud DB: A  
251 resource for mining sequences, genotypes, and phenotypes to accelerate potato breeding.  
252 *Plant Genome*. Wiley; 2014;7:lantgenome2013.12.0042.
- 253 19. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, et al. The *Arabidopsis*  
254 Information Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Res*.  
255 [academic.oup.com](http://academic.oup.com); 2012;40:D1202–10.
- 256 20. Woodhouse MR, Cannon EK, Portwood JL 2nd, Harper LC, Gardiner JM, Schaeffer ML, et  
257 al. A pan-genomic approach to genome databases using maize as a model system. *BMC Plant*  
258 *Biol*. 2021;21:385.
- 259 21. Wimalanathan K, Friedberg I, Andorf CM, Lawrence-Dill CJ. Maize GO Annotation-Methods,  
260 Evaluation, and Review (maize-GAMER). *Plant Direct*. 2018;2:e00052.
- 261 22. Collins K, Zhao K, Jiao C, Xu C, Cai X, Wang X, et al. SpinachBase: a central portal for  
262 spinach genomics. Database [Internet]. [academic.oup.com](http://academic.oup.com); 2019;2019. Available from:  
263 <http://dx.doi.org/10.1093/database/baz072>

- 264 23. Ouyang S, Zhu W, Hamilton J, Lin H, Campbell M, Childs K, et al. The TIGR Rice Genome  
265 Annotation Resource: improvements and new features. *Nucleic Acids Res.* 2007;35:D883–7.
- 266 24. Yuan Q, Ouyang S, Wang A, Zhu W, Maiti R, Lin H, et al. The institute for genomic research  
267 Osa1 rice genome annotation database. *Plant Physiol.* 2005;138:18–26.
- 268 25. Wood V, Lock A, Harris MA, Rutherford K, Bähler J, Oliver SG. Hidden in plain sight: what  
269 remains to be discovered in the eukaryotic proteome? *Open Biol. The Royal Society*;  
270 2019;9:180241.
- 271 26. Bossi F, Fan J, Xiao J, Chandra L, Shen M, Dorone Y, et al. Systematic discovery of novel  
272 eukaryotic transcriptional regulators using sequence homology independent prediction. *BMC*  
273 *Genomics.* Springer; 2017;18:480.
- 274 27. Lobb B, Tremblay BJ-M, Moreno-Hagelsieb G, Doxey AC. An assessment of genome  
275 annotation coverage across the bacterial tree of life. *Microb Genom [Internet]*. 2020;6. Available  
276 from: <http://dx.doi.org/10.1099/mgen.0.000341>
- 277 28. Camon E, Magrane M, Barrell D, Lee V, Dimmer E, Maslen J, et al. The Gene Ontology  
278 Annotation (GOA) Database: sharing knowledge in UniProt with Gene Ontology. *Nucleic Acids*  
279 *Res.* 2004;32:D262–6.
- 280 29. Tilkov S, Vinoski S. Node.js: Using JavaScript to Build High-Performance Network  
281 Programs. *IEEE Internet Comput.* 2010;14:80–3.
- 282 30. Jain, Bhansali, Mehta. AngularJS: A modern MVC framework in JavaScript. *Journal of*  
283 *Global Research in Computer [Internet]*. [jgrcs.info](http://www.jgrcs.info); 2014; Available from:  
284 <http://www.jgrcs.info/index.php/jgrcs/article/download/952/610>
- 285 31. Tello-Ruiz MK, Naithani S, Gupta P, Olson A, Wei S, Preece J, et al. Gramene 2021:  
286 harnessing the power of comparative genomics and pathways for plant research. *Nucleic Acids*  
287 *Res.* 2021;49:D1452–63.
- 288 32. Committee on Industrialization of Biology: A Roadmap to Accelerate the Advanced  
289 Manufacturing of Chemicals, Board on Chemical Sciences and Technology, Board on Life  
290 Sciences, Division on Earth and Life Studies, National Research Council. *Industrialization of*  
291 *Biology: A Roadmap to Accelerate the Advanced Manufacturing of Chemicals.* Washington  
292 (DC): National Academies Press (US); 2015.

293

## 294 Figures

295 **Figure 1** Status of genome function annotations.

296 Each pie chart shows the proportion of genes that are annotated to a domain of Gene Ontology  
297 (GO): molecular function, biological process, or cellular component. Green indicates genes that

298 have at least one experimentally validated GO annotation, blue indicates genes that are  
299 annotated but none are experimentally annotated, and gray indicates genes that do not have  
300 any GO annotations. The species are sorted by the percentage of genes with experimental  
301 evidence. A) selected model organisms; B) bioenergy models and crops [1]; C) other plant  
302 species with the highest percentage of genes with experimental evidence in UniProt.

303

304 **Figure 2** An example species-specific annotation web page shown for *Arabidopsis thaliana*. It  
305 consists of 3 parts: 1) a table that consists of data sources; 2) pie charts showing the proportion  
306 of each type of genes; and 3) a table showing the numbers of genes in each category, which  
307 can be toggled to show/hide.

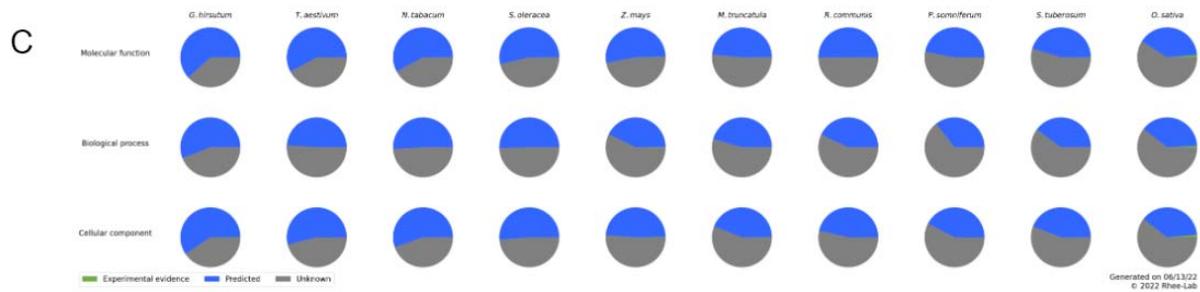
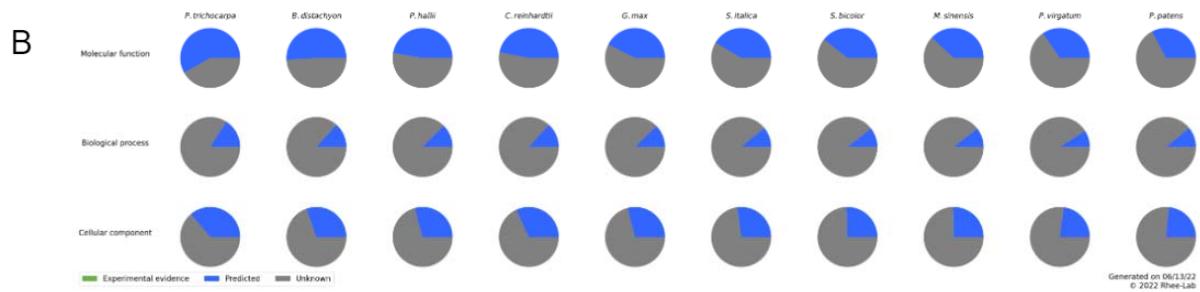
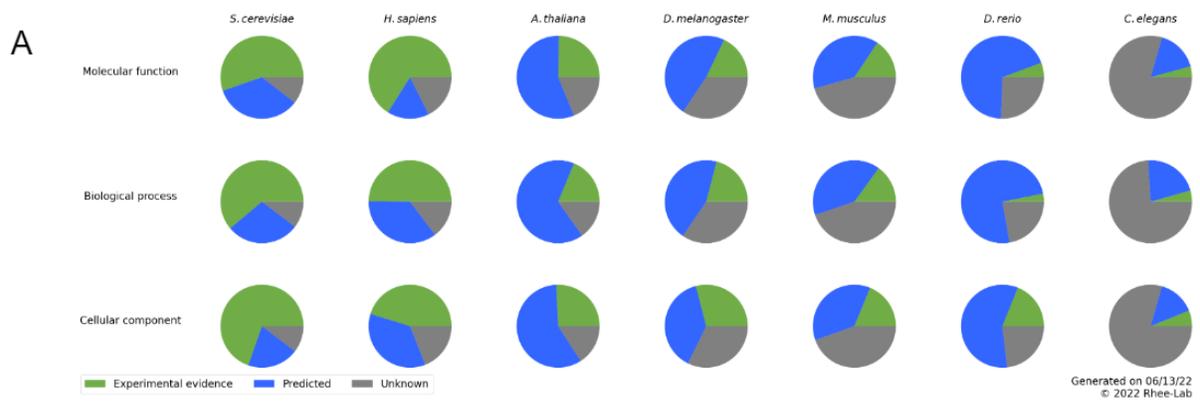
## 308 Supplemental Information

### 309 Additional file 1

310 **Table S1** Source data. A table describing the data sources, versions downloaded, and URLs

311

312

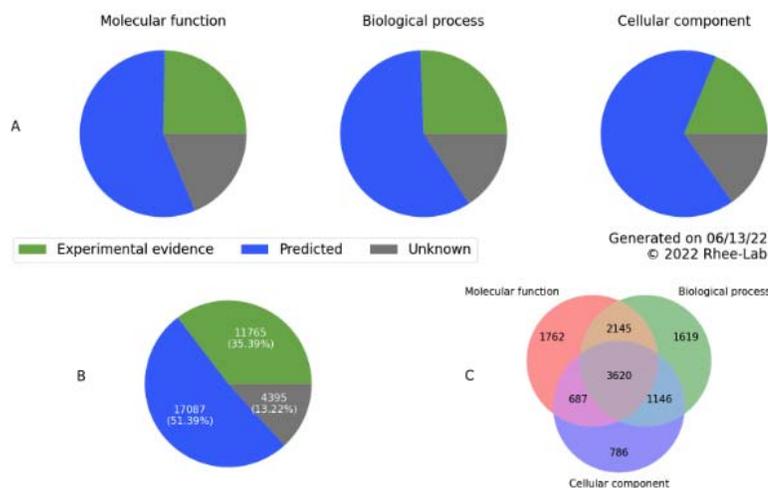


## Arabidopsis thaliana

GAF Source	<a href="#">Gene Ontology</a>
Sequence Source	<a href="#">The Arabidopsis Information Resource (TAIR)</a>
Download Date	6/13/22

[Back to Overview](#)

### Status of gene function elucidation and annotation



[Copy Link to Image](#) [Download "Experimental Evidence" Gene List](#) [Download "Predicted" Gene List](#) [Download "Unknown" Gene List](#)

### Stats of Arabidopsis thaliana (Taxon: 3702)

[Toggle Table](#) [Download Table](#)

<b>Num. of Genes</b>	33247
<b>Num. of Genes w/ Experimental evidence</b>	11775
<b>Num. of Genes w/ Computational evidence</b>	17077
<b>Num. of Genes with no GO annotations</b>	4395
<b>Num. of Molecular function (Experimental evidence)</b>	8217 (24.72%)
<b>Num. of Molecular function (Predicted)</b>	18802 (56.55%)
<b>Num. of Molecular function (Unknown)</b>	6228 (18.73%)
<b>Num. of Biological process (Experimental evidence)</b>	8558 (25.74%)
<b>Num. of Biological process (Predicted)</b>	19433 (58.45%)
<b>Num. of Biological process (Unknown)</b>	5256 (15.81%)
<b>Num. of Cellular component (Experimental evidence)</b>	6241 (18.77%)
<b>Num. of Cellular component (Predicted)</b>	21949 (66.02%)
<b>Num. of Cellular component (Unknown)</b>	5057 (15.21%)

Generation Date: 6/13/22

### Archived charts

[Link to archives.](#)