

Updates to the Alliance of Genome Resources central infrastructure

The Alliance of Genome Resources Consortium¹

¹A full list of members is provided at the end of this article.

The Alliance of Genome Resources (Alliance) is an extensible coalition of knowledgebases focused on the genetics and genomics of intensively studied model organisms. The Alliance is organized as individual knowledge centers with strong connections to their research communities and a centralized software infrastructure, discussed here. Model organisms currently represented in the Alliance are budding yeast, *Caenorhabditis elegans*, *Drosophila*, zebrafish, frog, laboratory mouse, laboratory rat, and the Gene Ontology Consortium. The project is in a rapid development phase to harmonize knowledge, store it, analyze it, and present it to the community through a web portal, direct downloads, and application programming interfaces (APIs). Here, we focus on developments over the last 2 years. Specifically, we added and enhanced tools for browsing the genome (JBrowse), downloading sequences, mining complex data (AllianceMine), visualizing pathways, full-text searching of the literature (Textpresso), and sequence similarity searching (SequenceServer). We enhanced existing interactive data tables and added an interactive table of paralogs to complement our representation of orthology. To support individual model organism communities, we implemented species-specific “landing pages” and will add disease-specific portals soon; in addition, we support a common community forum implemented in Discourse software. We describe our progress toward a central persistent database to support curation, the data modeling that underpins harmonization, and progress toward a state-of-the-art literature curation system with integrated artificial intelligence and machine learning (AI/ML).

Keywords: database; knowledgebase; software; text mining; data integration; *Drosophila*; yeast; *Caenorhabditis elegans*; zebrafish; mouse

Introduction

As has been discussed at length elsewhere (e.g. [Oliver et al. 2016](#); [Wood et al. 2022](#)), model organism knowledgebases [aka model organism databases (MODs)] provide daily utility to researchers for the design and interpretation of experiments, to computational biologists for curated data sets, and to genomic researchers for annotated genomes. Some of the major uses of the MODs have been 1-stop shopping for all information about a particular gene or obtaining cleansed data sets with standard metadata for computational analyses.

The Alliance of Genome Resources (referred to herein as the Alliance) is a consortium of MODs and the Gene Ontology Consortium (GOC). The mission of the Alliance is to support comparative genomics to investigate the genetic and genomic basis of human biology, health, and disease. To promote sustainability of the core community data resources that make up the Alliance, we implemented an extensible “knowledge commons” platform for comparative genomics built with modular, reusable infrastructure components that can support informatics resource needs across a wide range of species ([Howe et al. 2018](#); [Alliance of Genome Resources 2022](#); [Bult and Sternberg 2023](#)). In 2022, the Alliance was recognized as a Core Global Biodata Resource by the Global Biodata Coalition ([Anderson et al. 2017](#)).

Specifically, the Alliance of Genome Resources is organized as 2 interdependent units: Alliance Central and the Alliance Knowledge Centers. Alliance Central is responsible for developing and maintaining the software for data access and for the

coordination of data harmonization and data modeling activities across our members. A primary goal of Alliance Central is to reduce redundancy in systems administration and software development for model organism knowledgebases and to deploy a unified “look and feel” for access to, and display of, common data types and annotations across diverse model organisms and human, following findability, accessibility, interoperability, and reuse (FAIR) guiding principles. Model organism-specific knowledgebases serve as Alliance Knowledge Centers. Knowledge Centers are responsible for expert curation and submission of data to Alliance Central using Alliance Central infrastructure. Knowledge Centers also are responsible for organism-specific user support activities and for providing access to data types not yet supported by Alliance Central. The founding Alliance Knowledge Centers are *Saccharomyces* Genome Database (SGD; [Engel et al. 2022](#)), WormBase ([Davis et al. 2022](#); [Sternberg et al. 2024](#)), FlyBase ([Gramates et al. 2022](#)), Mouse Genome Database ([Ringwald et al. 2022](#)), the Zebrafish Information Network ([Bradford et al. 2023](#)), Rat Genome Database ([Vedi et al. 2023](#)), and the GOC ([Gene Ontology Consortium 2023](#)). The newest member, Xenbase ([Fisher et al. 2023](#)), joined the Alliance consortium in 2022.

Here, we describe our progress toward harmonizing information provided by our member resources, our development of a software infrastructure for ingest, curation, storage, analysis, and output of such information, and development of an efficient literature curation system. We start by describing new features in our web portal at AllianceGenome.org.

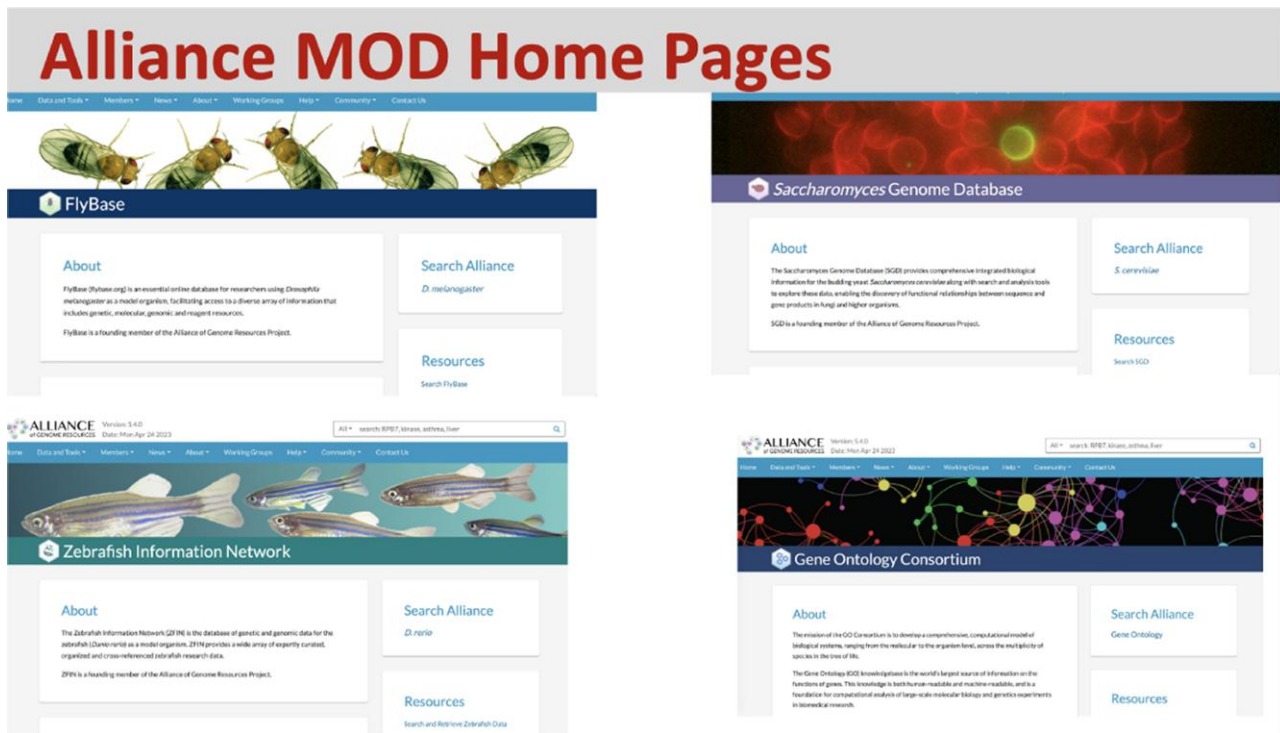


Fig. 1. MOD landing pages at the Alliance portal. A common look and feel that allows community-specific content.

The web portal Community homepages

The Alliance website features landing pages for each model organism in the Alliance consortium. These pages are accessed from the “Members” drop-down menu in the header on every Alliance page. These pages feature MOD-specific content such as meetings, news, and other MOD-specific resource links. A common template allows users to find the same types of information in each landing page (Fig. 1). As MODs transition their data and web services to the Alliance, their member pages will evolve into portals hosting additional MOD-specific data, tools, and links to organism-specific resources.

Xenopus in the Alliance

Xenbase, the *Xenopus* knowledgebase (Fisher et al. 2023), is the first knowledgebase to join the Alliance since the founding members initiated the consortium. *Xenopus* is an amphibian frog species used extensively in biomedical research and in particular for experimental embryology, cell biology, and disease modeling with genome editing (Carotenuto et al. 2023; Kostiuk and Khokha 2021). As a nonmammalian air-breathing tetrapod, *Xenopus* represents a valuable evolutionary transition between rodents and zebrafish for comparative genomic studies. Xenbase is built on the same underlying data schema (structure) as FlyBase (Chado). Two different *Xenopus* species are used interchangeably as a model system: *Xenopus tropicalis* is a diploid that is the preferred system for genome editing and genetics, whereas *Xenopus laevis* is an allotetraploid preferred for use in cell biology studies, microinjection, and microsurgery-style experimentation. *Xenopus tropicalis* has 1:1 relationships between most genes and human orthologs (excluding paralogs; Mitros et al. 2019), whereas *X. laevis* has 2 copies of most human orthologs. The allotetraploid formed via hybridization of 2 different frog species (Session et al. 2016), and the

complexities of genome evolution that subsequently occurred increase the difficulty of identifying orthology of the 2 *X. laevis* genes to their diploid relatives, including humans. Mapping of the diploid *X. tropicalis* genes to their human orthologs was performed as with the other organisms in the Alliance (see below). Because this method does not yet work in the context of an allotetraploid, the Alliance imports the *X. tropicalis* to *X. laevis* paralogy mappings from Xenbase, where they have been established through a combination of synteny analysis and manual curation; this was one major challenge in adding *Xenopus* to the Alliance.

Xenbase created software to upload content on a regular schedule formatted for the current Alliance data ingest schema. Currently, these data include orthology, the *Xenopus* anatomical ontology, standard gene information, gene expression data, publications, GO term associations, disease associations, anatomical phenotypes, and genome details. *Xenopus* genes can be found using the Alliance landing page search tool with *Xenopus* genes flagged by Xtr and Xla notations. The 2 copies of the genes in *X. laevis*, the allotetraploid, are further tagged as “(symbol).L” and “(symbol).S” to denote the genes on the long (L) and short (S) chromosome pairs of this species (e.g. *pax6.L* and *pax6.S*). Alliance release 6.0.0 has Xenbase data for 54,000 genes, 19,000 disease associations, over 45,000 gene expression records, and more than 7,000 anatomical phenotypes. Expression and phenotype data will be available in about a year.

In addition to the rich data made available to the Alliance from *Xenopus* research, this effort also served as a valuable test case for understanding the level of effort and complexities engendered in the addition of new knowledgebases to the Alliance and the functionality and adaptability of ingest system components.

New gene page section: paralogy

Gene pages now include a paralogy section populated with data from the Drosophila Research & Screening Center (DRSC)

Paralogy 



Gene symbol	Rank 	Alignment Length (aa) 	Similarity %	Identity %	Method Count	Method								
						Ensembl Compara	HGNc	InParanoid	OMA	OrthoFinder	OrthoInspector	PANTHER	PhylomeDB	SonicParanoid
<i>hlh-27</i>	1	268	99	99	3 of 8	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<i>hlh-28</i>	2	277	55	39	4 of 8	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<i>hlh-29</i>	3	279	54	38	4 of 8	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<i>hlh-26</i>	4	274	46	32	4 of 8	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<i>ref-1</i>	5	353	38	25	2 of 8	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>

Fig. 2. Paralog table for *C. elegans* *hlh-25*. The table presents a ranking of paralogs for the *hlh-25* gene, based on a weighted scoring algorithm that incorporates sequence conservation metrics. It lists the gene symbols, provides the alignment length in amino acids, and quantifies the similarity and identity percentages of genes paralogous to *hlh-25*. The methodology count, indicating the number of algorithms supporting the paralogous relationship, is also included. In this ranking, *hlh-27* is identified as the primary paralog due to its high similarity and identity scores, despite being recognized by fewer methods than *hlh-28*.

Integrative Ortholog Prediction Tool (DIOPT) version 9.1 developed by the DRSC (Hu et al. 2011, 2021). The assembly of protein sets and algorithmic inferences of their orthology from various sources was first centralized by the DRSC and then exported to the Alliance Central. We include the same data sources used for orthology, when these resources also provide paralogy information. Specifically, these resources have performed well on the standardized benchmarking from the Quest for Orthologs (QfO) Consortium (Nevers et al. 2022). Orthologous Matrix (OMA; Altenhoff et al. 2021) and PANTHER (Thomas et al. 2022) data sets were retrieved through the QfO benchmark portal (<https://orthology.benchmarkservice.org>), and Compara data were acquired directly from the EBI Compara FTP site. In addition, the DRSC conducted local analyses using InParanoid (Persson and Sonnhammer 2022), OrthoFinder (Emms and Kelly 2019), OrthoInspector (Nevers et al. 2019), and SonicParanoid (Cosentino and Iwasaki 2019) using the UniProt 2020 reference proteome set (UniProt Consortium 2023), the same set used in the downloaded data sets, to ensure consistency. Direct data submissions from PhylomeDB (Fuentes et al. 2022) and the SGD (Engel et al. 2022) were also integrated into the data set.

The new paralogy section comprises a table (Fig. 2), similar to the orthology table, that contains the gene symbol of related paralogs, a calculated rank, alignment length as the number of aligned amino acids, percentage of similarity and identity, and a count of the algorithms or methods that call the paralogous match. The ranking score was developed to sort the paralogs by overall similarity and was reviewed by curators to display optimally an acceptable rank order for well-studied sets of paralogs. The ranking score considers several factors, including alignment length, percent identity, and the number of paralogy methods that identify the paralog. Additional information for rank determination and alignment length are available to the users via a clickable help icon located next to those column headers.

The paralog section was released with Alliance version 6.0.0. Forthcoming updates will include the ability to sort and filter the table by column values and the availability of these data via our bulk downloads page. The existing tables on the gene pages for Function, Disease, and Expression all contain checkboxes for “Compare Ortholog Genes” that allow users to search across species for these features. We will add the additional checkbox “Compare Paralog Genes” to provide similar functionality for paralogous genes in a future Alliance release.

JBrowse sequence detail widget

A recent Alliance 6.0.0 release includes a new “Sequence Detail” section of all gene pages that uses JBrowse and JavaScript libraries to display an interactive widget that allows users to download DNA and amino acid sequences of genes in several possible configurations: genomic sequence highlighted with UTRs, coding and intronic regions, CDS regions, and translated protein for example (Fig. 3). In the next few releases, we will extend the functionality of the widget variant detail pages, where both the wild-type and variant sequences will be provided. When the variant occurs in the context of a protein coding gene, changes to the coding sequence and resulting translated protein will also be displayed and available for download.

Model organism BLAST

For more than 2 decades, some of the MOD members of the Alliance have hosted their own custom BLAST interfaces (Altschul et al. 1990; e.g. FlyBase Consortium 1999) that have allowed users to search custom databases related to those model organisms, e.g. subsets of related species or molecular clones, and display BLAST hits in Genome Browsers aligned with current gene models. We are now developing an updated and integrated Alliance BLAST, powered by SequenceServer (Priyam et al. 2019), that optimizes sequence analysis across model organisms. We have begun to update BLAST for the individual MODs. The new WormBase BLAST is now available online and can currently be accessed via the tools menu on wormbase.org. The results are linked to Genome Browsers and Alliance gene pages (Fig. 4). This tight connection allows users to navigate seamlessly between their BLAST results and the wealth of information available within the Alliance, enhancing the efficiency and depth of genetic research. For example, users can retrieve BLAST results for a sequence of interest and then easily navigate across Genome Browsers for different organisms, with a comparison to different tracks revealing how that sequence aligns with gene models, variants, and experimental tools (Fig. 5). From a project perspective, developing Alliance BLAST with a common cloud-optimized infrastructure will increase efficiency by reducing the cost of compute overhead and eliminating the need to manage separate MOD systems, which will then allow more focus on developing new functionality to support researchers. Our focus in the upcoming year is directed toward enhancing the user interface (UI), reflecting our commitment to providing an intuitive platform

Sequence Details ?

Transcript: Mode gene

```
>NM_001031689.3-gene
AGGCTTGCTGTGTCCTCGGCCCGCTCGGCGCGCC
GGCCGCCCTTACCTGCAGGCTCTTCTCCGCGCGG
CTGGGACCGGGCGCAGACAGACACTGGCCATGAC
AGCTGGACGTACGGGGCTGGTGTGCTGGCCATAT
TGGGCCCCAGACAGGTGAGCGCTGGGAGTCGGGTT
TTCCCTGTGCTGCTCCGCTCCCTCCCAATTCAGACTATTAGAAC
TCTGTAAGAAACCATCGGGATTTAAGTGGaaagagcacagggttgggggacCAAAGACCTGCCTGTTGTGCTAACTTGACCA
CAGACGAGTCTCCTACCTTTGGGCTTCAGTTTTGGACGATGATCTCCCAAGTTCTTTTACAGTTGAAAAATTTACTGATT
GCAGTTGCACCTCCGAAGTGGAGTAGTTGAAGGCATCTGaatgtcctcttttttttttttttttcgaagaAGATGCTCTGT
AGTCTctgtaaaatttaattttgaagactTAGTTCTCAAAAATGCAGTGGTGAATCCTCTTTTCGTCAGGTTAGAATT
TTAGTCTGTGCTGTGCTGCTGAGGAGTGAACCTCTCGGATTTTTGTTCTGTGATGTGCTTTCTTGAGGAGA
AGCAGCATCCATTCGCTTCAAAGGATTTATCAGAAGGGTTCacgaacaaaaaaaaaagaagaaagggttaGGAATCAGTC
CTGATCGAGTTCACGGGTCAGCCCTGATTTGGCTGTTGAACAGGATATTAAAGCCTAGAAGACAGATTGcagttcagag
aaagaaaaattgaggttagttattttgattatttagtAGGTCCTCCACTGCTAGAGATTTAGAATTTGAGTCACCATCCATAA
ATTCAAGTATAACTGTTGAGTGCCTCTCTTGTTTGGGATACTGGAGAGGAATAAAGACAGACAAGTTGCCAGTGTTCGTG
GAGCTTCTTACACTGCTGATGGGaaagatagattaaaaaacagcaataaataaataaaactggCATGTGAAGAAAACATA
```

- gene
- CDS
- cDNA
- protein
- genomic
- genomic +500bp up and down stream
- gene with collapsed introns
- gene with 500bp up and down stream
- gene with 500bp up and down stream and collapsed introns

Lowercase bases have been soft masked by NCBI Genomes to mark repetitive sequences.

Fig. 3. Sequence detail widget. Chosen views of a specific gene are readily available for copying as plain text or with highlights. 5' region of the human PLAA gene.

SequenceServer 2.0.0 Help & Support

BLASTN: 1 query, 1 database
[Edit search](#) | [New search](#)

Download FASTA, XML, TSV
[FASTA of all hits](#)
[FASTA of selected hit\(s\)](#)
[Alignment of all hits](#)
[Alignment of selected hit\(s\)](#)
[Standard tabular report](#)
[Full tabular report](#)
[Full XML report](#)

SequenceServer 2.0.0 using BLASTN 2.13.0+, query submitted on 2023-03-29 14:40:34 UTC
 Databases: C. elegans Genome Assembly (7 sequences, 100286401 characters)
 Parameters: task blastn, evaluate 1e-05, sc-match 2, sc-mismatch -3, gap-open 5, gap-extend 2, filter L;m;
 Please cite: <https://doi.org/10.1093/molbev/msz185>

Queries and their top hits: chord diagram

Query= Query_1 length: 487
[Graphical overview of hits](#) [SVG](#) [PNG](#)

Length distribution of matching sequences
Sequences producing significant alignments

# Similar sequences	Query coverage (%)	Total score	E value	Identity (%)
1. III	100	1308	0	100%
2. X	11	194	1.89*10 ⁻⁹	84.7%
3. I	23	249	6.58*10 ⁻⁹	87.5%
4. V	23	488	2.30*10 ⁻⁸	86%
5. IV	12	246	8.02*10 ⁻⁸	83.1%

hit 1, length: 13,783,801

[Select](#) | [Sequence](#) | [FASTA](#) | [Alignment](#) | [JBrowse](#) | [abi-1](#)
[Graphical overview of aligning region\(s\)](#) [SVG](#) [PNG](#)

a. Score: 879.53 (974), E value: 0, Identity: 487/487 (100%), Gaps: 0/487 (0%), Strand: + / -
 Query 1 CTGAAAAAATTTGCTTTTCGTGTTTTGACAAAACGTTTTCAAAAAAAAAAGGGAGCGAAAAATTCGACATAACTTATACAT 84
 Subject 5691938 CTGAAAAAATTTGCTTTTCGTGTTTTGACAAAACGTTTTCAAAAAAAAAAGGGAGCGAAAAATTCGACATAACTTATACAT 5691847

Fig. 4. Screenshot of results from the Alliance SequenceServer BLAST tool. The results have been enhanced relative to the default SequenceServer results page by the addition of links to Alliance JBrowse and to the corresponding gene page (in this case *C. elegans* abi-1) at the Alliance website for each BLAST hit.

for researchers in model organism genetics. We plan to produce more analysis tools as part of the evolving Alliance portal, thereby broadening the range of resources available for genetic research within the community.

AllianceMine

AllianceMine, a sophisticated, multifaceted search and retrieval tool that utilizes the InterMine software (Smith et al. 2012), offers a unified view of harmonized data, enabling advanced queries

across multiple species. For instance, gene lists can be processed as input and simultaneously query different annotations, such as “Show me genes associated with a (specific disease term)” (Fig. 6). The results from queries can be combined for further analysis and saved or downloaded in customizable file formats. Queries themselves can be customized by modifying predefined templates or by creating new templates to access a combination of specific data types. Thus, this powerful tool can be used in multiple ways, namely, for search, discovery, curation, and analysis.

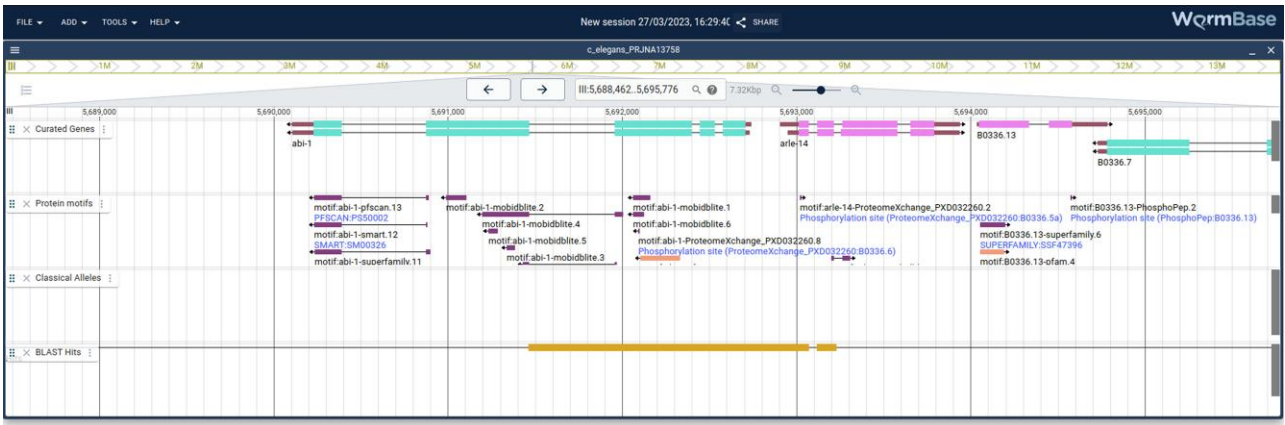


Fig. 5. Output of a BLAST search. After a user clicks on the JBrowse link for a BLAST hit, they are directed to the web service where they will see a track for the BLAST hit and how the hit aligns with other tracks.

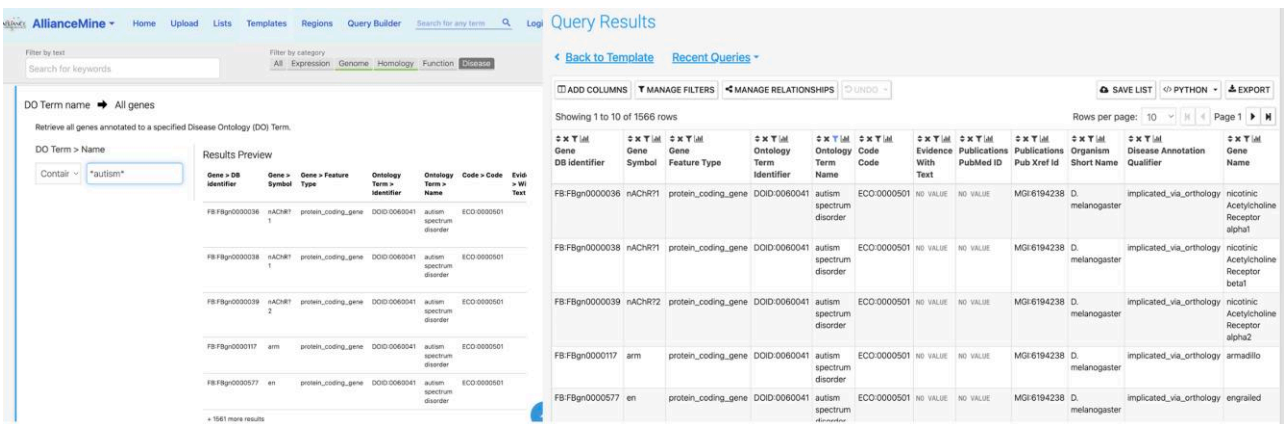


Fig. 6. AllianceMine example. Using a simple template, a disease ontology (DO) term, in this case “autism,” is chosen, and all genes associated with this DO term are returned in a downloadable table.

AllianceMine currently showcases harmonized data encompassing genes, diseases, GO, orthology, expression, alleles, variants, and FASTA formatted genome sequences. The tool also offers predefined queries or “templates” for cross-species searching. Continual optimization will ensure timely data synchronization with the main Alliance site, as well as integration of newly harmonized data types. Another aspect of improvement will be the addition of more templates, widgets, and precompiled lists, which can serve as logical input for templated queries.

SimpleMine

We designed SimpleMine for biologists to get essential information for a list of genes without any command-line or programming skill, or patience to learn the awesome power of AllianceMine discussed above. Users can submit a list of gene names or IDs to access more than 20 types of essential data with which they are associated. The results are 1 line per gene with detailed information separated by 4 types of separators: tab, comma, bar, and semicolon. Users can choose to display the output as HTML or to download a tab-delimited file. Alliance SimpleMine contains 10 species curated by the Alliance MODs. It provides easy gene name/ID conversion among MOD ID, public name (the commonly used name for public consumption), NCBI, PANTHER, Ensembl, and UniProtKB. Users can find summarized anatomic and

temporal expression patterns, variants, genetic, and physical interactions. Other essential gene information includes disease association and orthologs among all 10 species. The infrastructure of SimpleMine allows users to perform species-specific searches for lists of genes that take about 2 s to return results, or mixed-species searches that take about 10 s to complete.

Pathway displays with metabolites (GO Causal Activity Models)

We implemented a pathway display on Alliance gene pages that presents both GO Causal Activity Model (GO-CAM; Thomas *et al.* 2019) and Reactome pathway (Milacic *et al.* 2024) model. The display queries both the Reactome and GO application programming interfaces (APIs) and shows the number of pathways from each resource that contain the gene of interest. If a gene appears in multiple pathways, users can select which pathway to display. For the GO-CAM models, the viewer has been improved relative to previous releases of the Alliance website (Fig. 7). First, the layout has been improved to show clearly the overall causal flow through a pathway, from top to bottom and branching as necessary. Second, the viewer displays not only the activities of genes/proteins in a pathway but also metabolites, which is particularly useful for visualizing metabolic pathways. These metabolites may be either intermediates in a pathway or regulators of a

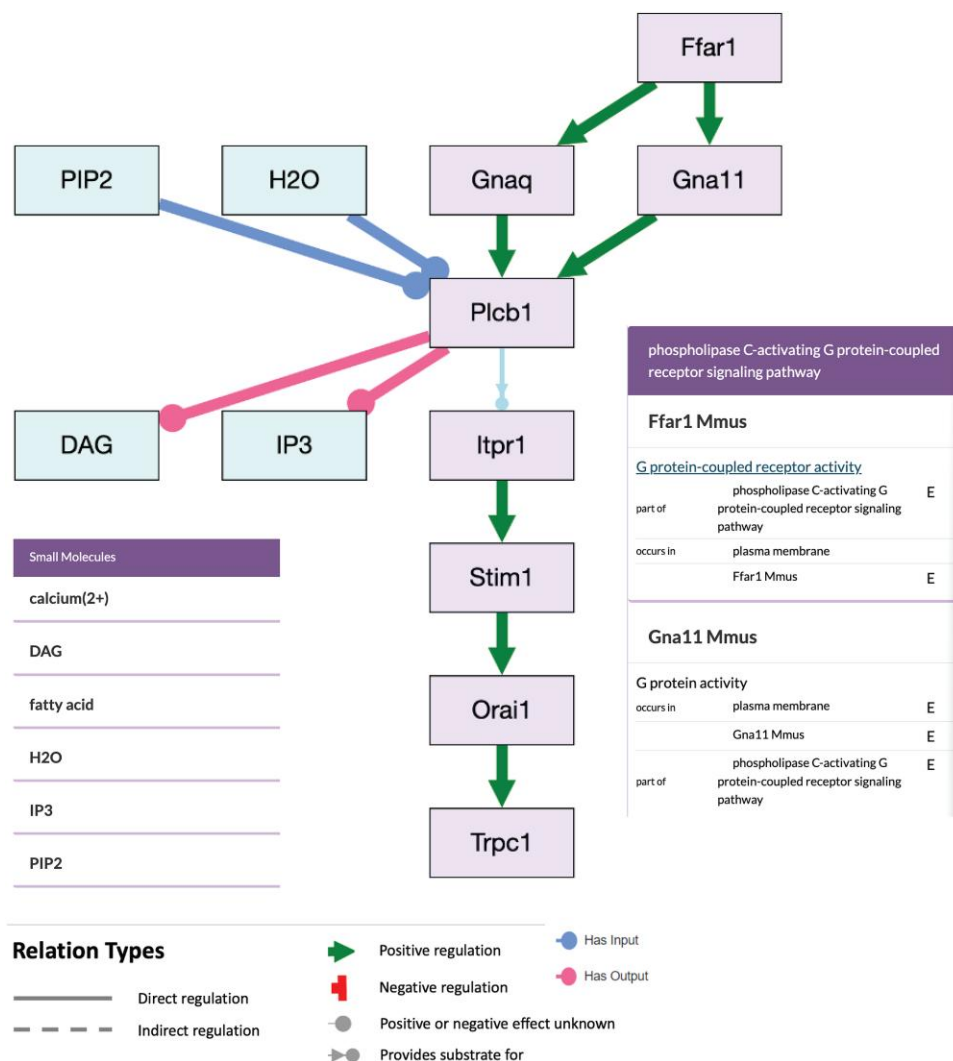


Fig. 7. Alliance pathway viewer. The pathway widget displays gene products (rectangles with gene names) and chemicals (rectangles with chemical abbreviations) and the flow of information and material between them (relations). These relations, shown in legend, indicate direct or indirect regulation that can be positive, negative, or of unknown effect direction. For metabolites that mediate the information flow between gene products, distinct shading distinguishes metabolites that are the inputs or outputs of a reaction.

protein activity. For signaling pathways, we distinguish between direct and indirect regulations and between positive, negative, or unknown effects.

Harmonized data models

The transition of data from individual MODs to the Alliance infrastructure requires data harmonization so that existing analogous MOD data classes (types/tables) can be loaded into Alliance databases using a consistent schema and language. The first step is for biocurators from each Alliance knowledge center to agree on which data classes are analogous and can be treated as a single, consolidated data class. The biocurators then align the properties (table columns) of the consolidated data class, including identifiers, types of values, and whether entity–property–value associations/triples require their own respective metadata and/or evidence records. To enable this process, the Linked Data Modeling Language (LinkML). We now have a standard workflow and common data modeling patterns that have streamlined the process, which we expect to complete in the next year. The LinkML specifications, authored in human-readable files, are

used to programmatically generate JavaScript Object Notation (JSON) schema specifications, which allow data quartermasters (DQMs) to move data to the persistent store. These specifications also inform curation software developers how to generate initial back-end (Java models and APIs) and front-end infrastructure (curation UI data tables and detail pages). Once DQMs have submitted their data files for a particular data class, the data are loaded into the persistent store and validated (see *Persistent store architecture* description below) and thus automatically populated into data tables and the curation interface. The data, having been harmonized, ingested, validated, and displayed to curators in the curation software, can now flow through to the public site according to the data pipeline described (see *Persistent store architecture* description below).

Many Alliance data classes have completely (or nearly completely) harmonized data models in LinkML (see https://github.com/alliance-genome/agr_curation_schema) including disease annotations, alleles, variants, expression annotations, and references. Although many other data classes have partially harmonized models, ongoing and future harmonization efforts will focus on completing harmonized models for the remaining

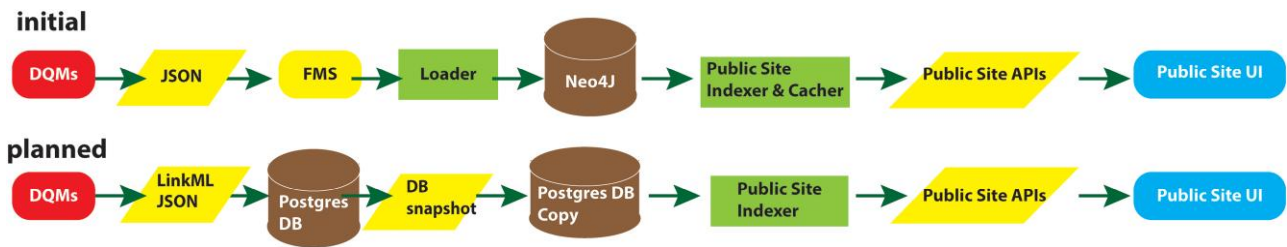


Fig. 8. Evolution of data flow. Graphical summary showing the design of short-term infrastructure initially deployed to support rapid delivery of unified data to the community and the planned production system. Red, data quartermasters at MODs; yellow, data; brown, database; green, transformations; blue, user interface.

curated data classes: genes, transcripts, proteins, nontranscribed genome features, affected genomic models (AGMs; strains, genotypes, and fish), phenotype annotations, molecular and genetic interactions, gene regulation annotations, high-throughput expression data set metadata (including for RNA-seq, single-cell RNA-seq, and proteomics data sets), species, reagents such as DNA clones and antibodies, images, persons, laboratories, companies, and various entity set classes like gene sets, which can be used for storing assay results and performing downstream analyses like ontology term enrichment, alignments, and other entity set processing calculations.

Persistent store architecture

We have designed a powerful database system that can handle most of the demands of our project including curation, analysis, and display of the data (Fig. 8). Specifically, we created a database using Postgres for long-term and persistent storage of Alliance curated data contributed by Alliance member MODs. In parallel to the existing (drop-and-reload) data pipeline (Alliance 2022), DQMs from each MOD now submit data according to our new LinkML schema in JSON format directly to the persistent store for ingestion, validation, and curation via create-read-update-delete (CRUD) operations enabled by a curation API library and Prime React UI. A data pipeline has been established to provide data from the persistent store Postgres database to our Alliance public website APIs and front-end web UIs and to other tools and services.

LinkML-based JSON files are ingested into Postgres with validation to ensure (1) recognition of submitted entities such as genes, alleles, AGMs (e.g. strains, genotypes), publications, experimental conditions, and ontology terms; (2) recognition of references to such entities in annotations and associations; (3) no entry of duplicate entities; and (4) proper handling of obsolete entities. Every file load is accompanied by a report (in Postgres and the curation UI) indicating (1) the recognized MD5 sum and size of the (uncompressed) file submitted; (2) the success or failure of the load; (3) the number of entities recognized in the submitted file; (4) the number of distinct entities loaded into Postgres; (5) the number and identity of entities (if any) that failed to load and the reason for the failure; (6) a link to download the submitted file; (7) the corresponding compatible LinkML model/schema version; and (8) the MOD data release version corresponding to the data in the file submitted. This information can be used by DQMs, curators, and developers to keep track of the fidelity of the data transfer and troubleshoot any issues that arise. Ontology (and other external resource) loads are updated nightly to ensure that the latest versions of such data are current. The source of truth for MOD data will be transitioned over to the Alliance infrastructure in phases,

beginning with a few data types from a few MODs and expanding over time to eventually include all (relevant) data types from all participating MODs; as part of this process, legacy issues with data are cleaned up.

To enable CRUD operations on persistent store data, curation APIs and a curation UI accessible with Okta authentication have been implemented (Fig. 9). Curators can now access data tables for the following data types: genes, alleles, variants, AGMs (e.g. strains, genotypes), publications [accessed via Alliance Bibliographic Central (ABC) APIs], experimental conditions, constructs, disease annotations, molecules [not already managed by Chemical Entities of Biological Interest (ChEBI)], ontology terms, and controlled vocabularies and their terms. CRUD operations have been fully enabled for disease annotations, experimental conditions, and controlled vocabularies, read-update operations have been enabled for alleles and variants, and read operations are enabled for the remaining data types. Work is underway to fully enable CRUD operations on all remaining data classes and their attributes including new data tables for transcripts, proteins, other (nongene) genome features, expression annotations, phenotype annotations, molecular interactions, genetic interactions, gene regulation annotations, antibodies, images, and more. In addition to data tables presenting all entries of a particular data class, the curation tool also has individual entity detail pages (for example, see an allele detail page at <https://curation.alliancegenome.org/#/allele/MGI:6446761>) for data entry and editing on a dedicated web page for 1 particular entity. The curation tool also enables user-specific and MOD-specific custom user settings and preferences to provide a UI most compatible with individual curators' workflows.

In the next year, the curation tool will include batch creation of data entities (e.g. annotations, reagents), batch editing, data history inspection and auditing, undo and review of latest changes, publication constraints (constrain data view and entry to publication currently being curated), customizations and MOD default settings for new entity creation and detail pages, incorporation of data entity and topic tagging information from the ABC literature store (see below), and incorporation of artificial intelligence (AI)/machine learning (ML) into the curation workflow.

For releases of persistent store data to the Alliance public website, Postgres database snapshots are taken and sent to a separate Postgres instance that feeds the data via the curation APIs (instantiated as a library) into the public site indexer where various data filtering and transformations occur before making those processed data available to our public website APIs via our Elasticsearch index. The Alliance public website UI, using existing UI infrastructure, is then modified or created to accommodate the data now flowing from the persistent store database.

Affected Genomic Models Table						
Curie ↑↓	Name ↑↓	Sub Type ↑↓	Taxon ↑↓	Data Provider ↑↓	Updated By ↑↓	Date Updated ↑↓
MGI:3720678	Tg(THY1-APP...	genotype	Mus musculu...	MGI		
WB:WBStrain000	FGP29	strain	Caenorhabdit...	WB		
WB:WBStrain000	IE4314	strain	Caenorhabdit...	WB		
MGI:5008182	Akr1a1 ^{Gt} (OST ...	genotype	Mus musculu...	MGI		
MGI:6492714	Atg7 ^{em1} (IMPC ...	genotype	Mus musculu...	MGI		
ZFIN:ZDB-FISH-1	crb2a ^{m289/m...}	fish	Danio rerio (...	ZFIN		

Fig. 9. Alliance curation tool. Screenshot of the Alliance curation tool interface showing an example of curated annotations of AGMs managed in the persistent store.

Security, stability, and backups

All services and data provided by the Alliance to its community are hosted on Amazon Web Services (AWS). This provides us with industry leading availability of up to 99.99% on services like EC2, which we use to host our virtual servers. We use additional AWS-managed services such as Elastic Beanstalk for application deployment, AWS Relational Database Service for hosting our relational (Postgres) databases, and Amazon OpenSearch Service for hosting our search indexes, which all provide automatic updates and maintenance for increased reliability. All files hosted at the Alliance of Genome Resources are stored in S3 buckets, which ensures industry leading durability and availability. Furthermore, we make daily backups of our relational databases and have processes in place that enable easy restore of those backups in case of failure or data corruption. All search indexes are derived from the persistent relational database and can be regenerated at any moment when required.

We make use of separated subnets between public-facing and private systems, and only services requiring public access are given public IP addresses, ensuring that public-facing services such as our curation interface can be accessed by our curators worldwide (through Okta Authentication), although the supporting back-end services such as the supporting databases can be kept private. Access to all services is furthermore restricted to allow access only to the required ports and services through the use of AWS Security Groups to control the allowed network traffic. AWS IAM users, groups, and roles are used to control the allowed AWS operations and access among Alliance developers. In all cases, the principle of least privilege is applied, so that the potential attack surface is reduced to a minimum (for example by not granting blanket AWS admin permissions to developers who do not have an AWS admin function). Access keys to any system can be revoked when misuse of those access keys is detected. We also configured our github repositories to be scanned automatically for accidental secret credential leakages through the use of GitGuardian software.

Literature acquisition

We designed and are implementing a literature system, ABC, that will support curation and, in the future, end users. The ABC supports the tasks of reference acquisition, triage, and curation workflow management. Specifically, the ABC is an ecosystem of online tools and supporting Alliance databases that manage all references and related metadata that are “in corpus” for the member MODs.

Literature acquisition at the Alliance begins with automated, organism-specific PubMed queries to retrieve candidate references for each MOD’s corpus. References matching the search criteria are then added to the ABC by assigning an Alliance reference identifier and importing associated bibliographic information to the database. Subsequently, curators manually sort references as either “in” or “out of corpus” based on the curation policies of the MOD and eliminate any false positive results from the initial search. While many thousands of papers are published each year, only some have information that is currently curated. For example, in 2022, the curatable literature size after triage was 3,181 for ZFIN, 3,221 for SGD, 2,130 for FlyBase, 1,419 for WormBase, and 437 for Xenbase. Once references are sorted, they enter MOD-specific curation workflows supported by task-specific ABC curator interfaces to, for example, add reference files, manually tag references with specific entities (e.g. genes, alleles, and data types) and topics (e.g. phenotypes, anatomic expression) using the Alliance Tags for Papers (ATP) ontology, and merge duplicate references. In addition to adding reference files manually, the full text of “in corpus” references included in the PubMed Central (PMC) open access set is also automatically downloaded. Curators may also use the ABC to add non-PubMed references. An additional key feature of the ABC is a search interface that allows curators to retrieve references based on various criteria including their in/out of corpus status, bibliographic data, and publication data range, if desired. Reference acquisition functionality can easily be extended to integrate additional MODs into the Alliance infrastructure.

To facilitate reference data exchange between the Alliance and MOD databases, the MODs provide a mapping file that associates MOD reference Compact Uniform Resource Identifiers (CURIEs) with PMIDs, e.g. ZFIN:ZDB-PUB-181026-2 - PMID:30352852. The MODs also provide reference CURIEs and data for references not included in PubMed but used by the MOD, such as internal curation references and those published in a journal not yet indexed at PubMed.

Over the past 25–30 years, Alliance member databases have independently developed methods to acquire, triage, and curate their respective literatures. Having implemented a common literature curation interface, database, and full-text acquisition system, the ABC is now poised to expand its functionality by incorporating ML methods developed by, and in production for, a subset of Alliance members to all groups. For example, automated pipelines that recognize entities (e.g. genes, alleles, and strains) as well as data types (e.g. phenotype, genetic interactions) can be developed for new groups with results stored centrally in the Alliance literature database. Incorporating more automated methods will allow faster association of the published literature with relevant biological concepts, information that can be displayed on future Alliance reference pages while the papers await detailed full curation. Centralized literature infrastructure will also support other curation pipelines, such as community curation by authors, which can then be more readily implemented for additional Alliance member communities, thus providing another avenue by which curated data can be swiftly included in the Alliance. Lastly, the common literature tool will allow Alliance biocurators to coordinate curation of multispecies references that will provide users a facile way to find and view cross-species research exploiting the strengths of each Alliance model organism, a primary goal of the Alliance.

Textpresso

Textpresso is a full-text literature search engine that gets power from its single-sentence scope, focus on a specific model organism (or topic), and categories of semantically or biologically related terms (Fig. 10; Müller et al. 2004, 2018). It has been used extensively by WormBase and SGD curators, as well as *C. elegans* and *Saccharomyces cerevisiae* researchers in addition to other MODs (Van Auken et al. 2012; Bowes et al. 2013).

The Alliance is committed to creating Textpresso instances tailored to the unique needs of each member database, all of which will be managed within the Alliance software ecosystem and connected to the ABC as a single reference data source. This will reduce the overhead of managing Textpresso at individual MODs while also simplifying development and deployment of new features. Users will benefit from simplified access to Textpresso from the Alliance website. We also plan to integrate Textpresso searches further into specific Alliance web pages such as gene or allele pages. Users will be able to obtain additional references to biological entities through Textpresso searches, adding information from potentially noncurated literature to the list of curated references currently linked on those pages. Textpresso will be available to Alliance biocurators and to the general public through MOD-customized websites and via APIs for programmatic access.

AI

The Alliance member MODs have a track record of implementing ML tools to enhance literature triage and curation efficiency. Notable examples include RGD's early adoption of standard software architectures such as Unstructured Information Management Architecture (UIMA, an Apache.org project) and



Fig. 10. Textpresso for SGD literature at the Alliance (<http://sgd-textpresso.alliancegenome.org/tpc/search>).

the development of the OntoMate system (Liu et al. 2015), an ontology-driven literature search engine, as well as WormBase's creation of Textpresso (Muller et al. 2004) and document classifiers for paper triage.

The rise of large language models (LLMs), such as BERT (Bidirectional Encoder Representations from Transformers) and ChatGPT, has transformed the natural language processing (NLP) landscape, but questions about their accuracy and “hallucinations” remain. The Alliance is developing LLMs for tasks such as document classification, named entity recognition (NER), sentence classification, computationally assisted triage, and curation and to build a natural language query system to simplify access to its underlying structured data.

Alliance members have developed AI/ML classifiers for determining with high accuracy whether papers returned from automated PubMed queries should be kept in their corpus or discarded (Jiang et al. 2020) and classifiers that can determine whether specific data types relevant for curation are present in a document (Fang et al. 2012). The Alliance is developing a central solution by providing these types of classifiers to all members.

Efforts are also underway to improve existing species-specific entity extraction and classification models, with a focus on incorporating human feedback in the loop and continuously training models based on data validated by professional biocurators and community curators. A centralized interface for “topic and entity tag” addition and validation during literature triage and curation is under development as part of the ABC. The interface allows curators to associate tags with publications and at the same time validate (or invalidate) results extracted from AI/ML methods. This interface will streamline the collection of valuable training and testing sets and will allow a more systematic approach to the creation and comparison of different AI/ML models.

Furthermore, the Alliance is adopting Evidence and Conclusion Ontology (ECO) terms to record systematically the type of evidence, e.g. neural network method evidence, and assertion method, e.g. automatic assertion, used for reference flagging and triage. This is especially relevant for topic and entity tags. Using ECO terms aligns with FAIR data principles and offers transparency in curation workflows.

APIs

APIs are a key component of Alliance Central's data service infrastructure for rapid, modular software development. We currently support a dozen APIs with hundreds of endpoints (Fig. 11). New APIs will be added as data harmonization and modeling of

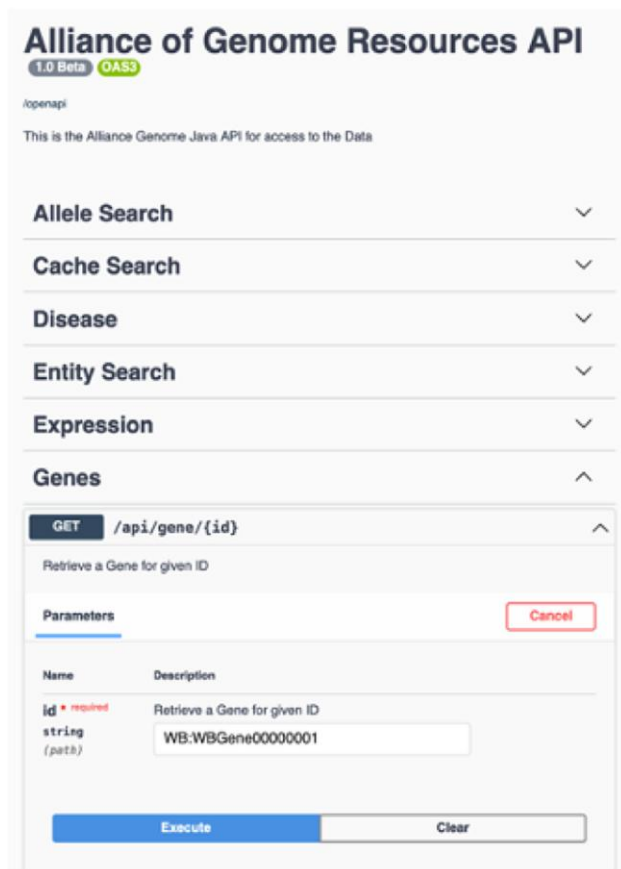


Fig. 11. Swagger interface for the Alliance APIs.

additional data entities are completed. We will expand public site APIs to generate all data needed for SimpleMine, AllianceMine, etc. from single endpoints. Current APIs include public site APIs (`agr_java_software` in the GitHub repo) and APIs available from a public Swagger UI page. Because the public APIs support only GET endpoints, they do not require authentication. All APIs that support both GET and PUT/POST/DELETE endpoints do require authentication. Some of the key API endpoints available at <https://www.alliancegenome.org/swagger-ui/> are `gene-summary`, `gene-disease`, `gene-interactions`, `homologs-species`, `allele-phenotypes`, `expression ribbon-summary`, etc.

Data preservation in external repositories

The Alliance of Genome Resources is committed to the long-term preservation of digital objects (annotations) and resources (e.g. ontologies and software) that are central to the management and integration of functional knowledge about the genomes of diverse model organisms. As part of this commitment, the annotations and resources generated by Alliance members are integrated into many long-standing external public bioinformatic resources (e.g. Ensembl, UniProt, and NCBI). Distribution of Alliance annotations from multiple sources provides a degree of redundancy that contributes to data stability and preservation. Alliance maintained ontologies and annotations are also deposited into third-party repositories that fulfill Open Science principles (see below). Leveraging community repositories ensures the data products and resources remain accessible to the research community even if the Alliance and/or its members cease operations.

Ontologies that Alliance members maintain are also available from long-term repositories including the OBO Foundry (<https://obofoundry.org/>) and Zenodo (zenodo.org).

Annotations related to gene expression, function, phenotype, disease associations, etc. that are generated by Alliance members and are available on the Alliance Data Downloads page are archived in Zenodo. Software developed as part of the Alliance of Genome Resources knowledge commons platform is available from GitHub (<https://github.com/alliance-genome>).

The external repositories used by the Alliance of Genome Resources include the *OBO Foundry* that was established in the early 2000s as a community-based initiative for development and maintenance of biological and biomedical ontologies using standardized practices. The Foundry is the ontology repository of choice for the Alliance because it is widely recognized as an authoritative source of well-maintained ontologies for biology and biomedical research. *Zenodo* is a general purpose repository maintained by CERN (European Council for Nuclear Research) for storing and sharing documents, data, and other digital research materials across many disciplines. *Zenodo* is a repository of choice for the Alliance, in part, because of the commitment by the European Commission to support *Zenodo* as long as CERN exists.

Outreach and interactions

The Alliance help desk

We established a common help desk email address (help@alliancegenome.org) that is featured prominently on the Alliance website header and footer under “Contact Us.” All inquiries submitted using this email are logged as tickets in the Alliance Jira software system. Members of the User Support Working Group respond to user questions and inquiries in a timely manner, typically within 48 h. Time to resolve user inquiries depends on the nature of the question or request. The Jira system tracks open tickets, forward tickets, tracks their active/resolved status, and classifies them by subject. We use the information, in part, to evaluate the design and utility of our UIs. For example, if particular questions or subjects arise frequently, we reevaluate the design and wording of the search form and/or results display that caused user confusion.

Online documentation

We provide extensive user documentation about using the Alliance data resources under the Help menu on the homepage (<https://www.alliancegenome.org/help>). The online documentation provides guidance on such topics as how to use the search functions, defines acceptable field parameters, and provides explanations of the displayed results. The User Support Working Group also works closely with the User Interface Working Group and the Developers to craft text for tooltips displayed on UIs.

Frequently Asked Question pages

The Frequently Asked Question (FAQ)/Known Issues page provides answers to commonly asked questions about the Alliance and also describes any known issues associated with a particular software release. The link to the FAQ page is featured prominently on the Alliance home page under the Help menu.

Illustrated tutorials and videos

We maintain several types of tutorial options that are accessible from the Help menu (<https://www.alliancegenome.org/tutorials>). The most requested types of tutorials are illustrated guides with screenshots on how to use various features of the Alliance web

The screenshot shows the Alliance of Genome Resources community forum home page. At the top, there is a navigation bar with 'Sign Up' and 'Log In' buttons, and a search icon. The main content is organized into two columns: 'Categories' and 'Latest'.

Categories:

- Alliance of Genome Resources** (29 posts): News and Announcements, Site Feedback, Data Discussion, General Discussion.
- Job Postings** (1.1k posts): Open positions and job announcements. Topics: Flies, Frogs, Mice, Rats, Worms, Yeast, Zebrafish, Other.
- Positions Wanted** (11 posts): Are you a graduate student, postdoc, or young faculty member looking for a position? Post your details and requirements here. Topics: Flies, Frogs, Mice, Rats, Worms, Yeast, Zebrafish.
- Meeting Announcements** (132 posts): Announcements and discussions about upcoming meetings. Topics: Flies, Frogs, Mammals/Human, Worms, Yeast, Zebrafish.
- Model Organism: Flies** (8 posts): Discussion related to *Drosophila melanogaster*. Topics: Reagents, FlyBase.
- Model Organism: Frogs** (4 posts): News and Announcements, Scientific Discussion, Stocks.

Latest:

- Welcome to Discourse** (0 replies, Nov '20): Worms.
- MMRRC Newly Available Strains July 2023 & MMRRC Newly Accepted Strains July 2023** (0 replies, 1d): Stocks.
- Drug-induced shrinkage of nematodes** (0 replies, 1d): Scientific Discussion.
- How to enter data in Kaplan Meier graph?** (0 replies, 1d): Methods & Reagents.
- Multi-purpose embryo extracts-Freon Free protocol** (0 replies, 4d): Methods & Reagents.
- Xenopus Developmental Biology 1-week course Sept 11-15, 2023** (0 replies, 4d): Frogs.
- Project Manager, Rare Disease Translational Center at JAX** (0 replies, 5d): Job Postings.

Fig. 12. Alliance community forum home page.

portal. When new functionality is released, we post to social media channels and issue “Tweertorials.” Short video tutorials are disseminated through the Alliance YouTube channel.

Alliance user community forum

The Alliance supports a centralized community discussion board implemented in Discourse (<https://community.alliancegenome.org/categories>; Fig. 12). Each model organism represented in the Alliance is represented as its own Discourse category with model organism-specific threads for news, discussion, and reagent information. The forum also includes categories for job postings, meeting announcements, and general information about the Alliance of Genome Resources. Alliance members with existing online community forums are migrating users to the Alliance Central forum.

Users are not required to register to access the forum but must register to post messages, questions, and announcements. On average, ~1,000 users a day access the forum. Posts include jobs open and sought, news, meeting announcements and discussion of research approaches, reagents, and interpretation.

Social media

In addition to a News and Events header that links to software release notes and other Alliance Central updates, the Alliance uses standard social media venues to engage with the user community, including Facebook (www.facebook.com/alliancegenome/), Twitter (now, X; twitter.com/alliancegenome), Mastodon (<https://genomic.social/@AllianceGenome>), and Bluesky (<https://bsky.app/profile/alliancegenome.bsky.social>).

Prospects and challenges

The tail of not-yet harmonized data

One challenge in the central Alliance infrastructure providing support for the union of MOD and GO features is the many unique

data set displays and tools that have evolved in the individual MODs over 2 decades. Among the 8 resources, this comprises 150 years of branch length! Although horizontal tool transfer has occurred, it is not complete. We are taking a few approaches to this problem. In some cases, where the data are stand-alone, we will simply move the data and code to the Alliance. In the short term, we will likely run tools off their existing servers. As tools age out, we will evaluate whether there is a broader mandate for that feature and, if so, implement it in the context of the Alliance.

There are types or aspects of our data that can be harmonized but have not yet been so. We adopted LinkML to help with harmonization because it provides a common language to represent structured data. The use of this language has spread to the point where our progress on harmonization is much more rapid.

AI

As discussed above, we are actively considering AI/ML applications throughout the project. Our practical approach is driven by us being subject matter experts. Because we have relied on human expert curation, we are in a unique position to evaluate and use the output of various AIs. Future plans include development of tools for creating training sets and a model manager for tracking ML models' performance. Integration with specialized bio-curation tools such as OntoMate and Textpresso is part of the strategy, with a vision of harmonizing AI/ML solutions across member sites.

We will also explore the use of AI/ML in gene function summarization. Included on gene pages at the Alliance are short textual gene summaries based on curated and structured data annotations that provide users a quick overview of gene function. The current automated system for generating gene summaries has produced more than 160,000 summaries (Alliance version 6.0.0) for 9 species, including humans (Kishore et al. 2020). However, to increase the coverage of genes further, we will explore the use

of LLMs. This is especially relevant for less-studied genes with few curated, structured data and for scaling and upkeep of the summaries to match the rate of new gene data from publications. Leveraging LLMs to generate gene summaries for less-studied genes, particularly those with limited curated data, offers the advantage of automatically uncovering relevant publications that may not have been previously curated. In principle, AI might be able to enhance or replace the automatically generated textual gene summaries for both well-studied and less-studied genes.

We will use prompt engineering and fine-tuning of LLMs to improve accuracy of the generated summaries. As part of a continual improvement process, we will ask professional biocurators to evaluate summaries, and we will develop a scoring system based on several features such as readability of summaries, inclusion of key gene data, and checking for inaccurate and false data. To improve and keep gene summaries up to date, we plan to retrieve newly published articles that contain gene data that were not available when the LLM was trained and add extracted relevant text from the identified articles to the LLM prompt. To do so, we will use tools such as Textpresso (Muller et al. 2004) and OntoMate (Liu et al. 2015).

Community curation

Some Alliance MODs employ community curation pipelines to engage authors in curation of their papers. For example, FlyBase utilizes the Fast-Track Your Paper (FTYP; Bunt et al. 2012; Larkin et al. 2021) tool that allows users to curate scientific papers, identify data types, and associate relevant genes with the reference. Authors using FTYP to ensure their papers appear quickly on the FlyBase website help highlight data needing manual curation and prioritize their publication for further curation.

Similarly, WormBase developed ACKnowledge (Author Curation to Knowledgebase; Arnaboldi et al. 2020), a semiautomated curation tool that lets authors curate their publications with the help of ML. Authors receive an email with a link to a form prepopulated by document-level classifiers that identify data types and several NER pipelines that extract lists of entities. Authors can correct and validate the extracted data using the form and submit curated information to WormBase. We will continue to provide these services to our community and develop a unified infrastructure that will help expand the service to other member communities.

Several Alliance members also collaborate with publishing groups, such as microPublication Biology (<https://www.micropublication.org/>) or the Genetics Society of America (<https://genetics-gsa.org/publications/>), to streamline prepublication data integrity verification and curation by curators and authors, enabling MODs to quality check and work with authors to correct data reporting before publication and promptly incorporate it into Alliance Knowledgebases upon article publication.

Dealing with satellite genomes and genetic models

In addition to the core genomes and associated data, our resources store and present information about the genes and genomes of relatively closely related organisms. For example, WormBase includes some genetically studied nematodes such as *Caenorhabditis briggsae* that benefit from the rich data models typical of *C. elegans*. Genetic screens and positional cloning (Inoue et al. 2007; Sharanya et al. 2012), CRISPR editing (Cohen and Sternberg 2019; Cohen et al. 2022; Ivanova and Moss 2023), and transcriptomic analyses (Jhaveri et al. 2022) are now routinely done in this species. For the Alliance to take on this responsibility

of WormBase, we need to support such satellite model organisms. Our plan is to support community gene structure annotation (e.g. for *Drosophila*, Sargent et al. 2020; for *C. elegans*, Moya et al. 2023) using the Apollo curation system designed specifically for such activity (Dunn et al. 2019).

High-throughput expression data and single-cell RNA-seq plans

We harmonized high-throughput expression metadata of mouse, rat, yeast, worm, fly, and zebrafish. Users can browse them with species, assay type (microarray, RNA-seq, tiling array, and proteomics), tissue, sex, and curated categories. We plan to add single-cell RNA-seq as a new assay type, making such data sets easily identifiable within our collection, with links to other resources, including Gene Expression Omnibus, EBI single-cell RNA-seq Expression Atlas, and CZI CellxGene, and to display the information above, we will implement a content-rich expression detail page that will provide a unified way to access all expression data associated with a specific gene, including link outs to other sources and MOD-specific single-cell RNA-seq gene expression graphs (Fig. 13).

Disease portal(s)

Providing users with ready and easy access to curated and harmonized model organism disease data and tools is crucial to accelerate research related to the pathogenesis of human disease. The Alliance has a wealth of disease-relevant data from 8 model organism species and human data, such as genes, alleles and variants implicated in disease, genotypes and strains that serve as disease models, and related data such as modifiers (herbals, chemicals, small molecules, etc.) that ameliorate or exacerbate the disease condition and may serve as candidates for potential drug development. To provide an easy entry point for clinical researchers and human geneticists to access the consolidated data and tools, we are in the process of designing and implementing a topic-specific resource—an Alzheimer's disease (AD) portal that will serve as a paradigm for other disease portals (Fig. 14). The AD portal will include orthologous genes in animal model systems, models with a mutation orthologous to one in a patient group, models with a specific set of phenotypes, and/or modifiers that have been shown to alter the disease condition. Building on the experience and pages developed for the AD portal, we will expand this paradigm to other disease portals. In addition to the specific disease portals, we also plan to provide "compare" functionalities across diseases. Features planned for the disease portal with AD as an example include a home page with an overview of the data in the portal, an autocomplete search box, links to other AD resources, and a list of the most recent papers from PubMed and/or from the ABC store (see example portal page below). The pages in the portal will be modeled on existing pages at the Alliance and will include gene summaries, alleles and variants, phenotypes, gene interactions, pathways, biological processes (based on GO), and gene expression. We also plan to provide visualizations of data analysis, for example, diseases that share genes and protein interactions that may point to common underlying molecular mechanisms. Up-to-date data sets, e.g. genes, strains, and modifiers (drugs, chemicals, herbals, etc. shown to either ameliorate or exacerbate phenotypes), will be available as downloadable files. Disease-specific data sets will also be available for query from AllianceMine. We will also provide up-to-date links to disease-specific literature and search capabilities through literature search engines such as the Textpresso instance dedicated to AD (<http://alzheimer.textpressocentral.org>), corpus

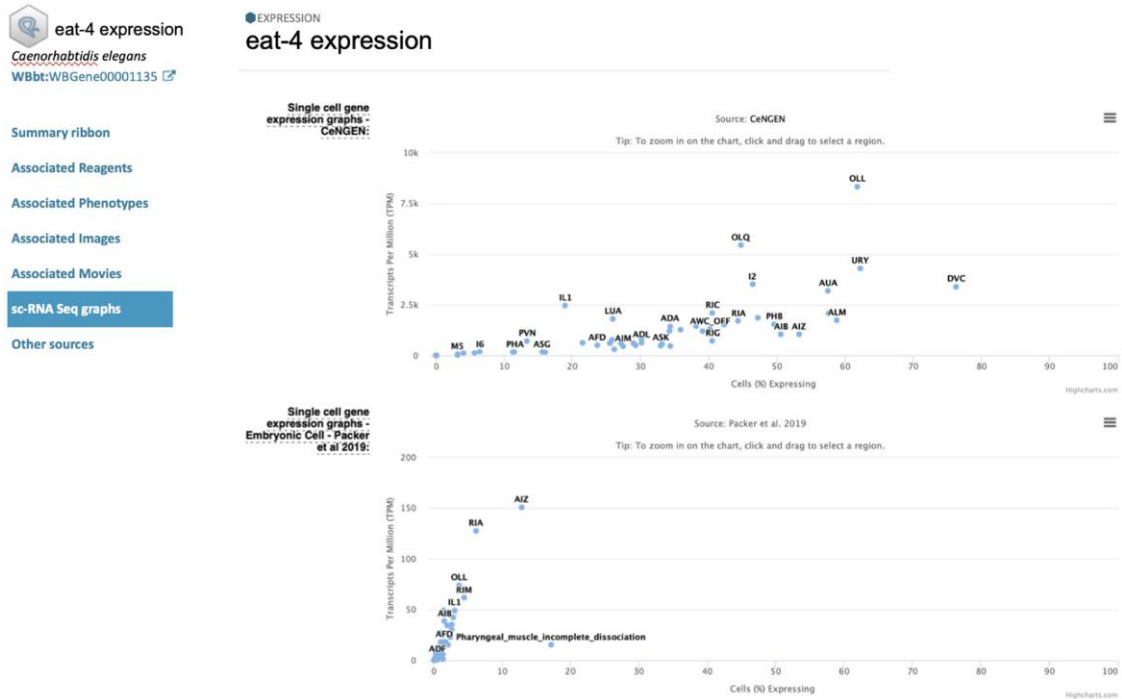


Fig. 13. Mockup of an expression detail page. This example shows one of the current features of [WormBase](#)—single-cell data from 2 studies—displayed on what will be part of an Alliance gene expression detail page.

Alzheimer's Disease Portal

Home Page

Alzheimer's Disease Portal
Bringing the power of model systems to the biomedical community

Tutorials

Access Data

3,000
Models

1,221
Genes

1,000
Alleles

96,000
Publications

9 model
species

Latest Papers (powered by Textpresso)

Yang et al., Long-term exercise pre-training attenuates Alzheimer's disease-related pathology in a transgenic rat model of Alzheimer's disease. *Genes*. 2022. Jun;44(3):1457-1477. PMID: 35229257

Shao et al., Treatment of Alzheimer's disease with framework nucleic acids. *Cell Prolif.* 2020. Apr;53(4):e12787. PMID: 32162733

Fang et al., Mitophagy inhibits amyloid- β and tau pathology and reverses cognitive deficits in models of Alzheimer's disease. *Nat Neurosci.* 2019 Mar;22(3):401-412. PMID: 30742114

Kiper et al., Use of Zebrafish Genetic Models to Study Etiology of the Amyloid-Beta and Neurofibrillary Tangle Pathways in Alzheimer's Disease. *Curr Neuropharmacol.* 2022 Mar 4;20(3):524-539. doi: PMID: 34030617

Community Resources

[Alzheimer's Society](#)

[Alzheimer's.gov](#)

[Alzheimer's Association](#)

[Alzheimer's Foundation of America](#)

...

Need Help? Contact Us: help@alliancegenome.org

Access Data Page

Category

Associated genes

Search Gene

APP

PTEN

....

Disease subtypes

AD1

AD2

...

Model species

Rat

Mouse

Zebrafish

...

Genes

Gene	Accession	Model	Species	Accession	Model	Species
APP

Alleles

Allele	Accession	Model	Species	Accession	Model	Species
...

Models

Model	Accession	Model	Species	Accession	Model	Species
...

Publications

PMID	Author(s)	Title	Journal	Year	Category
PMID 1234	J. Doe	The quick brown fox	Nature	2022	View
PMID 4567	J. Doe	Jump over the...	Science	2020	View

Fig. 14. Mockup of the AD portal showing the home page and the data access page. These views illustrate the type of information that will be available with a disease focus.

size—96,000 papers). Not all papers are curatable by the MODs given their extensive but not comprehensive data models, and thus, literature search will remain important.

The Alliance in the ecosystem of knowledgebases

The Alliance has a unique and complementary role relative to other informatics resources that support comparative biology.

For example, NCBI's new Comparative Genomics Resource (CGR; [Bornstein et al. 2023](#)) focuses on developing analysis tools and resources for *sequence-based* genome comparisons across a large number of species, and the Alliance focuses on standardized annotations, harmonized biological concepts, and comparison of *biological knowledge*. The CGR supports comparative sequence analysis for all eukaryotes whereas the Alliance is primarily

focused on model organisms used widely in biomedical research. These model organisms have a tremendous amount of highly valuable genetic, transgenic, and phenotypic data generated with multiple types of assays and are uniquely represented by the Alliance Knowledge Centers. The CGR uses the standardized gene summaries from the Alliance and follows nomenclature and ontology standards developed and maintained by Alliance members. For sequence analysis, the Alliance leverages sequence-based analysis tools developed and maintained by the CGR. Resource developers by and large appreciate the magnitude of the tasks we face in order to provide researchers with the information they need and strive to fill in the many gaps in services.

Data availability

All the data underlying this article are available at alliancegenome.org.

Acknowledgments

We thank our multiple communities for their patience and feedback about the prospect of the Alliance and their love of their own MODs. We also thank the members of our Scientific Advisory Board (Gary Bader, Alex Bateman, Helen Berman, Shawn Burgess, Andrew Chisholm, Phil Hieter, Brian Oliver, Calum Macrae, Titus Brown, Abraham Palmer, and Michelle Southard-Smith) for cogent advice and NHGRI Program Staff (Sandhya Xirasagar, Ajay Pillai, Valentina di Francesco, Sarah Hutchison, and Helen Thompson) for guidance.

Funding

The core funding for the Alliance is from the National Human Genome Research Institute and the National Heart, Lung and Blood Institute (U24HG010859). The curation of data and their harmonization is supported by National Human Genome Research Institute grants U24HG002659 (ZFIN), U24HG002223 (WormBase), U41HG000739 (FlyBase), U24HG001315 (SGD), U24HG000330 (MGD), P41HD064556 (Xenbase), U24HG011851 (Reactome + GO), and U41HG012212 (GO Consortium), as well as grant R01HL064541 from the National Heart, Lung and Blood Institute (RGD), P41HD062499 from the Eunice Kennedy Shriver National Institute of Child Health and Human Development (GXD), and the Medical Research Council UK grant MR/L001020/1 (WormBase). Additional effort was supported by the U.S. Department of Energy (DOE DE-AC02-05CH11231). Curation tools are supported in part by the National Library of Medicine NLM R01LM013871.

Conflicts of interest

The author(s) declare no conflicts of interest.

Literature cited

- Alliance of Genome Resources C. 2022. Harmonizing model organism data in the Alliance of Genome Resources. *Genetics* 220(4): iyac022. doi:10.1093/genetics/iyac022
- Altenhoff AM, Train CM, Gilbert KJ, Mediratta I, Mendes de Farias T, Moi D, Nevers Y, Radoykova HS, Rossier V, Warwick Vesztrocy A, et al. 2021. OMA orthology in 2021: website overhaul, conserved isoforms, ancestral gene order and more. *Nucleic Acids Res.* 49(D1):D373–D379. doi:10.1093/nar/gkaa1007.
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol.* 215(3):403–410. doi:10.1016/S0022-2836(05)80360-2.
- Anderson WP; Global Life Science Data Resources Working Group. 2017. Global life science data resources working, data management: a global coalition to sustain core data. *Nature* 543(7644):179. doi:10.1038/543179a.
- Arnaboldi V, Raciti D, Van Auken K, Chan JN, Müller HM, 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.
- Bornstein K, Gryan G, Chang ES, Marchler-Bauer A, Schneider VA. 2023. The NIH Comparative Genomics Resource: addressing the promises and challenges of comparative genomics on human health. *BMC Genomics* 24(1):575. doi:10.1186/s12864-023-09643-4.
- Bowes JB, Snyder KA, James-Zorn C, Ponferrada VG, Jarabek CJ, Burns KA, Bhattacharyya B, Zorn AM, Vize PD. 2013. The Xenbase literature curation process. *Database (Oxford)* 2013:bas046. doi:10.1093/database/bas046.
- Bradford YM, Van Slyke CE, Howe DG, Fashena D, Frazer K, Martin R, Paddock H, Pich C, Ramachandran S, Ruzicka L, et al. 2023. From multiallele fish to nonstandard environments, how ZFIN assigns phenotypes, human disease models, and gene expression annotations to genes. *Genetics* 224(1):iyad032. doi:10.1093/genetics/iyad032.
- Bult CJ, Sternberg PW. 2023. The alliance of genome resources: transforming comparative genomics. *Mamm Genome.* 34(4):531–544. doi:10.1007/s00335-023-10015-2.
- Bunt SM, Grumblin GB, Field HI, Marygold SJ, Brown NH, Millburn GH. 2012. FlyBase Consortium. Directly e-mailing authors of newly published papers encourages community curation. *Database (Oxford)* 2012:bas024. doi:10.1093/database/bas024.
- Carotenuto R, Pallotta MM, Tussellino M, Fogliano C. 2023. *Xenopus laevis* (Daudin, 1802) as a model organism for bioscience: a historic review and perspective. *Biology (Basel)* 12(6):890. doi:10.3390/biology12060890.
- Cohen S, Sternberg P. 2019. Genome editing of *Caenorhabditis briggsae* using CRISPR/Cas9 co-conversion marker *dpy-10*. *MicroPubl Biol.* 2019:000171. doi:10.17912/micropub.biology.000171.
- Cohen SM, Wrobel CJJ, Prakash SJ, Schroeder FC, Sternberg PW. 2022. Formation and function of dauer ascarosides in the nematodes *Caenorhabditis briggsae* and *Caenorhabditis elegans*. *G3 (Bethesda)* 12(3):jkac014. doi:10.1093/g3journal/jkac014.
- Cosentino S, Iwasaki W. 2019. SonicParanoid: fast, accurate and easy orthology inference. *Bioinformatics* 35(1):149–151. doi:10.1093/bioinformatics/bty631.
- Davis P, Zarowiecki M, Arnaboldi V, Becerra A, Cain S, Chan J, Chen WJ, Cho J, da Veiga Beltrame E, Diamantakis S, et al. 2022. WormBase in 2022-data, processes, and tools for analyzing *Caenorhabditis elegans*. *Genetics* 220(4):iyac003. doi:10.1093/genetics/iyac003.
- 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.
- Emms DM, Kelly S. 2019. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biol.* 20(1):238. doi:10.1186/s13059-019-1832-y.
- Engel SR, Wong ED, Nash RS, Aleksander S, Douglass AM, Karra E, Miyasato K, Simison SR, Skrzypek M, Weng MS, et al. 2022. New

- data and collaborations at the Saccharomyces Genome Database: updated reference genome, alleles, and the Alliance of Genome Resources. *Genetics* 220(4):iyab224. doi:[10.1093/genetics/iyab224](https://doi.org/10.1093/genetics/iyab224).
- Fang R, Schindelman G, Van Auken K, Fernandes J, Chen W, Wang X, Davis P, Tuli MA, Marygold SJ, Millburn G, et al. 2012. Automatic categorization of diverse experimental information in the bio-science literature. *BMC Bioinformatics* 13(1):16. doi:[10.1186/1471-2105-13-16](https://doi.org/10.1186/1471-2105-13-16).
- Fisher M, James-Zorn C, Ponferrada V, Bell AJ, Sundararaj N, Segerdell E, Chaturvedi P, Bayyari N, Chu S, Pells T, et al. 2023. Xenbase: key features and resources of the *Xenopus* model organism knowledgebase. *Genetics* 224(1):iyad018. doi:[10.1093/genetics/iyad018](https://doi.org/10.1093/genetics/iyad018).
- FlyBase Consortium. 1999. The FlyBase database of the *Drosophila* Genome Projects and community literature. *Nucleic Acids Res.* 27(1):85–88. doi:[10.1093/nar/27.1.85](https://doi.org/10.1093/nar/27.1.85).
- Fuentes D, Molina M, Chorostecki U, Capella-Gutiérrez S, Marcet-Houben M, Gabaldón T. 2022. PhylomeDB V5: an expanding repository for genome-wide catalogues of annotated gene phylogenies. *Nucleic Acids Res.* 50(D1):D1062–D1068. doi:[10.1093/nar/gkab966](https://doi.org/10.1093/nar/gkab966).
- Gene Ontology Consortium. 2023. The Gene Ontology knowledgebase in 2023. *Genetics* 224(1):iyad031. doi:[10.1093/genetics/iyad031](https://doi.org/10.1093/genetics/iyad031).
- Gramates LS, Agapite J, Attrill H, Calvi BR, Crosby MA, Dos Santos G, Goodman JL, Goutte-Gattat D, Jenkins VK, Kaufman T, et al. 2022. FlyBase: a guided tour of highlighted features. *Genetics* 220(4):iyac035. doi:[10.1093/genetics/iyac035](https://doi.org/10.1093/genetics/iyac035).
- Howe DG, Blake JA, Bradford YM, Bult CJ, Calvi BR, Engel SR, Kadin JA, Kaufman TC, Kishore R, Laulederkind SJF, et al. 2018. Model organism data evolving in support of translational medicine. *Lab Anim (NY)*. 47(10):277–289. doi:[10.1038/s41684-018-0150-4](https://doi.org/10.1038/s41684-018-0150-4).
- Hu Y, Comjean A, Rodiger J, Liu Y, Gao Y, Chung V, Zirin J, Perrimon N, Mohr SE. 2021. FlyRNAi.org—the database of the *Drosophila* RNAi screening center and transgenic RNAi project: 2021 update. *Nucleic Acids Res.* 49(D1):D908–D915. doi:[10.1093/nar/gkaa936](https://doi.org/10.1093/nar/gkaa936).
- Hu Y, Flockhart I, Vinayagam A, Bergwitz C, Berger B, Perrimon N, Mohr SE. 2011. An integrative approach to ortholog prediction for disease-focused and other functional studies. *BMC Bioinformatics* 12(1):357. doi:[10.1186/1471-2105-12-357](https://doi.org/10.1186/1471-2105-12-357).
- Inoue T, Aillon M, Poon S, Kim HK, Thomas JH, Sternberg PW. 2007. Genetic analysis of dauer formation in *Caenorhabditis briggsae*. *Genetics* 177(2):809–818. doi:[10.1534/genetics.107.078857](https://doi.org/10.1534/genetics.107.078857).
- Ivanova M, Moss EG. 2023. Orthologs of the *C. elegans* heterochronic genes have divergent functions in *C. briggsae*. *Genetics* 225(4):iyad177. doi:[10.1093/genetics/iyad177](https://doi.org/10.1093/genetics/iyad177).
- Jhaveri N, van den Berg W, Hwang BJ, Muller HM, Sternberg PW, Gupta BP. 2022. Genome annotation of *Caenorhabditis briggsae* by TEC-RED identifies new exons, paralogs, and conserved and novel operons. *G3 (Bethesda)* 12(7):jkac101. doi:[10.1093/g3journal/jkac101](https://doi.org/10.1093/g3journal/jkac101).
- Jiang X, Li P, Kadin J, Blake JA, Ringwald M, Shatkay H. 2020. Integrating image caption information into biomedical document classification in support of biocuration. *Database (Oxford)* 2020:baaa024. doi:[10.1093/database/baaa024](https://doi.org/10.1093/database/baaa024).
- 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](https://doi.org/10.1093/database/baaa037).
- Kostiuk V, Khokha MK. 2021. *Xenopus* as a platform for discovery of genes relevant to human disease. *Curr Top Dev Biol.* 145:277–312. doi:[10.1016/bs.ctdb.2021.03.005](https://doi.org/10.1016/bs.ctdb.2021.03.005).
- 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](https://doi.org/10.1093/nar/gkaa1026).
- Liu W, Laulederkind SJ, Hayman GT, Wang SJ, Nigam R, Smith JR, De Pons J, Dwinell MR, Shimoyama M. 2015. OntoMate: a text-mining tool aiding curation at the Rat Genome Database. *Database (Oxford)* 2015:bau129. doi:[10.1093/database/bau129](https://doi.org/10.1093/database/bau129).
- Milacic M, Rothfels K, Mathews L, Wright A, Jassal B, Shamovsky V, Trinh Q, Gillespie M, Sevilla C, Tiwari K, et al. 2024. The reactome pathway knowledgebase 2024. *Nucleic Acids Res.* 52(D1):D672–D678. doi:[10.1093/nar/gkad1025](https://doi.org/10.1093/nar/gkad1025).
- Mitros T, Lyons JB, Session AM, Jenkins J, Shu S, Kwon T, Lane M, Ng C, Grammer TC, Khokha MK, et al. 2019. A chromosome-scale genome assembly and dense genetic map for *Xenopus tropicalis*. *Dev Biol.* 452(1):8–20. doi:[10.1016/j.ydbio.2019.03.015](https://doi.org/10.1016/j.ydbio.2019.03.015).
- Moya ND, Stevens L, Miller IR, Sokol CE, Galindo JL, Bardas AD, Koh ESH, Rozenich J, Yeo C, Xu M, et al. 2023. Novel and improved *Caenorhabditis briggsae* gene models generated by community curation. *BMC Genomics* 24(1):486. doi:[10.1186/s12864-023-09582-0](https://doi.org/10.1186/s12864-023-09582-0).
- Müller HM, Kenny EE, Sternberg PW. 2004. Textpresso: an ontology-based information retrieval and extraction system for biological literature. *PLoS Biol.* 2(11):e309. doi:[10.1371/journal.pbio.0020309](https://doi.org/10.1371/journal.pbio.0020309).
- Müller HM, Van Auken KM, Li Y, Sternberg PW. 2018. Textpresso Central: a customizable platform for searching, text mining, viewing, and curating biomedical literature. *BMC Bioinformatics* 19(1):94. doi:[10.1186/s12859-018-2103-8](https://doi.org/10.1186/s12859-018-2103-8).
- Nevers Y, Jones TEM, Jyothi D, Yates B, Ferret M, Portell-Silva L, Codo L, Cosentino S, Marcet-Houben M, Vlasova A, et al. 2022. The Quest for Orthologs orthology benchmark service in 2022. *Nucleic Acids Res.* 50(W1):W623–W632. doi:[10.1093/nar/gkac330](https://doi.org/10.1093/nar/gkac330).
- Nevers Y, Kress A, Defosset A, Ripp R, Linard B, Thompson JD, Poch O, Lecompte O. 2019. OrthoInspector 3.0: open portal for comparative genomics. *Nucleic Acids Res.* 47(D1):D411–D418. doi:[10.1093/nar/gky1068](https://doi.org/10.1093/nar/gky1068).
- Oliver SG, Lock A, Harris MA, Nurse P, Wood V. 2016. Model organism databases: essential resources that need the support of both funders and users. *BMC Biol.* 14(1):49. doi:[10.1186/s12915-016-0276-z](https://doi.org/10.1186/s12915-016-0276-z).
- Persson E, Sonnhammer ELL. 2022. InParanoid-DIAMOND: faster orthology analysis with the InParanoid algorithm. *Bioinformatics* 38(10):2918–2919. doi:[10.1093/bioinformatics/btac194](https://doi.org/10.1093/bioinformatics/btac194).
- Priyam A, Woodcroft BJ, Rai V, Moghul I, Munagala A, Ter F, Chowdhary H, Pieniak I, Maynard LJ, Gibbins MA, et al. 2019. SequenceServer: a modern graphical user interface for custom BLAST databases. *Mol Biol Evol.* 36(12):2922–2924. doi:[10.1093/molbev/msz185](https://doi.org/10.1093/molbev/msz185).
- Ringwald M, Richardson JE, Baldarelli RM, Blake JA, Kadin JA, Smith C, Bult CJ. 2022. Mouse Genome Informatics (MGI): latest news from MGD and GXD. *Mamm Genome.* 33(1):4–18. doi:[10.1007/s00335-021-09921-0](https://doi.org/10.1007/s00335-021-09921-0).
- Sargent L, Liu Y, Leung W, Mortimer NT, Lopatto D, Goecks J, Elgin SCR. 2020. G-OnRamp: generating genome browsers to facilitate undergraduate-driven collaborative genome annotation. *PLoS Comput Biol.* 16(6):e1007863. doi:[10.1371/journal.pcbi.1007863](https://doi.org/10.1371/journal.pcbi.1007863).
- Session AM, Uno Y, Kwon T, Chapman JA, Toyoda A, Takahashi S, Fukui A, Hikosaka A, Suzuki A, Kondo M, et al. 2016. Genome evolution in the allotetraploid frog *Xenopus laevis*. *Nature* 538(7625):336–343. doi:[10.1038/nature19840](https://doi.org/10.1038/nature19840).
- Sharanya D, Thillainathan B, Marri S, Bojanala N, Taylor J, Flibotte S, Moerman DG, Waterston RH, Gupta BP. 2012. Genetic control of vulval development in *Caenorhabditis briggsae*. *G3 (Bethesda)* 2(12):1625–1641. doi:[10.1534/g3.112.004598](https://doi.org/10.1534/g3.112.004598).

- Smith RN, Aleksic J, Butano D, Carr A, Contrino S, Hu F, Lyne M, Lyne R, Kalderimis A, Rutherford K, *et al.* 2012. InterMine: a flexible data warehouse system for the integration and analysis of heterogeneous biological data. *Bioinformatics* 28(23):3163–3165. doi:10.1093/bioinformatics/bts577.
- Sternberg PW, Van Auken K, Wang Q, Wright A, Yook K, Zarowiecki M, Arnaboldi V, Becerra A, Brown S, Cain S, *et al.* 2024. WormBase 2024: status and transitioning to Alliance infrastructure. *Genetics*. iyae050. doi:10.1093/genetics/iyae050.
- Thomas PD, Ebert D, Muruganujan A, Mushayahama T, Albou LP, Mi H. 2022. PANTHER: making genome-scale phylogenetics accessible to all. *Protein Sci.* 31(1):8–22. doi:10.1002/pro.4218.
- Thomas PD, Hill DP, Mi H, Osumi-Sutherland D, Van Auken K, Carbon S, Balhoff JP, Albou LP, Good B, Gaudet P, *et al.* 2019. Gene Ontology Causal Activity Modeling (GO-CAM) moves beyond GO annotations to structured descriptions of biological functions and systems. *Nat Genet.* 51(10):1429–1433. doi:10.1038/s41588-019-0500-1.
- UniProt Consortium. 2023. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Res.* 51(D1):D523–D531. doi:10.1093/nar/gkac1052.
- Van Auken K, Fey P, Berardini TZ, Dodson R, Cooper L, Li D, Chan J, Li Y, Basu S, Muller HM, *et al.* 2012. WormBase Consortium. Text mining in the biocuration workflow: applications for literature curation at WormBase, dictyBase and TAIR. *Database (Oxford)* 2012:bas040. doi:10.1093/database/bas040.
- Vedi M, Smith JR, Thomas Hayman G, Tutaj M, Brodie KC, De Pons JL, Demos WM, Gibson AC, Kaldunski ML, Lamers L, *et al.* 2023. 2022 updates to the rat genome database: a findable, accessible, interoperable, and reusable (FAIR) resource. *Genetics* 224(1):iyad042. doi:10.1093/genetics/iyad042.
- Wood V, Sternberg PW, Lipshitz HD. 2022. Making biological knowledge useful for humans and machines. *Genetics* 220(4):iyac001. doi:10.1093/genetics/iyac001.

Editor: V. Wood

The Alliance of Genome Resources Consortium (alphabetical)

Suzanne A. Aleksander

Department of Genetics, Stanford University, Stanford, CA 94305

Anna V. Anagnostopoulos

The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME 04609, USA

Giulia Antonazzo

Department of Physiology, Development and Neuroscience, University of Cambridge, Downing Street, Cambridge CB2 3DY, UK

Valerio Arnaboldi

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

Helen Attrill

Department of Physiology, Development and Neuroscience, University of Cambridge, Downing Street, Cambridge CB2 3DY, UK

Andrés Becerra

European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK

Susan M. Bello

The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME 04609, USA

Olin Blodgett

The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME 04609, USA

Yvonne M. Bradford

Institute of Neuroscience, University of Oregon, Eugene, OR 97403

Carol J. Bult

The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME 04609, USA

Scott Cain

Informatics and Bio-computing Platform, Ontario Institute for Cancer Research, Toronto, ON M5G0A3, Canada

Brian R. Calvi

Department of Biology, Indiana University, Bloomington, IN 47408, USA

Seth Carbon

Environmental Genomics and Systems Biology, Lawrence Berkeley National Laboratory, Berkeley, CA

Juancarlos Chan

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

Wen J. Chen

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

J. Michael Cherry

Department of Genetics, Stanford University, Stanford, CA 94305

Jaehyoung Cho

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

Madeline A. Crosby

The Biological Laboratories, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, USA

Jeffrey L. De Pons

Medical College of Wisconsin—Rat Genome Database, Departments of Physiology and Biomedical Engineering, Medical College of Wisconsin, Milwaukee, WI 53226, USA

Peter D'Eustachio

NYU Grossman School of Medicine, New York, NY 10016

Stavros Diamantakis

European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK

Mary E. Dolan

The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME 04609, USA

Gilberto dos Santos

The Biological Laboratories, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, USA

Sarah Dyer

European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK

Dustin Ebert

Department of Population and Public Health Sciences, University of Southern California, Los Angeles, CA 90033, USA

Stacia R. Engel

Department of Genetics, Stanford University, Stanford, CA 94305

David Fashena

Institute of Neuroscience, University of Oregon, Eugene, OR 97403

Malcolm Fisher

Division of Developmental Biology, Cincinnati Children's Hospital Medical Center, 3333 Burnet Ave, Cincinnati, OH 45229, USA

Saoirse Foley

Department of Biological Sciences, Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh, PA 15203

Adam C. Gibson

Medical College of Wisconsin—Rat Genome Database, Departments of Physiology and Biomedical Engineering, Medical College of Wisconsin, Milwaukee, WI 53226, USA

Varun R. Gollapally

Medical College of Wisconsin—Rat Genome Database, Departments of Physiology and Biomedical Engineering, Medical College of Wisconsin, Milwaukee, WI 53226, USA

L. Sian Gramates

The Biological Laboratories, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, USA

Christian A. Grove

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

Paul Hale

The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME 04609, USA

Todd Harris

Informatics and Bio-computing Platform, Ontario Institute for Cancer Research, Toronto, ON M5G0A3, Canada

G. Thomas Hayman

Medical College of Wisconsin—Rat Genome Database, Departments of Physiology and Biomedical Engineering, Medical College of Wisconsin, Milwaukee, WI 53226, USA

Yanhui Hu

Department of Genetics, Howard Hughes Medical Institute, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, MA 02115, USA

Christina James-Zorn

Division of Developmental Biology, Cincinnati Children's Hospital Medical Center, 3333 Burnet Ave, Cincinnati, OH 45229, USA

Kamran Karimi

Department of Biological Sciences, University of Calgary, 507 Campus Dr NW, Calgary, AB T2N 4V8, Canada

Kalpna Karra

Department of Genetics, Stanford University, Stanford, CA 94305

Ranjana Kishore

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

Anne E. Kwitek

Medical College of Wisconsin—Rat Genome Database, Departments of Physiology and Biomedical Engineering, Medical College of Wisconsin, Milwaukee, WI 53226, USA

Stanley J.F. Laulederkind

Medical College of Wisconsin—Rat Genome Database, Departments of Physiology and Biomedical Engineering, Medical College of Wisconsin, Milwaukee, WI 53226, USA

Raymond Lee

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

Ian Longden

The Biological Laboratories, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, USA

Manuel Luyypaert

European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK

Nicholas Markarian

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

Steven J. Marygold

Department of Physiology, Development and Neuroscience, University of Cambridge, Downing Street, Cambridge CB2 3DY, UK

Beverley Matthews

The Biological Laboratories, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, USA

Monica S. McAndrews

The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME 04609, USA

Gillian Millburn

Department of Physiology, Development and Neuroscience, University of Cambridge, Downing Street, Cambridge CB2 3DY, UK

Stuart Miyasato

Department of Genetics, Stanford University, Stanford, CA 94305

Howie Motenko

The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME 04609, USA

Sierra Moxon

Environmental Genomics and Systems Biology, Lawrence Berkeley National Laboratory, Berkeley, CA

Hans-Michael Muller

Division of Biology and Biological Engineering 140-18, California Institute of Technology, Pasadena, CA 91125, USA

Christopher J. Mungall

Environmental Genomics and Systems Biology, Lawrence Berkeley National Laboratory, Berkeley, CA

Anushya Muruganujan

Department of Population and Public Health Sciences, University of Southern California, Los Angeles, CA 90033, USA

Tremayne Mushayahama

Department of Population and Public Health Sciences, University of Southern California, Los Angeles, CA 90033, USA

Robert S. Nash

Department of Genetics, Stanford University, Stanford, CA 94305

Paulo Nuin

Informatics and Bio-computing Platform, Ontario Institute for Cancer Research, Toronto, ON M5G0A3, Canada

Holly Paddock

Institute of Neuroscience, University of Oregon, Eugene, OR 97403

Troy Pells

Department of Biological Sciences, University of Calgary, 507 Campus Dr NW, Calgary, AB T2N 4V8, Canada

Norbert Perrimon

Department of Genetics, Howard Hughes Medical Institute, Harvard Medical School, 77 Avenue Louis Pasteur, Boston, MA 02115, USA

Christian Pich

Institute of Neuroscience, University of Oregon, Eugene, OR 97403

Mark Quinton-Tulloch

European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK

Daniela Raciti

Division of Biology and Biological Engineering 140-18,
California Institute of Technology, Pasadena, CA 91125, USA

Sridhar Ramachandran

Institute of Neuroscience, University of Oregon, Eugene, OR
97403

Joel E. Richardson

Institute of Neuroscience, University of Oregon, Eugene, OR
97403

Susan Russo Gelbart

The Biological Laboratories, Harvard University, 16 Divinity
Avenue, Cambridge, MA 02138, USA

Leyla Ruzicka

Institute of Neuroscience, University of Oregon, Eugene, OR
97403

Gary Schindelman

Division of Biology and Biological Engineering 140-18,
California Institute of Technology, Pasadena, CA 91125, USA

David R. Shaw

The Jackson Laboratory for Mammalian Genomics, Bar Harbor,
ME 04609, USA

Gavin Sherlock

Department of Genetics, Stanford University, Stanford, CA
94305

Ajay Shrivatsav

Department of Genetics, Stanford University, Stanford, CA
94305

Amy Singer

Institute of Neuroscience, University of Oregon, Eugene, OR
97403

Constance M. Smith

The Jackson Laboratory for Mammalian Genomics, Bar Harbor,
ME 04609, USA

Cynthia L. Smith

The Jackson Laboratory for Mammalian Genomics, Bar Harbor,
ME 04609, USA

Jennifer R. Smith

Medical College of Wisconsin—Rat Genome Database,
Departments of Physiology and Biomedical Engineering, Medical
College of Wisconsin, Milwaukee, WI 53226, USA

Lincoln Stein

Informatics and Bio-computing Platform, Ontario Institute for
Cancer Research, Toronto, ON M5G0A3, Canada

Paul W. Sternberg

Division of Biology and Biological Engineering 140-18,
California Institute of Technology, Pasadena, CA 91125, USA
(ORCID ID: 0000-0002-7699-0173)

Christopher J. Tabone

The Biological Laboratories, Harvard University, 16 Divinity
Avenue, Cambridge, MA 02138, USA

Paul D. Thomas

Department of Population and Public Health Sciences,
University of Southern California, Los Angeles, CA 90033, USA

Ketaki Thorat

Medical College of Wisconsin—Rat Genome Database,
Departments of Physiology and Biomedical Engineering, Medical
College of Wisconsin, Milwaukee, WI 53226, USA

Jyothi Thota

Medical College of Wisconsin—Rat Genome Database,
Departments of Physiology and Biomedical Engineering, Medical
College of Wisconsin, Milwaukee, WI 53226, USA

Monika Tomczuk

The Jackson Laboratory for Mammalian Genomics, Bar Harbor,
ME 04609, USA

Vitor Trovisco

Department of Physiology, Development and Neuroscience,
University of Cambridge, Downing Street, Cambridge CB2 3DY,
UK

Marek A. Tutaj

Medical College of Wisconsin—Rat Genome Database,
Departments of Physiology and Biomedical Engineering, Medical
College of Wisconsin, Milwaukee, WI 53226, USA

Jose-Maria Urbano

Department of Physiology, Development and Neuroscience,
University of Cambridge, Downing Street, Cambridge CB2 3DY, UK

Kimberly Van Auken

Division of Biology and Biological Engineering 140-18,
California Institute of Technology, Pasadena, CA 91125, USA

Ceri E. Van Slyke

Institute of Neuroscience, University of Oregon, Eugene, OR
97403

Peter D. Vize

Department of Biological Sciences, University of Calgary, 507
Campus Dr NW, Calgary, AB T2N 4V8, Canada

Qinghua Wang

Division of Biology and Biological Engineering 140-18,
California Institute of Technology, Pasadena, CA 91125, USA

Shuai Weng

Department of Genetics, Stanford University, Stanford, CA
94305

Monte Westerfield

Institute of Neuroscience, University of Oregon, Eugene, OR
97403

Laurens G. Wilming

The Jackson Laboratory for Mammalian Genomics, Bar Harbor,
ME 04609, USA

Edith D. Wong

Department of Genetics, Stanford University, Stanford, CA
94305

Adam Wright

Informatics and Bio-computing Platform, Ontario Institute for
Cancer Research, Toronto, ON M5G0A3, Canada

Karen Yook

Division of Biology and Biological Engineering 140-18,
California Institute of Technology, Pasadena, CA 91125, USA

Pinglei Zhou

The Biological Laboratories, Harvard University, 16 Divinity
Avenue, Cambridge, MA 02138, USA

Aaron Zorn

Division of Developmental Biology, Cincinnati Children's
Hospital Medical Center, 3333 Burnet Ave, Cincinnati, OH 45229,
USA

Mark Zytkovicz

The Biological Laboratories, Harvard University, 16 Divinity
Avenue, Cambridge, MA 02138, USA