

Status of genome function annotation in model organisms and crops

Bo Xue  | Seung Y. Rhee 

Department of Plant Biology, Carnegie Institution for Science, Stanford, California, USA

Correspondence

Seung Y. Rhee, Department of Plant Biology, Carnegie Institution for Science, Stanford, CA 94305, USA.

Email: srhee@carnegiescience.edu

Present addresses

Bo Xue and Seung Y. Rhee, Plant Resilience Institute, Michigan State University, East Lansing, MI 4882.

Funding information

National Science Foundation (NSF), Grant/Award Numbers: MCB-2052590, MCB-1916797, IOS-1546838; U.S. Department of Energy (DOE), Grant/Award Numbers: DE-SC0018277, DE-SC0020366, DE-SC0023160, DE-SC0021286

Abstract

Since the entry into genome-enabled biology several decades ago, much progress has been made in determining, describing, and disseminating the functions of genes and their products. Yet, this information is still difficult to access for many scientists and for most genomes. To provide easy access and a graphical summary of the status of genome function annotation for model organisms and bioenergy and food crop species, we created a web application (<https://genomeannotation.rheelab.org>) to visualize, search, and download genome annotation data for 28 species. The summary graphics and data tables will be updated semi-annually, and snapshots will be archived to provide a historical record of the progress of genome function annotation efforts. Clear and simple visualization of up-to-date genome function annotation status, including the extent of what is unknown, will help address the grand challenge of elucidating the functions of all genes in organisms.

KEYWORDS

bioenergy crops, food crops, gene function, Gene Ontology, genome annotation, model organisms

1 | INTRODUCTION

Rapid advances in DNA sequencing technologies made genome sequences widely available and revealed a plethora of genes encoded within the genomes (O'Leary et al., 2016). The timely invention and wide adoption of the Gene Ontology (GO) system transformed how gene and protein functions are described, quantified, and compared across many organisms (Ashburner et al., 2000; The Gene Ontology Consortium, 2021). A grand challenge in life sciences is to elucidate the functions of all the genes that have been identified through genome sequencing. One of the first steps in elucidating gene function systematically is to know which genes have unknown functions. Yet, a snapshot of the status of genome function annotation across species is still not readily available to scientists.

The status of genome function annotation is not easily accessible for several reasons. First, genome sequences and their annotations

are hosted across multiple databases that use different gene/protein/sequence identifier (ID) systems. Although some databases include cross-database references and provide tools to map IDs, such as UniProt's Retrieve/ID mapping and BioMart's ID conversion (Guberman et al., 2011), these tools are not available for all sequenced genomes. Second, gene function information is not generally annotated using the GO annotation and evidence code system in the literature and most databases. Third, genome function annotation databases generally only include annotated genes, and it is often not straightforward to identify unannotated genes. In enrichment analysis, unannotated genes play an important role in reducing ascertainment bias and facilitating the discovery of novel processes.

To provide scientists and students with an easy way to access and visualize the status of genome function annotations of model species and bioenergy and food crops, we created a web application that displays these data graphically and tabularly and allows searching and

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial](https://creativecommons.org/licenses/by-nc/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited and is not used for commercial purposes.

© 2023 The Authors. *Plant Direct* published by American Society of Plant Biologists and the Society for Experimental Biology and John Wiley & Sons Ltd.

downloading annotations using GO term IDs. The website retrieves data from multiple databases and generates plots that show the percentages of genes with experimental, computational, or no annotations. The snapshots are updated semi-annually, and past snapshots will be archived.

2 | RESULTS

To represent the status of genome function annotation, we selected three groups of organisms: model organisms, bioenergy model and crop species, and most-annotated plant species (Figure 1). Model organisms are important experimental tools for investigating

biological processes and represent key reference points of biological knowledge for other species (Ankeny & Leonelli, 2020; Fields & Johnston, 2005; A. M. Jones et al., 2008). This panel includes *Arabidopsis thaliana*, *Caenorhabditis elegans*, *Danio rerio*, *Drosophila melanogaster*, *Mus musculus*, *Saccharomyces cerevisiae*, and *Schizosaccharomyces pombe* (Figure 1a). We also included *Homo sapiens*, a species for which many model organisms are studied. Next, we selected bioenergy models and crops, which are important in expanding the renewable energy sector needed to combat the climate crisis and steward a more sustainable environment. For example, biomass is projected to become the biggest source of primary energy by 2050 (Reid et al., 2020). The bioenergy models and crops we selected include *Brachypodium distachyon*, *Chlamydomonas reinhardtii*,

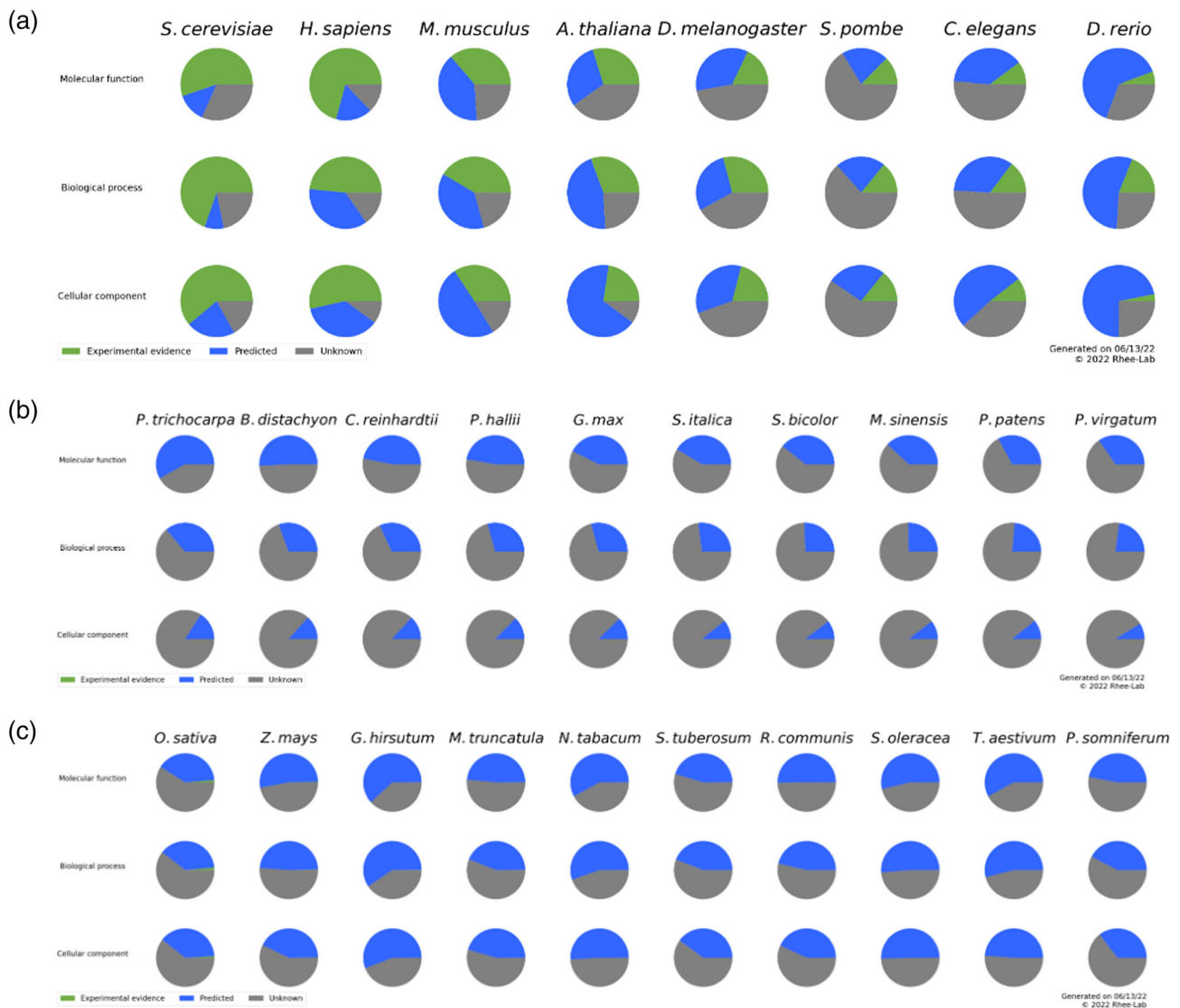


FIGURE 1 Status of genome function annotations. Each pie chart shows the proportion of genes that are annotated to a domain of Gene Ontology (GO): molecular function, biological process, or cellular component. Green indicates genes that have at least one experimentally validated GO annotation; blue indicates genes that are annotated with computationally predicted GO annotations; and gray indicates genes that do not have any GO annotations or have GO annotations with ND (no biological data available). The species are ordered by the average percentage of genes with experimental evidence (or genes with computationally predicted annotations if experimental evidence is not available) of all three GO domains. (a) Selected model organisms; (b) bioenergy models and crops; (c) other plant species with the highest percentage of genes with experimental evidence in UniProt.



Glycine max, *Miscanthus sinensis*, *Panicum hallii*, *Panicum virgatum*, *Physcomitrium patens*, *Populus trichocarpa*, *Sorghum bicolor*, and *Setaria italica* (Figure 1b). Finally, we included 10 additional plant species that have the highest number of GO annotations in UniProt (UniProt Consortium, 2019), which include: *Oryza sativa Japonica Group* (rice), *Gossypium hirsutum* (cotton), *Spinacia oleracea* (spinach), *Zea mays* (corn), *Medicago truncatula*, *Solanum tuberosum* (potato), *Ricinus communis* (castor bean), *Nicotiana tabacum* (tobacco), *Papaver somniferum* (opium poppy), and *Triticum aestivum* (wheat) (Figure 1c). These include the world's most important cereal crops, such as corn, rice, and wheat, and vegetable crops, such as potatoes (Food and Agriculture Organization of the United Nations, 2021).

There are several ways of accessing the status of genome function annotation for the 28 species on our [web application](#). From the front page, visitors can get a quick summary of the status of the genome function annotation as pie charts for the three groups of species (Figure 1). These pie charts show the percentage of genes that have: (1) annotations with experimental evidence (green); (2) only the annotations that are computationally generated (blue); or (3) no annotations or annotations as being unknown (gray) (Figure 1). Of the eight selected model organisms, *H. sapiens* has the highest percentage of genes whose functions are annotated with experimental evidence, followed by *S. cerevisiae* and *M. musculus*. *A. thaliana* has the lowest percentage of "unknown" genes. Among the model organisms, *C. elegans* is the least known species, with the greatest number of genes unannotated or annotated as having an unknown function. Most of the plant species have too few GO annotations based on experimental support to be visible in the pie charts. Visitors can get more detailed information about any of the species by clicking on the species name below the pie charts in the "Species details" sections. Each species page shows additional information about annotation status, including displaying the proportion of genes annotated to at least one GO domain (molecular function, cellular component, and biological process) (Ashburner et al., 2000; The Gene Ontology Consortium, 2021), as well as a Venn diagram showing the overlap of genes annotated to more than one GO domain (Figure 2). This page also has links to source data and a tabular format of the annotation summary for browsing and downloading. Users can also download lists of annotated genes directly from the summary table or search and download annotations using GO term IDs.

3 | DISCUSSION

Our website provides a convenient way to obtain and visualize the current state of genome function annotation for model organisms and crops for bioenergy, food, and medicine. These charts also serve as a proxy for illustrating how much is known and unknown. These snapshots will be updated on a semi-annual basis, and comparing the charts across time will reflect how biological knowledge changes over time. These snapshots can be useful in many contexts, including research projects, grant proposals, review articles, annual reports, and outreach materials.

In developing our web application, we encountered a few hurdles. First, there was not a single site where all the data were available. To obtain GO annotations for the 28 species, we had to visit several databases. An encouraging finding was that all sites that had GO annotations were using the GO Annotation File (GAF) format. Several tools and databases currently provide a single entry point for searching and retrieving GO annotations for multiple species. AmiGO by the GO consortium (Carbon et al., 2009) and QuickGO from the Gene Ontology Annotation (GOA) project (Binns et al., 2009) provide search functions for taxon-specific GO annotations. However, they require users to set up additional queries to export GAF files for species of interest. BioMart from the Ensembl project (Kinsella et al., 2011) is another tool that users can query GO annotations for multiple species using evidence codes. But their output fields do not include publications or references that are linked to the supporting evidence and require users to manually select output fields to conform to the GAF format. Second, our website includes genes that are unannotated, which are often missing in gene function annotations and enrichment analyses (Higgins et al., 2022). Currently, extracting genes that are not annotated requires many steps that differ across species. Including the unannotated genes in a genome in GAF files would facilitate many downstream applications. Third, visualizations of the status of genome annotations do not yet exist, which our website provides using pie charts and tabular summaries.

To our surprise, some plant species with well-maintained, species-specific databases seem to have a low number of experimentally supported GO-annotated genes in UniProt. The Arabidopsis Information Resource (TAIR) (Lamesch et al., 2012) and Sol Genomics Network (SGN) (Fernandez-Pozo et al., 2015) are the only two taxon-specific plant databases that provide GAF files with experimental evidence codes. Most plant genome databases stop at computationally generating GO annotations, and some well-studied species do not appear to have dedicated databases. Apart from *N. tabacum* and *P. somniferum*, all plant species on our website are included in the most recent version of Phytozome version 13, though their GO terms are assigned only computationally (Goodstein et al., 2012). The SGN (Fernandez-Pozo et al., 2015) hosts genome annotations for Solanaceae species, including *N. tabacum* and *S. tuberosum*. An annotation file for *N. tabacum* is available (Edwards et al., 2017), which is assigned with computational support coming from InterProScan (P. Jones et al., 2014). Spud DB (Hirsch et al., 2014) provides GO annotations for *S. tuberosum*, but they are generated with InterProScan and by finding best hits to the Arabidopsis proteome (TAIR10) (Lamesch et al., 2012). MaizeGDB (Woodhouse et al., 2021) provides GO annotation for *Z. mays* that is assigned with GO annotation tools including Argot2.5, FANN-GO, and PANNZER (Wimalanathan et al., 2018), which are all computational annotations. SpinachBase provides centralized access to *S. oleracea* genome, and their GO annotations are generated computationally with Blast2GO (Collins et al., 2019). *O. sativa Japonica Group* GO annotations can be found on the Rice Genome Annotation Project (Ouyang et al., 2007), which are assigned with Protein Basic Local Alignment Search Tool (BLASTP) searches against Arabidopsis GO-curated proteins (Yuan et al., 2005). Gramene

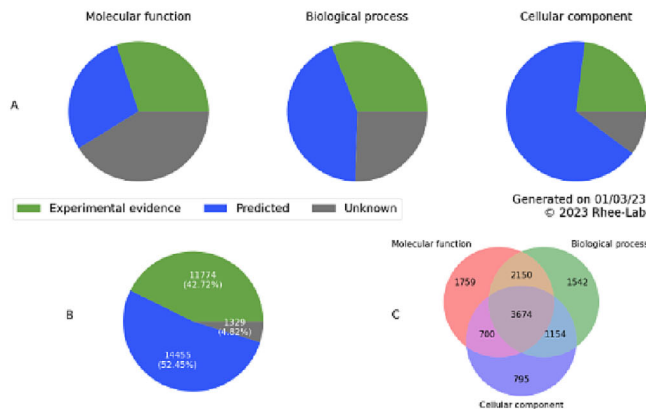
A. thaliana
Graphic Summary
Tabular Status
Search and Download

Arabidopsis thaliana

GAF Source [Gene Ontology](#)
Sequence Source [The Arabidopsis Information Resource \(TAIR\)](#)
Download Date 01/03/23

[Back to Overview](#) [Copy Link to Image](#)

Status of gene function elucidation and annotation



[Download "Experimental Evidence" Annotations](#) [Download "Predicted" Annotations](#) [Download "Unknown" Gene List](#)

Stats of *A.thaliana* (Taxon: 3702)

[Show/Hide Table](#) [Download Table](#)

Num. of Genes	27558
Num. of Genes w/ Experimental evidence	11774
Num. of Genes w/ Computational evidence	14455
Num. of Genes with no GO annotations	1329
Num. of Molecular function (Experimental evidence)	8283 (30.06%)
Num. of Molecular function (Predicted)	7918 (28.73%)
Num. of Molecular function (Unknown)	11357 (41.21%)
Num. of Biological process (Experimental evidence)	8520 (30.92%)
Num. of Biological process (Predicted)	12048 (43.72%)
Num. of Biological process (Unknown)	6990 (25.36%)
Num. of Cellular component (Experimental evidence)	6323 (22.94%)
Num. of Cellular component (Predicted)	18438 (66.91%)
Num. of Cellular component (Unknown)	2797 (10.15%)

Generation Date: 01/03/23

Search and Download Gene Annotations

[Show/Hide Form](#)

GO ID(s)

Input GO IDs, e.g. GO:0005515 (Separate each ID with whitespace or comma, only IDs that are in "GO:XX.X" format will be submitted).

Evidence Type

Choose...
Experimental Evidence
Predicted

Select GO annotation evidence type(s) (To select multiple items in a list, hold down the Ctrl (PC) or Command (Mac) key).

Aspect

Choose...
Molecular Function
Biological Process
Cellular Component

Select GO annotation aspect(s) (To select multiple items in a list, hold down the Ctrl (PC) or Command (Mac) key).

[Download](#)

FIGURE 2 An example species-specific annotation web page shown for *Arabidopsis thaliana*. It consists of three parts: (1) a table comprising data sources; (2) pie charts showing the proportion of each type of annotations; (3) a table showing the numbers of genes in each category, which users can toggle to either show or hide; and (4) a form to search and download annotations using Gene Ontology identifiers (GO IDs).



(Tello-Ruiz et al., 2021) hosts genome data for many species, but we could not find GO annotations with evidence codes. We were not able to find species-specific databases that provide GO annotations for *T. aestivum*, *G. hirsutum*, *M. truncatula*, *P. somniferum*, or *R. communis*. PLAZA (Van Bel et al., 2022) is a platform that integrates structural and functional genome annotation for many plant species. They do incorporate annotations supported by experimental evidence from GO (Acids research & 2021, 2021) and GOA (Huntley et al., 2015), and additional annotations are computationally generated using InterProScan or assigned with empirically validated GO annotations from orthologs. More efforts are needed in experimentally validating functional annotations made from computational approaches, characterizing genes of unknown function, and curating experimentally supported function descriptions in the literature into structured annotations such as GO annotations, which will be crucial for accelerating gene function discovery.

The data summarized on this web application can be linked to their sources, which can be used for a variety of investigations. Successful examples include exploring why certain proteins remain unannotated (Wood et al., 2019), developing pipelines to infer function without relying on sequence similarity (Bossi et al., 2017), and assessing annotation coverage across bacterial proteomes (Lobb et al., 2020). One of the biggest challenges in life sciences today is the limited understanding of what most genes encoded in genomes do. One of the first steps in the systemic elucidation of gene function is to know and easily access the set of genes with annotations as well as those without any annotations. Our website provides this functionality for 28 of the most intensely studied species. As our society transitions into biology-enabled manufacturing (National Research Council, 2015), fundamental knowledge of how genes and their products function at various scales will be crucial, as our society transitions into a new bio-economy era.

4 | METHODS

4.1 | Selecting species and data retrieval

For the model organisms, gene function annotations were downloaded as GAF files from the GO consortium website (May 16, 2022 release). For *S. pombe*, data were downloaded directly from PomBase (Harris et al., 2022). Genes found in a genome were retrieved as General Feature Format (GFF) files from the source indicated on the GO annotation download page. If information about the type of genes was provided in the GFF file, we only included the “protein-encoding” genes. A detailed description of the files used to generate charts on our website, including data for the other categories of species, can be found in Data S1.

Genome annotations and gene lists for bioenergy models and crops were downloaded from Phytozome version 13. Although some species in this category had GO annotations in the GO consortium database, the sequence IDs for genes could not easily be mapped to Phytozome IDs. To maintain consistency within this

category, all annotation files were downloaded from Phytozome. All Phytozome GO annotations are computationally generated (Goodstein et al., 2012). Gene lists were also retrieved from Phytozome version 13.

For the last category of plant species, we selected the most annotated plant species from the UniProt GO annotation database (Camon et al., 2004) and GAF files hosted on the GO consortium website. We downloaded these species reference proteomes from the UniProt release 2022_2 and retrieved the number of corresponding genes.

Using the evidence codes provided by GAF files, we generated the numbers of genes annotated with GO supported by experimental evidence. If a gene has at least one GO term annotated using any of the following evidence codes: EXP (Inferred from Experiment), IDA (Inferred from Direct Assay), IPI (Inferred from Physical Interaction), IMP (Inferred from Mutant Phenotype), IGI (Inferred from Genetic Interaction), or IEP (Inferred from Expression Pattern), we categorized the gene as having “Experimental Evidence” for function. Genes that have at least one annotated GO term but no terms that have the evidence codes described above are categorized as “Predicted.” Because Phytozome has only computationally generated GO annotations, all of their genes are categorized as having their functions “Predicted.” By subtracting the annotated genes from the total number of genes, we retrieved the number of genes without any GO annotations, which were included in the “Unknown” category. Finally, GO annotations with the ND (no biological data available) were moved to the “Unknown” category. These numbers were used to generate pie charts to show the proportions of genes in each function annotation category for each species.

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

4.2 | Creating the web application

To create a web application for hosting our charts, we used Node.js (Tilkov & Vinoski, 2010) for our server-side environment, which provides the application program interface (API) for the front end to retrieve the plots generated by Python. The front end of the website uses AngularJS (Jain et al., 2014).

AUTHOR CONTRIBUTIONS

Seung Y Rhee conceived the project and Bo Xue implemented it. Bo Xue and Seung Y Rhee wrote and edited the manuscript.

ACKNOWLEDGMENTS

We thank members of the Rhee Lab for discussions and suggestions on the project and Kristen Yawitz for editing the manuscript. This work was supported, in part, by the U.S. Department of Energy, Office of Science, Office of Biological and Environmental

Research, Genomic Science Program grant nos. DE-SC0018277, DE-SC0020366, DE-SC0023160, and DE-SC0021286, and the National Science Foundation grants MCB-2052590, MCB-1916797, and IOS-1546838. This work was done on the ancestral land of the Muwekma Ohlone Tribe, which was and continues to be of great importance to the Ohlone people.

CONFLICT OF INTEREST STATEMENT

The authors declare no competing interests.

PEER REVIEW

The peer review history for this article is available in the Supporting Information for this article.

DATA AVAILABILITY STATEMENT

Data used in this study are all publicly available. GO annotation files were downloaded from (<http://current.geneontology.org/annotations/index.html> 2022-05-16 release, accessed June 13, 2022) and Phytozome (<https://data.jgi.doe.gov/refine-download/phytozome> version 13, accessed June 13, 2022). Gene data were downloaded from sources indicated on the GO (<http://current.geneontology.org/products/pages/downloads.html> accessed June 13, 2022), Phytozome, Pombase (<http://www.pombase.org/> “pombase-2022-10-01” release), and UniProt (<https://www.uniprot.org/> accessed June 13, 2022). Data S1 provides detailed information about annotation and gene source databases for each species, downloaded versions, and URLs. Scripts used to process the data and generate the graphs are written in Python 3 and are available at GitHub (<https://github.com/TheRheeLab/AnnotationStats>).

ORCID

Bo Xue  <https://orcid.org/0000-0002-7957-3033>

Seung Y. Rhee  <https://orcid.org/0000-0002-7572-4762>

REFERENCES

- Ankeny, R. A., & Leonelli, S. (2020). Model organisms. In *Elements in the philosophy of biology*. Cambridge University Press. <https://doi.org/10.1017/9781108593014>
- Ashburner, M., Ball, C. A., Blake, J. A., Botstein, D., Butler, H., Cherry, J. M., Davis, A. P., Dolinski, K., Dwight, S. S., Eppig, J. T., Harris, M. A., Hill, D. P., Issel-Tarver, L., Kasarskis, A., Lewis, S., Matese, J. C., Richardson, J. E., Ringwald, M., Rubin, G. M., & Sherlock, G. (2000). Gene Ontology: Tool for the unification of biology. *Nature Genetics*, 25(1), 25–29. <https://doi.org/10.1038/75556>
- Binns, D., Dimmer, E., Huntley, R., Barrell, D., O'Donovan, C., & Apweiler, R. (2009). QuickGO: A web-based tool for Gene Ontology searching. *Bioinformatics*, 25(22), 3045–3046. <https://doi.org/10.1093/bioinformatics/btp536>
- Bossi, F., Fan, J., Xiao, J., Chandra, L., Shen, M., Dorone, Y., Wagner, D., & Rhee, S. Y. (2017). Systematic discovery of novel eukaryotic transcriptional regulators using sequence homology independent prediction. *BMC Genomics*, 18(1), 480. <https://doi.org/10.1186/s12864-017-3853-9>
- Camon, E., Magrane, M., Barrell, D., Lee, V., Dimmer, E., Maslen, J., Binns, D., Harte, N., Lopez, R., & Apweiler, R. (2004). The Gene Ontology Annotation (GOA) Database: Sharing knowledge in Uniprot with Gene Ontology. *Nucleic Acids Research*, 32(Database issue), 262D–266D. <https://doi.org/10.1093/nar/gkh021>
- Carbon, S., Ireland, A., Mungall, C. J., Shu, S. Q., Marshall, B., Lewis, S., the AmiGO Hub, & the Web Presence Working Group. (2009). AmiGO: Online access to ontology and annotation data. *Bioinformatics*, 25(2), 288–289. <https://doi.org/10.1093/bioinformatics/btn615>
- Collins, K., Zhao, K., Jiao, C., Xu, C., Cai, X., Wang, X., Ge, C., Dai, S., Wang, Q., Wang, Q., Fei, Z., & Zheng, Y. (2019). SpinachBase: A central portal for spinach genomics. *Database: The Journal of Biological Databases and Curation*, 2019, baz072. <https://doi.org/10.1093/database/baz072>
- Edwards, K. D., Fernandez-Pozo, N., Drake-Stowe, K., Humphry, M., Evans, A. D., Bombarely, A., Allen, F., Hurst, R., White, B., Kernodle, S. P., Bromley, J. R., Sanchez-Tamburrino, J. P., Lewis, R. S., & Mueller, L. A. (2017). A reference genome for *Nicotiana tabacum* enables map-based cloning of homeologous loci implicated in nitrogen utilization efficiency. *BMC Genomics*, 18(1), 448. <https://doi.org/10.1186/s12864-017-3791-6>
- Fernandez-Pozo, N., Menda, N., Edwards, J. D., Saha, S., Teclé, I. Y., Strickler, S. R., Bombarely, A., Fisher-York, T., Pujar, A., Foerster, H., Yan, A., & Mueller, L. A. (2015). The Sol Genomics Network (SGN)—From genotype to phenotype to breeding. *Nucleic Acids Research*, 43(Database issue), D1036–D1041. <https://doi.org/10.1093/nar/gku1195>
- Fields, S., & Johnston, M. (2005). Cell biology. Whither model organism research? *Science*, 307(5717), 1885–1886. <https://doi.org/10.1126/science.1108872>
- Food and Agriculture Organization of the United Nations. (2021). FIGURE 21: World production of crops, main commodities. FAO Statistical Yearbook 2021. https://doi.org/10.4060/cb4477en-fig_21
- Goodstein, D. M., Shu, S., Howson, R., Neupane, R., Hayes, R. D., Fazo, J., Mitros, T., Dirks, W., Hellsten, U., Putnam, N., & Rokhsar, D. S. (2012). Phytozome: A comparative platform for green plant genomics. *Nucleic Acids Research*, 40(D1), D1178–D1186. <https://doi.org/10.1093/nar/gkr944>
- Guberman, J. M., Ai, J., Arnaiz, O., Baran, J., Blake, A., Baldock, R., Chelala, C., Croft, D., Cros, A., Cutts, R. J., Di Génova, A., Forbes, S., Fujisawa, T., Gadaleta, E., Goodstein, D. M., Gundem, G., Haggarty, B., Haider, S., Hall, M., ... Kasprzyk, A. (2011). BioMart Central Portal: An open database network for the biological community. *Database: The Journal of Biological Databases and Curation*, 2011, bar041.
- Harris, M. A., Rutherford, K. M., Hayles, J., Lock, A., & Bähler, J. (2022). Fission stories: Using PomBase to understand *Schizosaccharomyces pombe* biology. *Genetics*, 220(4), iyab222. <https://academic.oup.com/genetics/article-abstract/220/4/iyab222/6481557>, <https://doi.org/10.1093/genetics/iyab222>
- Higgins, D. P., Weisman, C. M., Lui, D. S., D'Agostino, F. A., & Walker, A. K. (2022). Defining characteristics and conservation of poorly annotated genes in *Caenorhabditis elegans* using WormCat 2.0. *Genetics*, 221(4), iyac085. <https://doi.org/10.1093/genetics/iyac085>
- Hirsch, C. D., Hamilton, J. P., Childs, K. L., Cepela, J., Crisovan, E., Vaillancourt, B., Hirsch, C. N., Habermann, M., Neal, B., & Buell, C. R. (2014). Spud DB: A resource for mining sequences, genotypes, and phenotypes to accelerate potato breeding. *The Plant Genome*, 7(1), lantgenome2013.12.0042. <https://doi.org/10.3835/plantgenome2013.12.0042>
- Huntley, R. P., Sawford, T., Mutowo-Muullenet, P., Shypitsyna, A., Bonilla, C., Martin, M. J., & O'Donovan, C. (2015). The GOA database: Gene Ontology annotation updates for 2015. *Nucleic Acids Research*, 43(D1), D1057–D1063. <https://doi.org/10.1093/nar/gku1113>
- Jain, N., Mangal, P., & Mehta, D. (2014). AngularJS: A modern MVC framework in JavaScript. *Journal of Global Research in Computer*



- Sciences, 5(12), 17–23. <http://www.jgrcs.info/index.php/jgrcs/article/download/952/610>
- Jones, A. M., Chory, J., Dangl, J. L., Estelle, M., Jacobsen, S. E., Meyerowitz, E. M., Nordborg, M., & Weigel, D. (2008). The impact of Arabidopsis on human health: Diversifying our portfolio. *Cell*, 133(6), 939–943. <https://doi.org/10.1016/j.cell.2008.05.040>
- Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen, J., Mitchell, A., Nuka, G., Pesseat, S., Quinn, A. F., Sangrador-Vegas, A., Scheremetjew, M., Yong, S.-Y., Lopez, R., & Hunter, S. (2014). InterProScan 5: Genome-scale protein function classification. *Bioinformatics*, 30(9), 1236–1240. <https://doi.org/10.1093/bioinformatics/btu031>
- Kinsella, R. J., Kähäri, A., Haider, S., Zamora, J., Proctor, G., Spudich, G., Almeida-King, J., Staines, D., Derwent, P., Kerhornou, A., Kersey, P., & Flicek, P. (2011). Ensembl BioMart: A hub for data retrieval across taxonomic space. *Database: The Journal of Biological Databases and Curation*, 2011, bar030.
- Lamesch, P., Berardini, T. Z., Li, D., Swarbreck, D., Wilks, C., Sasidharan, R., Muller, R., Dreher, K., Alexander, D. L., Garcia-Hernandez, M., Karthikeyan, A. S., Lee, C. H., Nelson, W. D., Ploetz, L., Singh, S., Wensel, A., & Huala, E. (2012). The Arabidopsis Information Resource (TAIR): Improved gene annotation and new tools. *Nucleic Acids Research*, 40(D1), D1202–D1210. <https://doi.org/10.1093/nar/gkr1090>
- Lobb, B., Tremblay, B. J.-M., Moreno-Hagelsieb, G., & Doxey, A. C. (2020). An assessment of genome annotation coverage across the bacterial tree of life. *Microbial Genomics*, 6(3), e000341. <https://doi.org/10.1099/mgen.0.000341>
- National Research Council. (2015). *Industrialization of biology: A roadmap to accelerate the advanced manufacturing of chemicals*. National Academies Press.
- O'Leary, N. A., Wright, M. W., Brister, J. R., Ciufu, S., Haddad, D., McVeigh, R., Rajput, B., Robbertse, B., Smith-White, B., Ako-Adjei, D., Astashyn, A., Badretdin, A., Bao, Y., Blinkova, O., Brover, V., Chetvermin, V., Choi, J., Cox, E., Ermolaeva, O., ... Pruitt, K. D. (2016). Reference sequence (RefSeq) database at NCBI: Current status, taxonomic expansion, and functional annotation. *Nucleic Acids Research*, 44(D1), D733–D745. <https://doi.org/10.1093/nar/gkv1189>
- Ouyang, S., Zhu, W., Hamilton, J., Lin, H., Campbell, M., Childs, K., Thibaud-Nissen, F., Malek, R. L., Lee, Y., Zheng, L., Orvis, J., Haas, B., Wortman, J., & Buell, C. R. (2007). The TIGR rice genome annotation resource: Improvements and new features. *Nucleic Acids Research*, 35(Database), D883–D887. <https://doi.org/10.1093/nar/gkl976>
- Reid, W. V., Ali, M. K., & Field, C. B. (2020). The future of bioenergy. *Global Change Biology*, 26(1), 274–286. <https://doi.org/10.1111/gcb.14883>
- Tello-Ruiz, M. K., Naithani, S., Gupta, P., Olson, A., Wei, S., Preece, J., Jiao, Y., Wang, B., Chougule, K., Garg, P., Elser, J., Kumari, S., Kumar, V., Contreras-Moreira, B., Naamati, G., George, N., Cook, J., Bolser, D., D'Eustachio, P., ... Ware, D. (2021). Gramene 2021: Harnessing the power of comparative genomics and pathways for plant research. *Nucleic Acids Research*, 49(D1), D1452–D1463. <https://doi.org/10.1093/nar/gkaa979>
- The Gene Ontology Consortium. (2021). The Gene Ontology resource: Enriching a GOld mine. *Nucleic Acids Research*, 49(D1), D325–D334. <https://doi.org/10.1093/nar/gkaa1113>
- Tilkov, S., & Vinoski, S. (2010). Node.js: Using JavaScript to build high-performance network programs. *IEEE Internet Computing*, 14(6), 80–83. <https://doi.org/10.1109/MIC.2010.145>
- UniProt Consortium. (2019). UniProt: A worldwide hub of protein knowledge. *Nucleic Acids Research*, 47(D1), D506–D515. <https://doi.org/10.1093/nar/gky1049>
- Van Bel, M., Silvestri, F., Weitz, E. M., Kreft, L., Botzki, A., Coppens, F., & Vandepoele, K. (2022). PLAZA 5.0: Extending the scope and power of comparative and functional genomics in plants. *Nucleic Acids Research*, 50(D1), D1468–D1474. <https://doi.org/10.1093/nar/gkab1024>
- Wimalanathan, K., Friedberg, I., Andorf, C. M., & Lawrence-Dill, C. J. (2018). Maize GO annotation—Methods, evaluation, and review (maize-GAMER). *Plant Direct*, 2(4), e00052. <https://doi.org/10.1002/pld3.52>
- Wood, V., Lock, A., Harris, M. A., Rutherford, K., Bähler, J., & Oliver, S. G. (2019). Hidden in plain sight: What remains to be discovered in the eukaryotic proteome? *Open Biology*, 9(2), 180241. <https://doi.org/10.1098/rsob.180241>
- Woodhouse, M. R., Cannon, E. K., Portwood, J. L. 2nd, Harper, L. C., Gardiner, J. M., Schaeffer, M. L., & Andorf, C. M. (2021). A pan-genomic approach to genome databases using maize as a model system. *BMC Plant Biology*, 21(1), 385. <https://doi.org/10.1186/s12870-021-03173-5>
- Yuan, Q., Ouyang, S., Wang, A., Zhu, W., Maiti, R., Lin, H., Hamilton, J., Haas, B., Sultana, R., Cheung, F., Wortman, J., & Buell, C. R. (2005). The institute for genomic research Osa1 rice genome annotation database. *Plant Physiology*, 138(1), 18–26. <https://doi.org/10.1104/pp.104.059063>

SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

How to cite this article: Xue, B., & Rhee, S. Y. (2023). Status of genome function annotation in model organisms and crops. *Plant Direct*, 7(7), e499. <https://doi.org/10.1002/pld3.499>