


The Arabidopsis Information Resource in 2024

Leonore Reiser ,* Erica Bakker, Sabarinath Subramaniam, Xingguo Chen, Swapnil Sawant, Kartik Khosa, Trilok Prithvi, Tanya Z. Berardini*

Phoenix Bioinformatics, Newark, CA 94560, USA

*Corresponding author: 39899 Balentine Drive, Suite 200, Newark, CA 94560, USA. Email: lreiser@arabidopsis.org; *Corresponding author: 39899 Balentine Drive, Suite 200, Newark, CA 94560, USA. Email: tberardini@arabidopsis.org

Since 1999, The Arabidopsis Information Resource (www.arabidopsis.org) has been curating data about the *Arabidopsis thaliana* genome. Its primary focus is integrating experimental gene function information from the peer-reviewed literature and codifying it as controlled vocabulary annotations. Our goal is to produce a “gold standard” functional annotation set that reflects the current state of knowledge about the Arabidopsis genome. At the same time, the resource serves as a nexus for community-based collaborations aimed at improving data quality, access, and reuse. For the past decade, our work has been made possible by subscriptions from our global user base. This update covers our ongoing biocuration work, some of our modernization efforts that contribute to the first major infrastructure overhaul since 2011, the introduction of JBrowse2, and the resource's role in community activities such as organizing the structural reannotation of the genome. For gene function assessment, we used gene ontology annotations as a metric to evaluate: (1) what is currently known about Arabidopsis gene function and (2) the set of “unknown” genes. Currently, 74% of the proteome has been annotated to at least one gene ontology term. Of those loci, half have experimental support for at least one of the following aspects: molecular function, biological process, or cellular component. Our work sheds light on the genes for which we have not yet identified any published experimental data and have no functional annotation. Drawing attention to these unknown genes highlights knowledge gaps and potential sources of novel discoveries.

Keywords: plant genomics; model organism database; community resource; biocuration

Introduction

The Arabidopsis Information Resource (TAIR; <http://arabidopsis.org>) is a comprehensive online digital research resource for the biology of *Arabidopsis thaliana* (Huala et al. 2001; Garcia-Hernandez et al. 2002; Berardini et al. 2015; Reiser et al. 2022). The TAIR database contains information about genes, proteins, gene expression, alleles, mutant phenotypes, germplasms, clones, genetic markers, genetic and physical maps, publications, and the research community.

TAIR is a curated database; data are processed by Ph.D.-level plant biologists who ensure their accuracy. Curation adds value to the large-scale genomic data by incorporating information from diverse sources and making accurate associations between related data. Data from manual literature curation, such as protein localization, biochemical function, gene expression, and phenotypes, are added to the corpus of knowledge presented for each locus in the genome. TAIR aims to produce a “gold standard” functionally annotated plant genome that plant biologists can use as a reference for understanding gene function in crop species and other plants of importance to humans (Berardini et al. 2015). The resource also provides data analysis and visualization tools whose usage has recently been described (Reiser et al. 2022).

Initially funded for 14 years by the US National Science Foundation, TAIR has been sustained by subscriptions from academic institutions, corporations, research institutes, and individuals since 2014 (Reiser et al. 2016).

For 25 years, a generation of scientists has relied upon TAIR for up-to-date, high-quality information and tools provided by

scientists and software developers who interact with and respond to the needs of the community. It is a true model organism database used not only by scientists whose primary research organism is *A. thaliana* but also by the broader biological research community that uses knowledge gained in this organism to inform their understanding of their organisms. This update covers work done by the resource's staff in the last few years in the areas of genome functional and structural annotation, tool improvement, findable, accessible, interoperable and reusable (FAIR) data advocacy, and community service.

Functional annotation of *A. thaliana* genes using the gene ontology

Since 2001, TAIR curators have been using gene ontology (GO) for manual curation of Arabidopsis gene functions from the literature. GO, which describes the biological roles, molecular activities, and subcellular localization of gene products, has emerged as the de facto standard for gene functional annotation (Ashburner et al. 2000; Gene Ontology Consortium et al. 2023). GO curation is the process of extracting and codifying experimental knowledge into annotations that can be used in computational analyses. GO annotations are primarily used to predict functions of unknown genes and newly sequenced genomes, and for gene set analyses for hypothesis generation. TAIR's ultimate goal is to maintain a gold standard annotated reference plant genome (Berardini et al. 2015) that serves as a baseline for predicting gene function in other species, as well as a comparator to other genomes.

GO annotation represents just one aspect of the functional data TAIR curates from the literature. TAIR curators also use the plant ontology (PO) to capture gene expression information, craft gene summaries, as well as adding allele and phenotype information, all of which are linked to individual genes. As of 2023 October, 13,439 loci have curated summaries, 7,684 loci have one or more phenotypes, 23,123 loci have a total of over 550,000 gene expression annotations to PO terms for gene structure and growth and developmental stages, and 25,500 loci have been linked to primary literature. These counts include information for both sequenced and genetic loci. As genetic loci are cloned, we merge the relevant related records into those of the now known sequenced locus. The locus detail pages in TAIR present a comprehensive view of each locus that includes the data curated from the literature as well as other data sources that help build a more complete picture of an individual locus' function.

The functional annotation datasets generated by TAIR support a number of analytic and predictive tools for Arabidopsis and other plant species. For example, TAIR's GO annotation data are part of the underlying dataset the Protein ANalysis THrough Evolutionary Relationships (PANTHER) gene set enrichment tool (Mi et al. 2013) that is widely used for statistical analysis of over/underrepresented gene sets in Arabidopsis or as part of network analysis tools such as Knetminer (Hassani-Pak et al. 2021). As noted below, orthology-based assignment of gene function to other plant species also often relies upon Arabidopsis reference genome annotations. Additionally, these manually curated datasets (singularly or in combination) will likely prove even more useful to train large language models.

The importance of curating experimental data

Successful computational methods for inferring gene function invariably rely upon a foundational dataset grounded in experimental evidence. Since many new plant genomes generate their GO annotations based on similarity to Arabidopsis, having a well-annotated genome supported by experimental data is essential to producing high-quality computationally annotated genomes. Each GO annotation includes an evidence code which allows a user to trace whether the supporting evidence is experimental or non-experimental. Annotations with experimental evidence codes are supported by wet lab work, using either low throughput (e.g. BiFC experiments) or high throughput (e.g. proteomics data) methods. Non-experimental annotations are supported by methods that include computational pipelines such as InterPro2GO (Jones et al. 2014) that use mappings between domains and functions to assign terms (evidence code of Inferred from Electronic Annotation [IEA]) and phylogeny-based methods like Phylogenetic Annotation INference Tool (PAINT) (Gaudet et al. 2011), in which annotations are transferred based on descent from a common ancestor (evidence code of Inferred from Biological aspect of Ancestor [IBA]). In evaluating GO annotations and analysis results that use those annotations, researchers should consider the type of evidence as well as the specificity of the GO term. IEA annotations tend to use more general terms, whereas experimentally supported functions tend to use more specific terms.

GO annotation datasets for Arabidopsis change over time

GO annotation datasets are subject to change and those changes can affect the analysis and interpretation of data (Jacobson et al. 2018; The Gene Ontology Consortium 2019). As with all biological knowledge, what we know about gene function can change

over time as new functions are discovered and published. Annotations are also subjected to periodic quality checks to ensure the validity of the data. For example, IEA annotations are removed from the GO after one year and, where possible, replaced with updated (and presumably better) data. All changes to experimentally based annotations are reviewed by curators. Changes to datasets can also occur because of changes to the ontologies themselves (e.g. term inserts, deletes, and merges) that necessitate re-examination of the gene-term associations. At TAIR, annotations are updated on a weekly basis on the website and exported on a quarterly basis to the GO where those annotations are merged with Arabidopsis annotations from other sources such as UniProt (The UniProt Consortium et al. 2023) and the GO Consortium (GOC) (Gaudet et al. 2011). TAIR does integrate the *A. thaliana* annotations made by UniProt and the GOC on a regular basis, as new files are released by these groups. For these reasons, we strongly advise researchers to use the most current annotation data sets either from the TAIR website (https://www.arabidopsis.org/download/index-auto.jsp?dir=%2Fdownload_files%2FGO_and_PO_Annotations%2FGene_Ontology_Annotations) or the GOC website (<http://geneontology.org/docs/download-go-annotations/>) for any downstream applications such as gene set enrichment analysis.

Current status of gene function annotation in Arabidopsis

We used GO annotation as a proxy to assess the percentage of genes for which some functional information is available. For this analysis, we focused on the proteome (27,657 total, annotation version: Araport11) because the majority of annotations are made to protein coding genes, and it is the dataset used for orthology-based predictions of gene function.

Before delving into the actual numbers for Arabidopsis, it is important to define what we mean by “known” and “unknown” genes. For each aspect of the GO, we define a “known” gene as one having at least one annotation to that aspect that is supported by an experimental or non-experimental evidence code. “Unknown” genes are defined by having a GO annotation to the “root” term of the ontology, using the ND (No biological Data available) evidence code. For example, a gene product that has no predicted or experimental biological activity would be annotated to the GO term molecular function (GO:0003674), with the evidence code ND, to indicate that the molecular activity of the gene product is unknown at the time of literature review. For each aspect of the GO, we determined (1) the fraction of the proteome that was unknown (UNK), experimentally determined (EXP), and non-experimentally determined (non-EXP) and (2) the numbers and identities of the “unknown” gene set.

Figure 1 shows a stacked histogram where each bar represents an aspect of the GO. For each aspect, the UNK (yellow), EXP (blue), and non-EXP (red) percentages are shown. The most well-annotated aspect is the GO cellular component with 90% of the proteome having either experimental or predicted localization. This is likely due to the relative ease of assaying protein localization and large numbers of proteomics datasets available providing experimental evidence, as well as the relatively facile ability to predict localization based on structural features such as nuclear localization sequences or transmembrane domains. The least well described aspect is GO molecular function, for which 39% of the genome is unknown. Again, this is not surprising considering the difficulty of systematically assessing molecular activities (e.g. a specific enzymatic activity) relative to the generalized biological processes for which those molecular activities are necessary, such as “cell proliferation”.

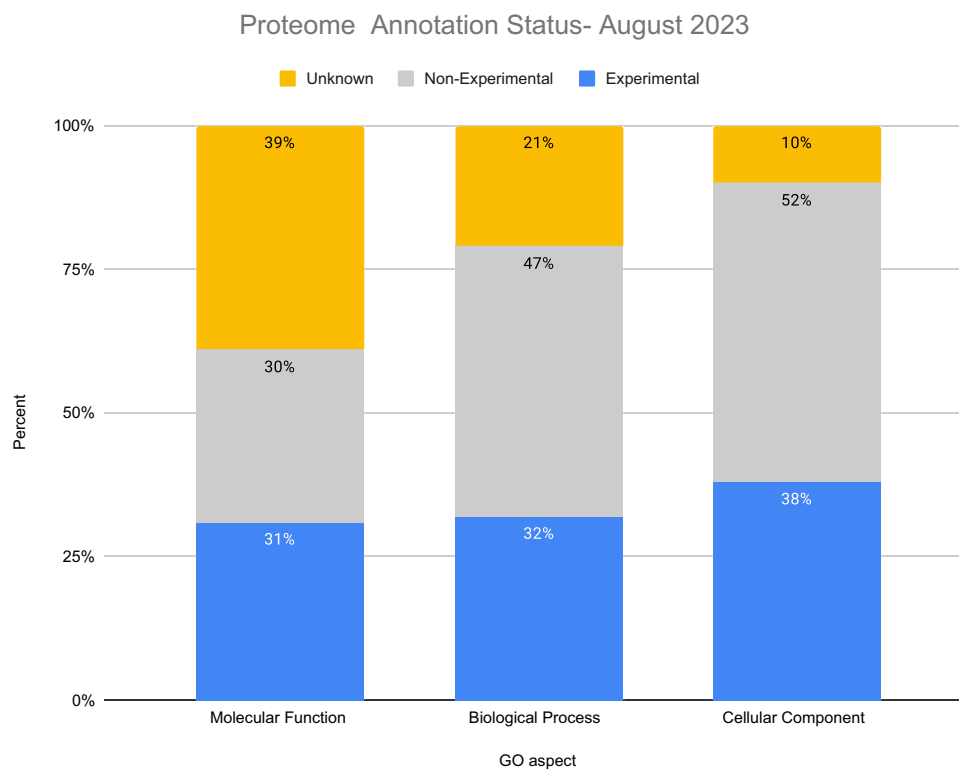


Fig. 1. Histogram showing the annotation status of the Arabidopsis proteome by GO aspect and GO evidence class. The unknown set includes proteins with annotations to the root ontology term using the evidence code ND. The experimental set includes proteins with at least one annotation using one of these evidence codes: Inferred from Direct Assay, Inferred from Expression Pattern, Inferred from Genetic Interaction, Inferred from Mutant Phenotype, Inferred from Physical Interaction, inferred from High-throughput Direct Assay, inferred from High-throughput Expression Pattern, or inferred from EXPeriment (EXP). The non-experimental set includes proteins ONLY having annotations of at least one of the following evidence codes: IEA, Inferred from Sequence or Structural Similarity, Non-traceable Author Statement, Traceable Author Statement, Inferred by Curator, Inferred from Reviewed Computational Analysis, IBA, and Inferred from Sequence Model.

To identify the set of unknown genes, we sought the intersection of unknowns for each GO aspect. Figure 2 shows a Venn diagram displaying the intersection of unknown genes from each aspect. A total of 1,224 protein coding genes lack any functional annotations at all. An additional 3,374 lack annotations for biological process and molecular function and thus can also be classified as unknown. A total of 375 of these unknowns are annotated as “hypothetical protein” and may not actually be real genes. Regularly updated versions of this list are available on the TAIR website (<https://conf.phoenixbioinformatics.org/pages/viewpage.action?pageId=22807120>).

Other groups have also used GO as a proxy for known-ness. A different representation of the GO annotation status of Arabidopsis (and other genomes) can be accessed via the Genome Annotation Status Charts (<https://genomeannotation.rheelab.org/>) generated by Xue and Rhee (2023). Observed differences between the data presented here and in the snapshot may be due to differences in the source files (we used our curation database whereas the Xue paper uses the gene association format files from the GOC) and gene set (whole genome vs proteome only). The smaller set of unknowns presented here is likely because we limited our analysis to protein coding genes.

Characteristics of the unknown gene set

There are both biological and non-biological reasons why some proteins remain unknown. Among the 4,598 “unknown” genes, 375 are annotated as being “hypothetical proteins” in the most recent public annotation version (Araport11) meaning their existence is questionable. A small number of these “hypothetical

protein” genes are not present in the in-progress structural re-annotation of the *A. thaliana* genome (see community-driven project below) and new genes were added so these numbers will likely be adjusted in the next genome release. Biological reasons might include genetic redundancy or difficult to assess phenotypes. Some unknown proteins might belong to members of plant-specific gene families and therefore would not have been included in the phylogenetic-based inference work done as part of the PAINT project because that focus is on curating families with representatives from the human genome. These plant-specific proteins might have novel functions that have yet to be represented within the GO.

We used the PhyloGenes resource (Zhang et al. 2020) to examine the phylogenetic classifications of the unknown proteins. First, we mapped the unknown protein IDs to their corresponding entries in the PhyloGenes database and then identified the set of corresponding PANTHER families. We then queried for the stored highest level taxonomic distribution for each family, that is, the highest order taxon that includes all of the families. Among the 4,598 proteins, 4,070 were mapped to 1,650 distinct PANTHER families (PantherDB v.17). Figure 3 shows the distribution of these PANTHER families with unknown Arabidopsis protein members in each of the taxonomic groups. While some of the families have members from all of Eukaryota, 70% belong to “plant-specific” taxonomic groupings ranging from Viridiplantae to Brassica specific. Most have no associated GO information but some have annotations based on domains. For example, the PTHR34269 family, <http://www.phylogenies.org/tree/PTHR34269>, has curated members from Arabidopsis and is a plant-specific

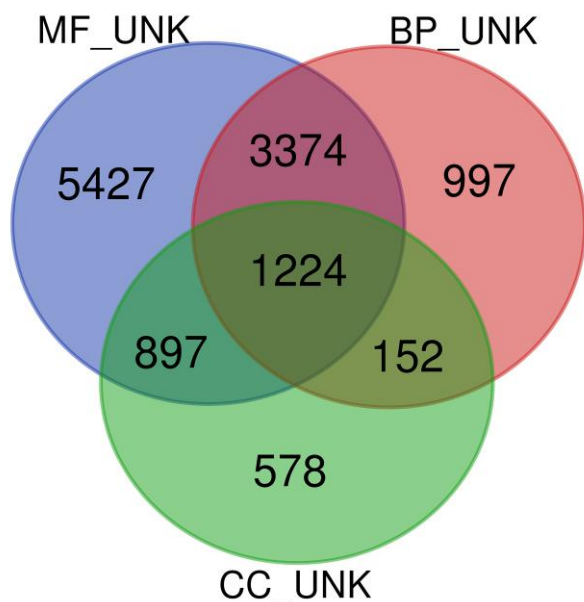


Fig. 2. Venn diagram illustrating the overlap among proteins having ND annotations to each aspect. Files containing Arabidopsis Genome Initiative (AGI) locus IDs for each aspect (MF_UNK, unknown molecular functions; BP_UNK, unknown biological process; CC_UNK, unknown cellular component; INSERT REF for files) were uploaded to the Vlaams Instituut voor Biotechnologie (VIB) Venn Diagram Generator (<http://bioinformatics.psb.ugent.be/webtools/Venn/>).

(spanning eudicots) family. This family is characterized by the presence of a B3 domain (Swaminathan et al. 2008) and includes well studied members of the Auxin Response Factor/LEAFY COTYLEDON 2/[ABSCISIC ACID INSENSITIVE 3] – [VP1/ABI3-LIKE and Reproductive Meristem sub families. B3 domain functions in (sequence specific) DNA binding. Therefore, it is likely that other unannotated members of the family also have that molecular activity and that annotation would be supported by sequence/phylogenetic analysis. There are many other families such as PTHR10826 (<http://www.phylogenies.org/tree/PTHR10826>) that span eukaryotes but plant genes (subfamilies) lack annotations even though there are experimental annotations for other non-plant species. This may be because the plant genes are sufficiently diverged into subfamilies and therefore annotations cannot be propagated without supporting evidence from plants. Another reason is that gene functions may have been described experimentally for other plant species in the literature but, because experimental plant GO annotation is limited, that data have not yet been captured as GO annotations. More comprehensive curation of plant gene function would enable propagation of IBA annotations if there was novelty in the plant lineage.

Human and financial resources are other limitations as many functions may be known but have not been captured as a computable GO annotation or what is known is based on prediction and not experimentation. Rocha et al. (2023), in creating their Unknownome database, assigned a “knowness” score clusters based on a number of factors including GO evidence weights. Again, even well conserved plant-specific families may rank low in known-ness by their metric. For example, the Dirigent protein family which is found in the Tracheophyta (https://unknownome.mrc-lmb.cam.ac.uk/cluster_details/UKP06412/) has well defined experimental functions for some of the member proteins from Arabidopsis and other species (Paniagua et al. 2017) but has a standard known-ness of 0.0. Probably the reason this cluster comes up with such a low “known” score is because it is a plant-specific gene family (restricted to

Tracheophyta; see <http://www.phylogenies.org/tree/PTHR46442>). PAINT annotations prioritize PANTHER families with human genes, therefore anything that is plant specific is unlikely to be curated based on biological ancestry.

How to fill in the knowledge gaps?

Any approach that uses GO annotation as a proxy for known-ness has limitations; the most significant being the incomplete nature of GO annotations. The output of experimental data/literature vastly outpaces the ability of curators to process the data both in terms of sheer numbers of genes described and number of articles published. This is especially true for plants where only a fraction of the knowledge published in the literature has been captured. In order to increase the functional annotation coverage of the genome and provide a more comprehensive functional annotation dataset, we need to increase throughput through a two-pronged approach that includes (1) curation at the time of publication and (2) curation of the “backlog” of papers. At TAIR we have been trying to tackle both approaches. To address the first issue, we (and others) have developed strategies and tools to encourage authors to curate their data as they publish (Berardini et al. 2012; Rutherford et al. 2014; Arnaboldi et al. 2020; Larkin et al. 2021; Reiser et al. 2022). We and others are also exploring the use of machine learning and artificial intelligence to assist in data extraction and curation from primary literature (Müller et al. 2018; Kishore et al. 2020).

Community curation: become a GOATherder

Since 2008, we have developed tools that enable researchers to contribute GO and PO annotations to TAIR and expand the gene function knowledgebase beyond what our curators can do. Since then we have processed 12,232 annotations from 1,692 papers submitted by 1,130 community members. From 2013 to 2020, we supported TAIR's Online Annotation Submission Tool (TOAST) to facilitate community curation of Arabidopsis genes (Li et al. 2012). In 2020, we replaced this tool with the Generic Online Annotation Tool (GOAT; <https://goat.phoenixbioinformatics.org/>). As with TOAST, GOAT is a literature-based curation tool, meaning it is designed for curating experimental gene function data on a per paper basis. Users can contribute annotations for their own or other people's published works. The GOAT prototype was developed over two years as capstone projects for two cohorts from the Rochester Institute of Technology.

GOAT is a simple web application that allows for basic GO and PO annotations as well as the addition of comments suitable for incorporation into gene summaries (Fig. 4). GOAT uses ORCID authentication (<https://orcid.org/>), so users must register or have an ORCID ID to begin. Once logged in, users enter the DOI or PubMed ID for the article they wish to curate. Then they can add as many genes as they want using one of the allowed name types (AGI Locus ID, UniProt ID, or RNA central ID). To annotate a gene, they must first select the “subject” gene product (from the supplied list) and then a type of annotation (GO biological process, GO molecular function, GO cellular component, PO structure, PO developmental stage, protein–protein interaction, or comment). Based on that selection, users can then search for GO or PO terms within that subset of the ontology. They can then pick an appropriate evidence code from the Evidence and Conclusion Ontology (Nadendla et al. 2022) for the experiment that supports the assertion/annotation. The interface is intuitive and a tutorial is available on TAIR's YouTube Channel (<https://www.youtube.com/watch?v=t5oB51yX6Lobrief>). Community annotations are reviewed by a TAIR curator to make sure that they are consistent with annotation rules and best practices before integration into TAIR and eventual consumption by GOC and other resources.

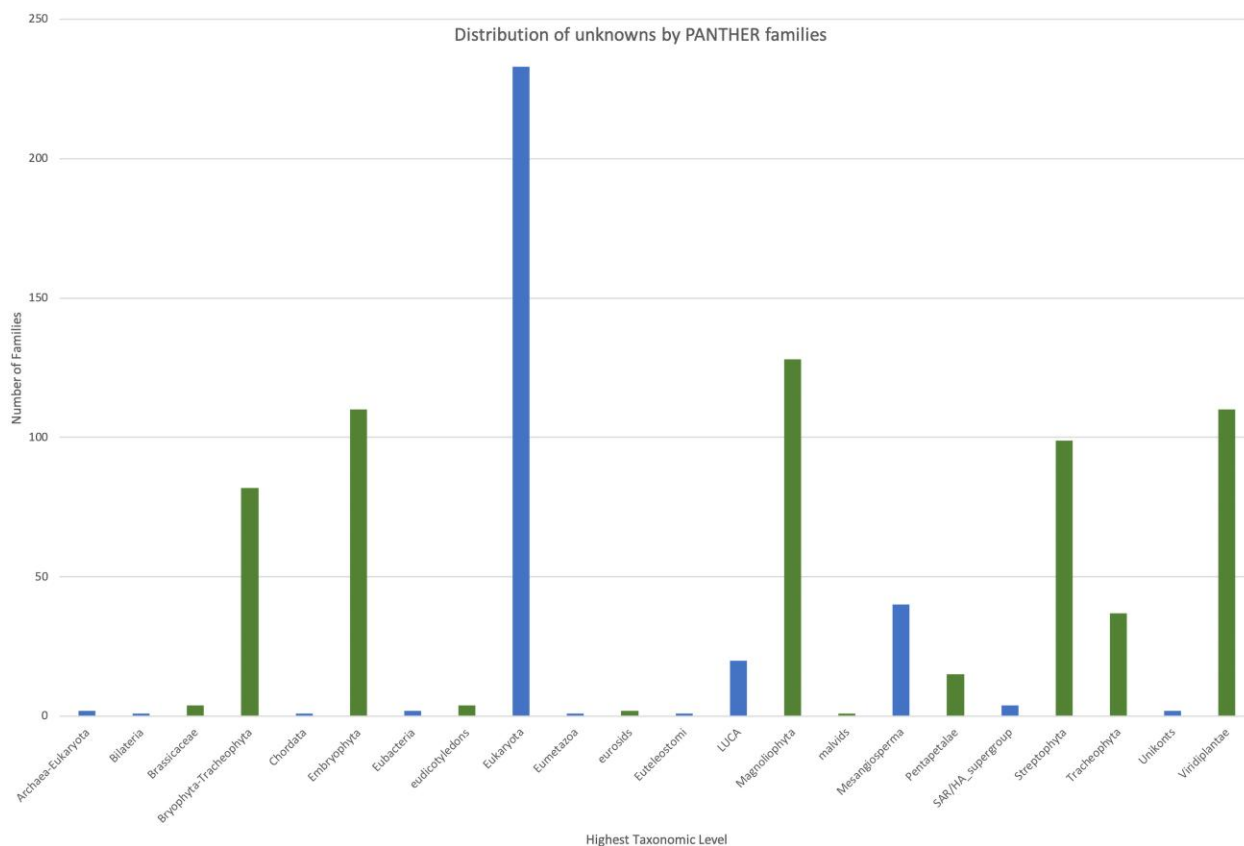


Fig. 3. Histogram displaying the distribution of PANTHER 17 gene families containing Arabidopsis unknown proteins grouped by highest taxonomic classification. Plant-specific families are indicated by green bars.

GOAT was designed for flexibility and can be used to annotate ANY gene, from any organism, as long as it has an RNA central or UniProt ID. This provides the potential to fill in gaps for gene function in Arabidopsis or any plant species. Curating non-Arabidopsis gene function allows the capture of aspects of plant biology that either do not exist in Arabidopsis (such as nodulation or wood formation) or are more well described in other species. By extending curation of experimentally defined functions to other species, we can narrow the knowledge gap and create a better representation of plant gene function across species.

Back end changes to improve database speed and stability

The Locus Detail pages are the most highly used pages at TAIR. For the time period spanning 2022 November 1 to 2023 October 31, our usage analytics show that locus pages were accessed 7,970,973 times. They consolidate all things gene function related: GO annotations, symbols and full names, summary, publications, alleles, germplasms, stock and clone information, RNA, protein and expression information, gene family and homolog data, as well as links out to external resources with complementary information about the genes. We were running into major issues with long page load times because of the substantial amount of data being retrieved from multiple tables in the TAIR Oracle database and then being aggregated. To solve this problem, we denormalized the Oracle data, and stored it in an Amazon Web Services Simple Storage Service (AWS S3) bucket as individual JavaScript Object Notation files, which considerably speed up data processing and retrieval time and reduced computational costs.

AWS S3 is a scalable, durable, and cost-effective object storage service provided by Amazon. It is a self-managed service, meaning that the organization using it (Phoenix Bioinformatics, in our case) does not have to maintain physical servers. This leads to significant reductions in maintenance costs, as there is no need to manage, upgrade, or replace server hardware. S3's high scalability means that it can handle a growing volume of data and requests without compromising performance. AWS S3 is highly distributed, which means that data are redundantly stored across multiple data centers. This eliminates the risk of a single point of failure. One of the most crucial benefits of this migration to AWS S3 is the dramatic improvement in data retrieval speed. The data retrieval time for a typical Locus Detail page was reduced from 1 min (using the previous technology) to an impressive 300 ms when utilizing AWS S3. This is a substantial enhancement in user experience, making the Locus Detail pages much more responsive. In summary, by transitioning from traditional data storage and retrieval methods to denormalized data stored in AWS S3 buckets, we have not only achieved significant performance improvements but also reduced costs and improved the overall reliability and availability of the data.

We are currently working on a complete refactoring of the website using current technology and framework to replace the early 2000s-era software, with new modules deployed on the beta.arabidopsis.org website as they are completed.

Moving to a better genome browser, JBrowse2

Since 2020 May, JBrowse has been serving as TAIR's primary genome browser. It is the source of the map images that we provide on both the Locus and Gene Detail pages. All new data tracks have been added exclusively to JBrowse since its introduction at TAIR.

Fig. 4. Screenshot of GOAT data submission interface after logging in via ORCID. a) Users add DOI or PubMed ID for the paper they are curating. b) Users enter locus identifiers and any gene names/symbols. Users can add more genes by clicking the “Add Another Gene” button. c) Users must enter at least one annotation for at least one gene (specified in the above list). d) Users can add as many annotations as desired. e) They can choose different types of annotations from the drop down menu. The type of annotation determines the set of GO or PO terms available as well as the types of evidence (Method). Once all annotations are entered, the user is prompted to review the submission (F) before submitting. Submissions are then reviewed by a TAIR curator before being imported into TAIR and integrated to the GO database on a quarterly basis.

The Javascript-based JBrowse was first released in 2009 (Skinner et al. 2009) and is no longer being actively developed. To reduce the overall maintenance load associated with supporting four genome browsers (SeqViewer, GBrowse, JBrowse, and JBrowse2), support for the much older technology stack-powered SeqViewer and GBrowse will be discontinued. As more features and plugins are added to JBrowse2 and that platform becomes even more stable, we anticipate that we will shift exclusively to JBrowse2 as TAIR’s genome browser and sunset JBrowse support as well. User feedback to our announcement of sunsetting SeqViewer and GBrowse has highlighted several features of SeqViewer that are particularly valued by the community.

- 1) The ability to download gene sequences that have genome coordinates, UTRs, exons, introns, and start and stop codons as well as intergenic regions marked. This can be done with the SeqLighter plugin in TAIR’s JBrowse and the feature has been requested from the JBrowse2 developers.
- 2) The ability to copy a DNA sequence from sequence viewer nucleotide view and paste to DNA editor applications like A plasmid Editor (ApE), keeping the lower case (intron and non-coding) and uppercase (exon) characters. This feature has been requested from the JBrowse2 developers.

JBrowse2 is a complete rewrite of JBrowse 1 with a similar user interface but a modern software architecture (Diesh et al. 2023). This more modern browser is under active development and maintenance, and features the capability for viewing genomic structural

variants and evolutionary relationships among genes and genomes with syntenic visualizations. JBrowse2 is built on a contemporary tech stack and boasts optimized algorithms and a streamlined code-base, making it significantly faster than its predecessor. This means quicker load times, smoother navigation, and an overall enhanced user experience. The intuitive user interface makes the platform easier to navigate for seasoned users but also lowers the learning curve for newcomers. After testing JBrowse2 in beta mode for several months at TAIR, the tool is now available on the main website (jbrowse2.arabidopsis.org/index.html). Most of the data tracks that are available in the TAIR JBrowse are present in JBrowse2 (Fig. 5). A few remain untransferred due to a current lack of the appropriate plugins for visualization and for that reason we will continue to maintain and update the original JBrowse.

A new feature introduced in JBrowse2 enables visualization capability for syntenic datasets. This feature allows researchers to compare and contrast gene order and orientation across multiple genomes in a visually intuitive manner. By overlaying syntenic regions on the reference genome, JBrowse 2 provides a comprehensive view of genomic conservation, facilitating insights into evolutionary events such as gene duplications, inversions, and translocations. This integration of syntenic datasets into JBrowse 2 empowers scientists with a powerful tool for deciphering the complexities of genome evolution. In the initial release, we provide access to the *A. thaliana* and *Arabidopsis lyrata* genomes for syntenic comparisons (Fig. 6). Additional syntenic datasets for over 30 plant species, both monocot and dicot, will be made available as they are generated.

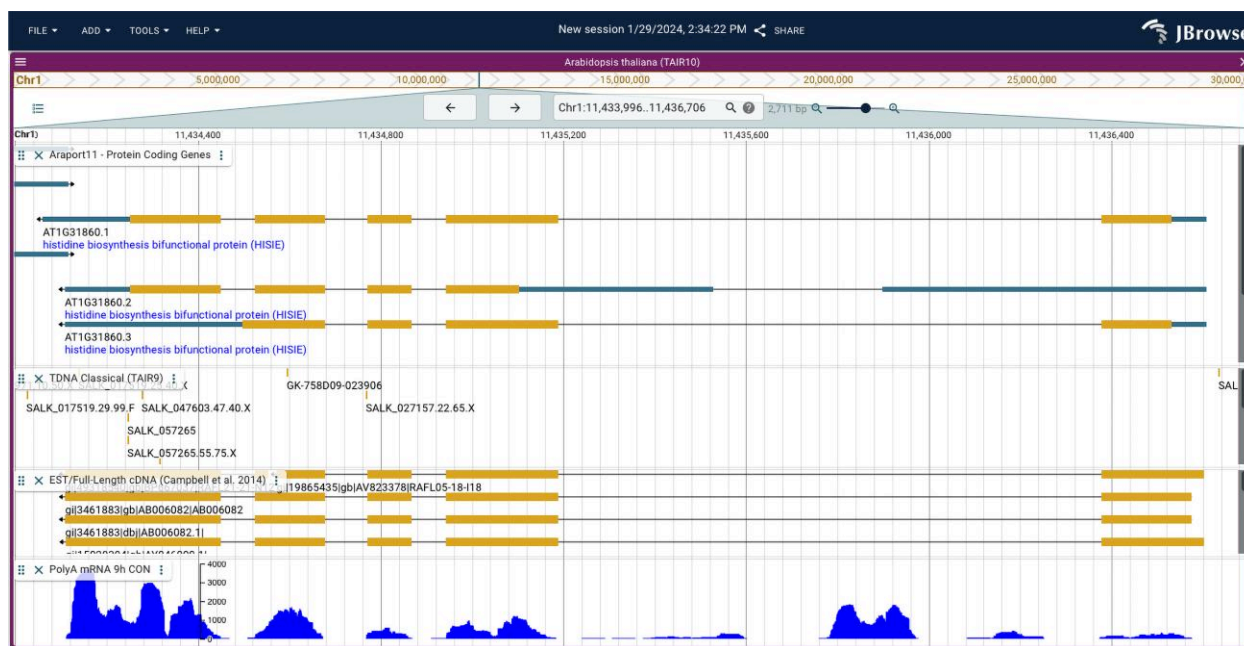


Fig. 5. JBrowse2 interface showing the locus At1g31860, the structure of the three different gene models, locations of T-DNA insertions, supporting cDNAs, and one some mRNA-seq expression data as a coverage track.

TAIR as a community hub: coordinating the reannotation of the genome

Since its inception, TAIR has functioned as a community hub for projects of broad interest, impact, and importance, such as maintaining project lists for the NSF 2010 project and the early days of the Multinational Arabidopsis Steering Committee (MASC). More recently, TAIR has played a key role in motivating, organizing, and managing the latest reannotation of the reference genome based on the Col-0 ecotype.

Brief history of genome releases

When the *Arabidopsis thaliana* genome was first sequenced and published in 2000, it marked the first complete plant genome available to the scientific community ([Arabidopsis Genome Initiative 2000](#)). After the initial annotation of the Col-0 genome, 10 subsequent versions followed. The first four (TIGR2 through TIGR5) were done by what was then called The Institute for Genome Research (TIGR; [Haas et al. 2003, 2005](#)). The next five (TAIR6 through TAIR10) were funded by an NSF grant for the TAIR project ([Swarbreck et al. 2008; Lamesch et al. 2012](#)). Araport11 ([Cheng et al. 2017](#)), the most recent version that was released in June 2016, was funded by a grant to the J. Craig Ventner Institute (JCVI) for the Araport project. Over the years, additional experimental results like ESTs and RNAseq were incorporated into prediction pipelines with manual review following the automated predictions. The manual review process was not only necessary but essential in increasing the quality of each reannotation. The resulting products were incorporated into GenBank's RefSeq section and from there were available to the broader bioinformatics and research community for use.

Initiation of V12

With the advances in sequencing, assembly, and annotation technologies, it was glaringly apparent that an update to the Araport11 release was needed. Attempts to find directed funding for this effort were not successful and so another approach was needed.

With the help of Nicholas Provart at the University of Toronto, Tanya Berardini, TAIR's Director, convened a Zoom meeting of about 20 interested members from the Arabidopsis sequencing and genome assembly community in 2022 October to assess if there was broad support and community buy-in for a community resourced approach. Scientific expertise would be provided by the community and project management, and technical support and tool hosting would be provided by TAIR. Five phases of the annotation process were identified ([Fig. 7](#)): reference sequence assembly, automatic annotation, manual review, submission to the International Nucleotide Sequence Database Collaboration, and dissemination/integration of the new reference into community tools and resources.

The response was overwhelmingly in favor of proceeding with this plan. While the science of the genome reannotation will be reported in another publication (The Arabidopsis Col-CC Reannotation Team, in preparation), we wish to share some background on the organization of the effort as part of this update.

Volunteer recruitment for all phases

Past annotation and reannotation efforts were done with dedicated funding to TIGR/TAIR/JCVI. Without dedicated funding for the 12th version, we were determined to make the most out of the Arabidopsis community's expertise and goodwill to create a resource for the entire scientific community. Scientists across the globe were either recruited or came forward to contribute their skills in bioinformatics, annotation, assembly, automated annotation, manual review of genes, systematic reannotation of transposons and transposable elements, lncRNAs, rRNAs, and repeat elements. As much as possible, we tried to provide clear expectations and deadlines for completion of the assigned tasks. All contributors will be co-authors of the reannotation publication which will be submitted for review and publication after project completion. Contributions to the effort will be acknowledged using CRediT (Contributor Roles Taxonomy; <https://credit.niso.org/>), as they are for this publication.

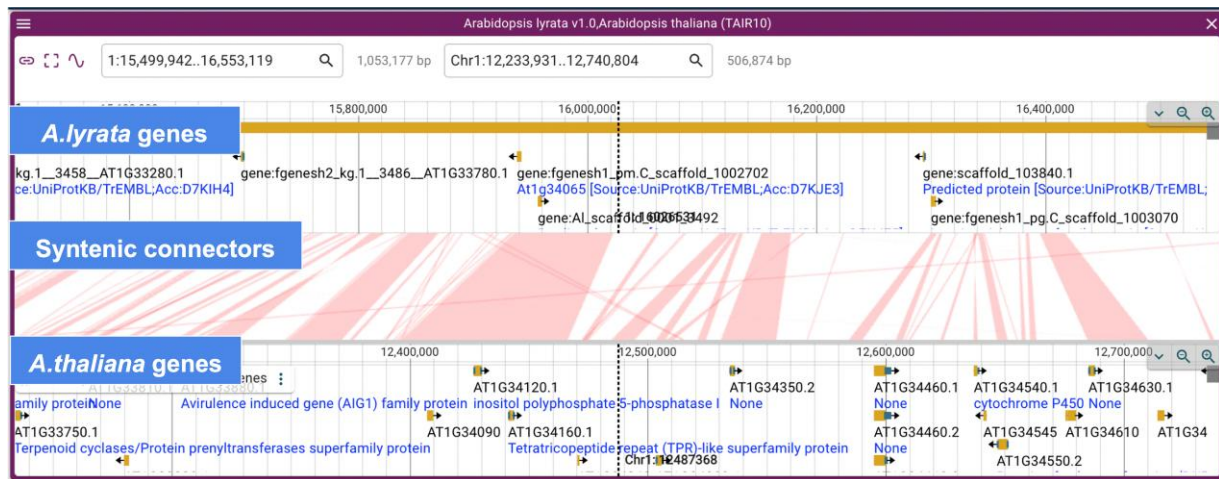


Fig. 6. JBrowse 2 visualization of syntenic comparison between genomic regions of *A. thaliana* and *A. lyrata*. Individual tracks show *A. lyrata* protein coding genes (version 1.0), “Connectors” indicating syntenic regions between *A. thaliana* and *A. lyrata*, and *A. thaliana* protein coding genes (Araport11 release). Syntenic comparison was performed using MC-Scan and the output PAF file was uploaded into JBrowse2 to create the above syntenic panel.



Fig. 7. Phases of genome reannotation.

Apollo hosting

TAIR decided to use Apollo (Dunn et al. 2019) as our community curation tool for reviewing the results of the NCBI automated annotation pipeline. After initially setting up a very small test instance, the final Apollo server was an AWS EC2 instance with 16G memory, 2 virtual Central Processing Units, and 1,000 GB storage space, which provided not only enough storage space for all of the evidence tracks but also enough capacity to perform file manipulation and transformation. Most of the evidence tracks were provided to us in General Feature Format (GFF). Some needed to be transformed to bigwig format for more intuitive visualization and gene model validation.

Manual review training

We were able to draw from the deep experience in the broader genome annotation community that has used Apollo for similar projects. Specifically, we reused and adapted teaching slides and guidelines from the maize community (shared by Marcela Tello Ruiz at Gramene), the Apollo developer group (shared by Monica Muñoz Torres, create while with Berkeley Bioinformatics Open-source Projects (BBOP)), and workflow management from the i5K project (shared by Monica Poelchau and Chris Childers at the US Department of Agriculture). We conducted six-1.5 h training sessions over 5 weeks. Over 70 participants from 10 countries attended at least one of the sessions. Weekly Zoom office hours (one and a half hours a week, with two different times to accommodate global time zones) were established for “live” feedback and troubleshooting and for almost 3 months of manual review, at least one community member took advantage of the discussion time.

Communication

With a globally distributed volunteer force, it was essential to have both synchronous and asynchronous communication channels. We established a central website for tracking progress and

milestones (<https://conf.phoenixbioinformatics.org/display/COM/A.+thaliana+Col-0+v12+reannotation+effort>). Phoenix hosted a dedicated Slack channel for the manual review part of the project. Updates were shared by email, Slack, and TAIR’s X account. The combination of all of these venues was necessary to ensure that information was distributed in a timely fashion and reached the needed audiences. Regular updates kept the community motivated, involved, and informed.

In the process of organizing and executing the reannotation project, we have identified useful tools, resources, and strategies, as well as potential pitfalls to avoid. We plan to share the resources and lessons gleaned from this experience, to help other groups that may face similar challenges (i.e. lack of funding for genome annotation).

Consolidating online community resources

Another way that TAIR serves as a community resource is by identifying and sharing useful data resources. The Arabidopsis Community Resources Portal (<https://conf.phoenixbioinformatics.org/display/COM/Resources>) is a curated collection of databases, data sets, and other digital resources of interest beyond TAIR for Arabidopsis researchers. The initial list was curated by members of the MASC Bioinformatics subgroup. Each entry is tagged with searchable keywords such as “gene_expression” and “proteomics” or entire list can be browsed via the page tree structure. We welcome suggestions and contributions from community members to add to this resource.

Promoting FAIR standards as part of the AgBioData Consortium

TAIR also engages with the broader research community to promote better practices in data management and reuse. TAIR is a

founding member of the AgBioData Consortium (www.agbiodata.org) and supports efforts to ensure that agricultural and related data are FAIR (Wilkinson et al. 2016). Toward that end, we participate in consortium-wide working groups aimed at developing data and data management standards (Harper et al. 2018; Saha et al. 2022; Clarke et al. 2023; Deng et al. 2023). We have also drafted some guidelines for authors on how to make their Arabidopsis publications more FAIR (<https://conf.phoenixbioinformatics.org/pages/viewpage.action?pageId=22807345>; (Reiser et al. 2018)) and updated our list of recommendations on where to submit data including what data TAIR accepts and what it does not (<https://www.arabidopsis.org/submit/index.jsp>). We welcome feedback from the community.

Perspective on future direction

Almost 25 years after TAIR's inception, the resource continues to grow and adapt to the changing needs of the community and the constantly shifting landscape of the technology that supports its online delivery. Changes in funding model aside, TAIR's strength has always been and continues to be its deep connections with the scientific community that it serves. We will continue to nurture those ties and use them to guide TAIR's expansion into new areas with the essential services that researchers and students have relied on for so many years.

Long-term sustainability

As a core resource for plant biologists, it is essential to have secure, long-term funding. Since 2013, TAIR has been supported by community subscriptions, and has successfully transitioned away from episodic grant funding (Reiser et al. 2016). For the last decade, TAIR has been funded largely by over 225 academic institutional subscriptions (61% of TAIR's total subscription revenue), a few national subscriptions (25%), and corporate subscriptions (10%) that cover full access to the resource for tens of thousands of scientists all over the world. There are also a couple hundred individual academic subscribers who contribute about 3% of TAIR's total subscription revenue. Subscriptions have provided a stable source of funding that supports ongoing curation and some of the enhancements and improvements outlined here. Even with modest increases to cover inflation, the renewal rate for institutional subscriptions has been fairly stable (over 95%). We continue to offer the lowest rates that we can and provide free access for (1) teaching purposes (21 courses at 20 institutions in 2023 alone), (2) to US-based Historically Black Colleges and Universities, and (3) to countries classified as low income economies by the World Bank. At this point, almost half of TAIR's lifetime has been self-supported and we look forward to continuing to provide a valued, high-quality resource to the community.

Database citation

If TAIR is either generally useful or essential in your research, please cite this publication (or any of the older TAIR publications from the reference list) whenever you publish your own work. Model organism databases (MODs) provide a huge resource for the scientific community and their contributions are not recognized often enough in the published literature. Literature citation helps track not only TAIR's but other MOD's impact in a quantifiable manner. Such metrics are essential evidence in outreach efforts to funding agencies.

Data availability

The website URL is www.arabidopsis.org. The GO annotation file available at Berardini et al. (2022) can be used to reconstruct Fig. 1. The data files used for generating Fig. 2 are available at FigShare: Reiser et al. (2023). We will update the files on a regular basis with updates available through the TAIR website at https://www.arabidopsis.org/download/index-auto.jsp?dir=%2Fdownload_files%2FGenes%2FUnknown_Gene_Lists. Cumulative data files with information on gene function, publication links, germplasm, and phenotype information, as well as GFF files with updated gene symbol and full name information are released every quarter (beginning of January, April, July, and October). Subscriber Data Releases contain data updated within the last 12 months and are available at this URL to those with current subscriptions to TAIR: https://www.arabidopsis.org/download/index-auto.jsp?dir=/download_files/Subscriber_Data_Releases. Use of these files are governed by the Terms of Use, full text available here: http://www.arabidopsis.org/doc/about/tair_terms_of_use/417.

After a year, the Subscriber Data Releases are moved into the Public Data Releases folders at this URL: https://www.arabidopsis.org/download/index-auto.jsp?dir=/download_files/Public_Data_Releases.

All files in the Public_Data_Releases folder are made available to the public under the CC-BY 4.0 license (<https://creativecommons.org/licenses/by/4.0/>).

Acknowledgments

The GOAT software was originally developed by the following Rochester Institute of Technology students: Austin Hartnett, John King, Ian Montgomery, Nick Peretti, Arron Reed, Ben Grawi, and Gavin Nishizawa, and integrated into Phoenix's technology stack by Qian Li. The authors thank the members of the plant biology research community for their continued support, feedback, and data contributions.

Funding

This work was funded by national, academic, institutional, corporate, and individual subscriptions to TAIR. Open access to this publication is funded by subscriptions.

Conflicts of interest

The author(s) declare no conflict of interest.

Author contributions

Conceptualization: LR and TZB; data curation: LR, EB, S. Subramaniam, and TZB; formal analysis: LR and TZB; investigation: LR and TZB; project administration: TZB; software: XC, KK, S. Sawant, S. Subramaniam, and TP; supervision: TZB and TP; visualization: LR, S. Subramaniam, and TZB; writing—original draft: LR, TZB, TP, S. Sawant, and S. Subramaniam; writing—review and editing: all authors.

Literature cited

Arabidopsis Genome Initiative. 2000. Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*. *Nature*. 408(6814):796–815. doi:10.1038/35048692.

- Arnaboldi V, Raciti D, Van Auken K, Chan JN, Müller H-M, Sternberg PW. 2020. Text mining meets community curation: a newly designed curation platform to improve author experience and participation at WormBase. Database (Oxford). 2020:baaa006. doi:10.1093/database/baaa006.
- Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT, et al. 2000. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. Nat Genet. 25(1):25–29. doi:10.1038/75556.
- Berardini TZ, Li D, Muller R, Chetty R, Ploetz L, Singh S, Wensel A, Huala E. 2012. Assessment of community-submitted ontology annotations from a novel database-journal partnership. Database (Oxford). 2012:bas030–bas030. doi:10.1093/database/bas030.
- Berardini T, Reiser L, Huala E. 2022. TAIR functional annotation data (TAIR_Data_20220331) [Data set]. Zenodo. doi:10.5281/zenodo.7843882.
- Berardini TZ, Reiser L, Li D, Mezheritsky Y, Muller R, Strait E, Huala E. 2015. The Arabidopsis Information Resource: making and mining the “gold standard” annotated reference plant genome. Genesis. 53(8):474–485. doi:10.1002/dvg.22877.
- Cheng C-Y, Krishnakumar V, Chan AP, Thibaud-Nissen F, Schobel S, Town CD. 2017. Araport11: a complete reannotation of the *Arabidopsis thaliana* reference genome. Plant J. 89(4):789–804. doi:10.1111/tpj.13415.
- Clarke JL, Cooper LD, Poelchau MF, Berardini TZ, Elser J, Farmer AD, Ficklin S, Kumari S, Laporte M-A, Nelson RT, et al. 2023. Data sharing and ontology use among agricultural genetics, genomics, and breeding databases and resources of the AgBioData Consortium. arXiv:2307.08958 [cs.DB]. <https://doi.org/10.48550/arXiv.2307.08958>, preprint: not peer reviewed.
- Deng CH, Naithani S, Kumari S, Cobo-Simon I, Quezada-Rodríguez EH, Skrabisova M, Gladman N, Correll MJ, Sikiru AB, Afuwape OO, et al. 2023. Genotype and phenotype data standardization, utilization and integration in the big data era for agricultural sciences. Database (Oxford). 2023:baad088. doi:10.1093/database/baad088.
- Diesh C, Stevens GJ, Xie P, Martinez TDJ, Hershberg EA, Leung A, Guo E, Dider S, Zhang J, Bridge C, et al. 2023. JBrowse 2: a modular genome browser with views of synteny and structural variation. Genome Biol. 24(1):74. doi:10.1186/s13059-023-02914-z.
- Dunn NA, Unni DR, Diesh C, Munoz-Torres M, Harris NL, Yao E, Rasche H, Holmes IH, Elsik CG, Lewis SE. 2019. Apollo: democratizing genome annotation. PLoS Comput Biol. 15(2):e1006790. doi:10.1371/journal.pcbi.1006790.
- Garcia-Hernandez M, Berardini TZ, Chen G, Crist D, Doyle A, Huala E, Knee E, Lambrecht M, Miller N, Mueller L, et al. 2002. TAIR: a resource for integrated Arabidopsis data. Funct Integr Genomics. 2(6):239–253. doi:10.1007/s10142-002-0077-z.
- Gaudet P, Livstone MS, Lewis SE, Thomas PD. 2011. Phylogenetic-based propagation of functional annotations within the Gene Ontology Consortium. Brief Bioinform. 12(5):449–462. doi:10.1093/bib/bbr042.
- Gene Ontology Consortium, Aleksander SA, Balhoff J, Carbon S, Cherry JM, Drabkin HJ, Ebert D, Feuermann M, Gaudet P, Harris NL, et al. 2023. The gene ontology knowledgebase in 2023. Genetics. 224(1):iyad031. doi:10.1093/genetics/iyad031.
- Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith RK Jr, Hannick LI, Maiti R, Ronning CM, Rusch DB, Town CD, et al. 2003. Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. Nucleic Acids Res. 31(19):5654–5666. doi:10.1093/nar/gkg770.
- Haas BJ, Wortman JR, Ronning CM, Hannick LI, Smith RK Jr, Maiti R, Chan AP, Yu C, Farzad M, Wu D, et al. 2005. Complete reannotation of the Arabidopsis genome: methods, tools, protocols and the final release. BMC Biol. 3(1):7. doi:10.1186/1741-7007-3-7.
- Harper L, Campbell J, Cannon EKS, Jung S, Poelchau M, Walls R, Andorf C, Arnaud E, Berardini TZ, Birkett C, et al. 2018. AgBioData Consortium recommendations for sustainable genomics and genetics databases for agriculture. Database (Oxford). 2018:bay088. doi:10.1093/database/bay088.
- Hassani-Pak K, Singh A, Brandizi M, Hearnshaw J, Parsons JD, Amberkar S, Phillips AL, Doonan JH, Rawlings C. 2021. KnetMiner: a comprehensive approach for supporting evidence-based gene discovery and complex trait analysis across species. Plant Biotechnol J. 19(8):1670–1678. doi:10.1111/pbi.13583.
- Huala E, Dickerman AW, Garcia-Hernandez M, Weems D, Reiser L, LaFond F, Hanley D, Kiphart D, Zhuang M, Huang W, et al. 2001. The Arabidopsis Information Resource (TAIR): a comprehensive database and web-based information retrieval, analysis, and visualization system for a model plant. Nucleic Acids Res. 29(1):102–105. doi:10.1093/nar/29.1.102.
- Jacobson M, Sedeño-Cortés AE, Pavlidis P. 2018. Monitoring changes in the gene ontology and their impact on genomic data analysis. Gigascience. 7(8):giy103. doi:10.1093/gigascience/giy103.
- Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G, et al. 2014. InterProScan 5: genome-scale protein function classification. Bioinformatics. 30(9):1236–1240. doi:10.1093/bioinformatics/btu031.
- Kishore R, Arnaboldi V, Van Slyke CE, Chan J, Nash RS, Urbano JM, Dolan ME, Engel SR, Shimoyama M, Sternberg PW, et al. 2020. Automated generation of gene summaries at the alliance of genome resources. Database (Oxford). 2020:baaa037. doi:10.1093/database/baaa037.
- Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, Muller R, Dreher K, Alexander DL, Garcia-Hernandez M, et al. 2012. The Arabidopsis Information Resource (TAIR): improved gene annotation and new tools. Nucleic Acids Res. 40(D1):D1202–D1210. doi:10.1093/nar/gkr1090.
- Larkin A, Marygold SJ, Antonazzo G, Attrill H, Dos Santos G, Garapati PV, Goodman JL, Gramates LS, Millburn G, Strelets VB, et al. 2021. FlyBase: updates to the *Drosophila melanogaster* knowledge base. Nucleic Acids Res. 49(D1):D899–D907. doi:10.1093/nar/gkaa1026.
- Li D, Berardini TZ, Muller RJ, Huala E. 2012. Building an efficient curation workflow for the Arabidopsis literature corpus. Database (Oxford). 2012:bas047–bas047. doi:10.1093/database/bas047.
- Mi H, Muruganujan A, Casagrande JT, Thomas PD. 2013. Large-scale gene function analysis with the PANTHER classification system. Nat Protoc. 8(8):1551–1566. doi:10.1038/nprot.2013.092.
- Müller H-M, Van Auken KM, Li Y, Sternberg PW. 2018. Textpresso central: a customizable platform for searching, text mining, viewing, and curating biomedical literature. BMC Bioinform. 19(1):94. doi:10.1186/s12859-018-2103-8.
- Nadendla S, Jackson R, Munro J, Quaglia F, Mészáros B, Olley D, Hobbs ET, Goralski SM, Chibucos M, Mungall CJ, et al. 2022. ECO: the evidence and conclusion ontology, an update for 2022. Nucleic Acids Res. 50(D1):D1515–D1521. doi:10.1093/nar/gkab1025.
- Paniagua C, Bilkova A, Jackson P, Dabravolski S, Riber W, Didi V, Houser J, Gigli-Bisceglia N, Wimmerova M, Budínská E, et al. 2017. Dirigent proteins in plants: modulating cell wall metabolism during abiotic and biotic stress exposure. J Exp Bot. 68(13):3287–3301. doi:10.1093/jxb/erx141.
- Reiser L, Bakker E, Subramaniam S, Chen X, Sawant S, Khosa K, Prithvi T, Berardini TZ. 2024. Supplemental Material for Reiser et al., 2023. GSA Journals. Dataset. doi:10.25386/genetics.24498637.v1.
- Reiser L, Berardini TZ, Li D, Muller R, Strait EM, Li Q, Mezheritsky Y, Vetushko A, Huala E. 2016. Sustainable funding for biocuration: the Arabidopsis Information Resource (TAIR) as a case study of

- a subscription-based funding model. Database (Oxford). 2016: baw018. doi:[10.1093/database/baw018](https://doi.org/10.1093/database/baw018).
- Reiser L, Harper L, Freeling M, Han B, Luan S. 2018. FAIR: a call to make published data more findable, accessible, interoperable, and reusable. *Mol Plant*. 11(9):1105–1108. doi:[10.1016/j.molp.2018.07.005](https://doi.org/10.1016/j.molp.2018.07.005).
- Reiser L, Subramaniam S, Zhang P, Berardini T. 2022. Using the Arabidopsis Information Resource (TAIR) to find information about Arabidopsis genes. *Curr Protoc*. 2(10):e574. doi:[10.1002/cpz1.574](https://doi.org/10.1002/cpz1.574).
- Rocha JJ, Jayaram SA, Stevens TJ, Muschalik N, Shah RD, Emran S, Robles C, Freeman M, Munro S. 2023. Functional unknowns: systematic screening of conserved genes of unknown function. *PLoS Biol*. 21(8):e3002222. doi:[10.1371/journal.pbio.3002222](https://doi.org/10.1371/journal.pbio.3002222).
- Rutherford KM, Harris MA, Lock A, Oliver SG, Wood V. 2014. Canto: an online tool for community literature curation. *Bioinformatics*. 30(12):1791–1792. doi:[10.1093/bioinformatics/btu103](https://doi.org/10.1093/bioinformatics/btu103).
- Saha S, Cain S, Cannon EKS, Dunn N, Farmer A, Hu Z-L, Maslen G, Moxon S, Mungall CJ, Nelson R, et al. 2022. Recommendations for extending the GFF3 specification for improved interoperability of genomic data. arXiv:2202.07782 [q-bio.OT]. <https://doi.org/10.48550/arXiv.2202.07782>, preprint: not peer reviewed.
- Skinner ME, Uzilov AV, Stein LD, Mungall CJ, Holmes IH. 2009. JBrowse: a next-generation genome browser. *Genome Res*. 19(9):1630–1638. doi:[10.1101/gr.094607.109](https://doi.org/10.1101/gr.094607.109).
- Swaminathan K, Peterson K, Jack T. 2008. The plant B3 superfamily. *Trends Plant Sci*. 13(12):647–655. doi:[10.1016/j.tplants.2008.09.006](https://doi.org/10.1016/j.tplants.2008.09.006).
- Swarbreck D, Wilks C, Lamesch P, Berardini TZ, Garcia-Hernandez M, Foerster H, Li D, Meyer T, Muller R, Ploetz L, et al. 2008. The Arabidopsis Information Resource (TAIR): gene structure and function annotation. *Nucleic Acids Res*. 36(Database):D1009–D1014. doi:[10.1093/nar/gkm965](https://doi.org/10.1093/nar/gkm965).
- The Gene Ontology Consortium. 2019. The gene ontology resource: 20 years and still GOing strong. *Nucleic Acids Res*. 47(D1):D330–D338. doi:[10.1093/nar/gky1055](https://doi.org/10.1093/nar/gky1055).
- The UniProt Consortium, Bateman A, Martin M-J, Orchard S, Magrane M, Ahmad S, Alpi E, Bowler-Barnett EH, Britto R, Bye-A-Jee H, et al. 2023. UniProt: the universal protein knowledgebase in 2023. *Nucleic Acids Res*. 51(D1):D523–D531. doi:[10.1093/nar/gkac1052](https://doi.org/10.1093/nar/gkac1052).
- Wilkinson MD, Dumontier M, Aalbersberg IJJ, Appleton G, Axton M, Baak A, Blomberg N, Boiten J-W, da Silva Santos LB, Bourne PE, et al. 2016. The FAIR guiding principles for scientific data management and stewardship. *Sci Data*. 3(1):160018. doi:[10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18).
- Xue B, Rhee SY. 2023. Status of genome function annotation in model organisms and crops. *Plant Direct*. 7(7):e499. doi:[10.1002/pld3.499](https://doi.org/10.1002/pld3.499).
- Zhang P, Berardini TZ, Ebert D, Li Q, Mi H, Muruganujan A, Prithvi T, Reiser L, Sawant S, Thomas PD, et al. 2020. PhyloGenes: an online phylogenetics and functional genomics resource for plant gene function inference. *Plant Direct*. 4(12):e00293. doi:[10.1002/pld3.293](https://doi.org/10.1002/pld3.293).

Editor: T. Harris