

1 Updates to the Alliance of Genome Resources Central Infrastructure

2 Alliance of Genome Resources Consortium

3 Revised for Genetics 2024

4 The Alliance of Genome Resources Consortium (alphabetical)

5 Suzanne A. Aleksander⁹, Anna V. Anagnostopoulos⁵, Giulia Antonazzo¹⁰, Valerio Arnaboldi¹,
 6 Helen Attrill¹⁰, Andrés Becerra², Susan M. Bello⁵, Olin Blodgett⁵, Yvonne M. Bradford¹¹, Carol J.
 7 Bult⁵, Scott Cain⁸, Brian R. Calvi⁴, Seth Carbon⁶, Juancarlos Chan¹, Wen J. Chen¹, J. Michael
 8 Cherry⁹, Jaehyoung Cho¹, Madeline A. Crosby³, Jeffrey L. De Pons⁷, Peter D'Eustachio¹⁵,
 9 Stavros Diamantakis², Mary E. Dolan⁵, Gilberto dos Santos³, Sarah Dyer², Dustin Ebert¹², Stacia
 10 R. Engel⁹, David Fashena¹¹, Malcolm Fisher¹⁶, Saoirse Foley¹³, Adam C. Gibson⁷, Varun R.
 11 Gollapally⁷, L. Sian Gramates³, Christian A. Grove¹, Paul Hale⁵, Todd Harris⁸, G. Thomas
 12 Hayman⁷, Yanhui Hu¹⁴, Christina James-Zorn¹⁶, Kamran Karimi¹⁷, Kalpana Karra⁹, Ranjana
 13 Kishore¹, Anne E. Kwitek⁷, Stanley J. F. Laulederkind⁷, Raymond Lee¹, Ian Longden³, Manuel
 14 Luypaert², Nicholas Markarian¹, Steven J. Marygold¹⁰, Beverley Matthews³, Monica S.
 15 McAndrews⁵, Gillian Millburn¹⁰, Stuart Miyasato⁹, Howie Motenko⁵, Sierra Moxon⁶, Hans-
 16 Michael Muller¹, Christopher J. Mungall⁶, Anushya Muruganujan¹², Tremayne Mushayahama¹²,
 17 Robert S. Nash⁹, Paulo Nuin⁸, Holly Paddock¹¹, Troy Pells¹⁷, Norbert Perrimon¹⁴, Christian
 18 Pich¹¹, Mark Quinton-Tulloch², Daniela Raciti¹, Sridhar Ramachandran¹¹, Joel E. Richardson¹¹,
 19 Susan Russo Gelbart³, Leyla Ruzicka¹¹, Gary Schindelman¹, David R. Shaw⁵, Gavin Sherlock⁹,
 20 Ajay Shrivatsav⁹, Amy Singer¹¹, Constance M. Smith⁵, Cynthia L. Smith⁵, Jennifer R. Smith⁷,
 21 Lincoln Stein⁸, Paul W. Sternberg¹, Christopher J. Tabone³, Paul D. Thomas¹², Ketaki Thorat⁷,
 22 Jyothi Thota⁷, Monika Tomczuk⁵, Vitor Trovisco¹⁰, Marek A. Tutaj⁷, Jose-Maria Urbano¹⁰,
 23 Kimberly Van Auken¹, Ceri E. Van Slyke¹¹, Peter D. Vize¹⁷, Qinghua Wang¹, Shuai Weng⁹,
 24 Monte Westerfield¹¹, Laurens G. Wilming⁵, Edith D. Wong⁹, Adam Wright⁸, Karen Yook¹, Pinglei
 25 Zhou³, Aaron Zorn¹⁶, Mark Zytkevich³

26
 27
 28
 29
 30 ¹Division of Biology and Biological Engineering 140-18, California Institute of Technology,
 31 Pasadena, CA 91125, USA.

32 ²European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust
 33 Genome Campus, Hinxton, Cambridge CB10 1SD, UK.

34 ³The Biological Laboratories, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138,
 35 USA.

36 ⁴Department of Biology, Indiana University, Bloomington, IN 47408, USA.

37 ⁵The Jackson Laboratory for Mammalian Genomics, Bar Harbor, ME, 04609, USA.

38 ⁶Environmental Genomics and Systems Biology, Lawrence Berkeley National Laboratory,
 39 Berkeley, CA

© The Author(s) 2024. Published by Oxford University Press on behalf of The Genetics Society of America.
 This is an Open Access article distributed under the terms of the Creative Commons Attribution License
 (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and
 reproduction in any medium, provided the original work is properly cited.

1 ⁷Medical College of Wisconsin - Rat Genome Database, Departments of Physiology and
2 Biomedical Engineering, Medical College of Wisconsin, Milwaukee, WI 53226, USA.

3 ⁸Informatics and Bio-computing Platform, Ontario Institute for Cancer Research, Toronto, ON
4 M5G0A3, Canada.

5 ⁹Department of Genetics, Stanford University, Stanford, CA 94305

6 ¹⁰Department of Physiology, Development and Neuroscience, University of Cambridge,
7 Downing Street, Cambridge CB2 3DY, UK.

8 ¹¹Institute of Neuroscience, University of Oregon, Eugene, OR 97403

9 ¹²Department of Population and Public Health Sciences, University of Southern California, Los
10 Angeles, CA 90033, USA

11 ¹³Department of Biological Sciences, Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh,
12 PA 15203

13 ¹⁴Department of Genetics, Howard Hughes Medical Institute, Harvard Medical School, 77
14 Avenue Louis Pasteur, Boston, MA 02115, USA

15 ¹⁵NYU Grossman School of Medicine, New York NY 10016

16 ¹⁶ Division of Developmental Biology, Cincinnati Children's Hospital Medical Center, 3333
17 Burnet Ave, Cincinnati, OH 45229, USA

18 ¹⁷ Department of Biological Sciences, University of Calgary, 507 Campus Dr NW, Calgary, AB
19 T2N 4V8, Canada

20
21

22 **correspondence: Paul W. Sternberg pws@caltech.edu**

23

24 **Running Title: Alliance of Genome Resources**

25

26 Keywords: database; knowledgebase; software; text mining; data integration; *Drosophila*; yeast;
27 *C. elegans*; zebrafish; mouse

28

29 **Abstract**

30 The Alliance of Genome Resources (Alliance) is an extensible coalition of knowledgebases
31 focused on the genetics and genomics of intensively-studied model organisms. The Alliance is
32 organized as individual knowledge centers with strong connections to their research
33 communities and a centralized software infrastructure, discussed here. Model organisms
34 currently represented in the Alliance are budding yeast, *C. elegans*, *Drosophila*, zebrafish, frog,
35 laboratory mouse, laboratory rat, and the Gene Ontology Consortium. The project is in a rapid
36 development phase to harmonize knowledge, store it, analyze it, and present it to the
37 community through a web portal, direct downloads, and Application Programming Interfaces

1 (APIs). Here we focus on developments over the last two years. Specifically, we added and
2 enhanced tools for browsing the genome (JBrowse), downloading sequences, mining complex
3 data (AllianceMine), visualizing pathways, full-text searching of the literature (Textpresso), and
4 sequence similarity searching (SequenceServer). We enhanced existing interactive data tables
5 and added an interactive table of paralogs to complement our representation of orthology. To
6 support individual model organism communities, we implemented species-specific “landing
7 pages” and will add disease-specific portals soon; in addition, we support a common community
8 forum implemented in Discourse software. We describe our progress towards a central
9 persistent database to support curation, the data modeling that underpins harmonization, and
10 progress towards a state-of-the art literature curation system with integrated Artificial
11 Intelligence and Machine Learning (AI/ML).
12

13 **Introduction**

14 As has been discussed at length elsewhere (e.g., Oliver et al. 2016; Wood et al. 2022), model
15 organism knowledgebases (a.k.a. model organism databases; MODs) provide daily utility to
16 researchers for the design and interpretation of experiments, to computational biologists for
17 curated datasets, and to genomic researchers for annotated genomes. Some of the major uses
18 of the MODs have been one-stop shopping for all information about a particular gene or
19 obtaining cleansed datasets with standard metadata for computational analyses.
20

21 The Alliance of Genome Resources (referred to herein as the Alliance) is a consortium of MODs
22 and the Gene Ontology Consortium (GOC). The mission of the Alliance is to support
23 comparative genomics to investigate the genetic and genomic basis of human biology, health,
24 and disease. To promote sustainability of the core community data resources that make up the
25 Alliance, we implemented an extensible “knowledge commons” platform for comparative
26 genomics built with modular, re-usable infrastructure components that can support informatics
27 resource needs across a wide range of species (Alliance of Genome Resources, 2022; Howe et
28 al. 2018; Bult and Sternberg, 2023). In 2022, the Alliance was recognized as a Core Global
29 Biodata Resource by the Global Biodata Coalition (Anderson et al 2017).
30

31 Specifically, the Alliance of Genome Resources is organized as two interdependent units:
32 Alliance Central and the Alliance Knowledge Centers. **Alliance Central** is responsible for
33 developing and maintaining the software for data access and for the coordination of data
34 harmonization and data modeling activities across our members. A primary goal of Alliance
35 Central is to reduce redundancy in systems administration and software development for model
36 organism knowledgebases and to deploy a unified ‘look and feel’ for access to, and display of,
37 common data types and annotations across diverse model organisms and human, following
38 Findability, Accessibility, Interoperability, and Reuse (FAIR) guiding principles. Model organism-
39 specific knowledgebases serve as **Alliance Knowledge Centers**. Knowledge Centers are
40 responsible for expert curation and submission of data to Alliance Central using Alliance Central
41 infrastructure. Knowledge Centers also are responsible for organism-specific user support
42 activities and for providing access to data types not yet supported by Alliance Central. The
43 founding Alliance Knowledge Centers are *Saccharomyces* Genome Database (Engel et al.
44 2022), WormBase (Davis et al. 2022; Sternberg et al., 2024), FlyBase (Gramates et al. 2022),

1 Mouse Genome Database (Ringwald et al. 2022), the Zebrafish Information Network (Bradford
2 et al. 2023), Rat Genome Database (Vedi et al. 2023), and the Gene Ontology Consortium
3 (Gene Ontology Consortium 2023). The newest member, Xenbase (Fisher et al, 2023), joined
4 the Alliance consortium in 2022.

5
6 Here we describe our progress toward harmonizing information provided by our member
7 resources, our development of a software infrastructure for ingest, curation, storage, analysis,
8 and output of such information, and development of an efficient literature curation system. We
9 start by describing new features in our web portal at AllianceGenome.org.

10 11 **The Web Portal**

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

19
20 **Xenopus in the Alliance.** Xenbase, the *Xenopus* knowledgebase (Fisher et al 2023), is the first
21 knowledgebase to join the Alliance since the founding members initiated the consortium.
22 *Xenopus* is an amphibian frog species used extensively in biomedical research, and in particular
23 for experimental embryology, cell biology, and disease modeling with genome editing
24 (Carotenuto et al. 2023; Kostiuik and Khokha 2021). As a non-mammalian air-breathing
25 tetrapod, *Xenopus* represents a valuable evolutionary transition between rodents and zebrafish
26 for comparative genomic studies. Xenbase is built on the same underlying data schema
27 (structure) as FlyBase (Chado). Two different *Xenopus* species are used interchangeably as a
28 model system: *X. tropicalis* is a diploid that is the preferred system for genome editing and
29 genetics, whereas *X. laevis* is an allotetraploid preferred for use in cell biology studies,
30 microinjection, and microsurgery-style experimentation. *X. tropicalis* has 1:1 relationships
31 between most genes and human orthologs (excluding paralogs) (Mitros et al. 2019), whereas *X.*
32 *laevis* has two copies of most human orthologs. The allotetraploid formed via hybridization of
33 two different frog species (Session et al 2016), and the complexities of genome evolution that
34 subsequently occurred increase the difficulty of identifying orthology of the two *X. laevis* genes
35 to their diploid relatives, including humans. Mapping of the diploid *X. tropicalis* genes to their
36 human orthologs was performed as with the other organisms in the Alliance (see below).
37 Because this method does not yet work in the context of an allotetraploid, the Alliance imports
38 the *X. tropicalis* to *X. laevis* paralogy mappings from Xenbase, where they have been
39 established through a combination of synteny analysis and manual curation; this was one major
40 challenges in adding *Xenopus* to the Alliance.

41
42 Xenbase created software to upload content on a regular schedule formatted for the current
43 Alliance data ingest schema. Currently these data include orthology, the *Xenopus* anatomical
44 ontology, standard gene information, gene expression data, publications, GO term associations,

1 disease associations, anatomical phenotypes, and genome details. *Xenopus* genes can be
2 found using the Alliance landing page search tool with *Xenopus* genes flagged by *Xtr* and *Xla*
3 notations. The two copies of the genes in *X. laevis*, the allotetraploid, are further tagged as
4 '(symbol).L' and '(symbol).S' to denote the genes on the long (L) and short (S) chromosome
5 pairs of this species (e.g., *pax6.L* and *pax6.S*). Alliance release 6.0.0 has Xenbase data for
6 54,000 genes, 19,000 disease associations, over 45,000 gene expression records and more
7 than 7,000 anatomical phenotypes. Expression and phenotype data will be available in about a
8 year.

9
10 In addition to the rich data made available to the Alliance from *Xenopus* research, this effort also
11 served as a valuable test case for understanding the level of effort and complexities engendered
12 in the addition of new knowledgebases to the Alliance, and the functionality and adaptability of
13 ingest system components.

14
15 **New gene page section: paralogy.** Gene pages now include a Paralogy section populated
16 with data from the Drosophila Research & Screening Center (DRSC) Integrative Ortholog
17 Prediction Tool (DIOPT) version 9.1 developed by the DRSC (Hu et al. 2011, 2020). The
18 assembly of protein sets and algorithmic inferences of their orthology from various sources was
19 first centralized by the DRSC and then exported to the Alliance Central. We include the same
20 data sources used for orthology, when these resources also provide paralogy information.
21 Specifically, these resources have performed well on the standardized benchmarking from the
22 Quest for Orthologs (QfO) Consortium (Nevers et al. 2022). Orthologous Matrix (OMA)
23 (Altenhoff et al. 2021) and PANTHER (Thomas et al. 2022) datasets were retrieved through the
24 QfO benchmark portal (<https://orthology.benchmarkservice.org>), and Compara data were
25 acquired directly from the EBI Compara FTP site. In addition, the DRSC conducted local
26 analyses using Inparanoid (Persson and Sonnhammer. 2022), OrthoFinder (Emms and Kelly
27 2019), OrthoInspector (Nevers et al. 2019), and sonicParanoid (Cosentino and Iwasaki 2019)
28 using the UniProt 2020 reference proteome set (UniProt Consortium 2023), the same set used
29 in the downloaded datasets, to ensure consistency. Direct data submissions from PhylomeDB
30 (Fuentes et al. 2022) and the *Saccharomyces* Genome Database (SGD; Engel et al. 2022) were
31 also integrated into the dataset.

32
33 The new paralogy section comprises a table (**Figure 2**), similar to the orthology table, that
34 contains the gene symbol of related paralogs, a calculated rank, alignment length as the
35 number of aligned amino acids, percentage of similarity and identity, as well as a count of the
36 algorithms or methods that call the paralogous match. The ranking score was developed to sort
37 the paralogs by overall similarity, and was reviewed by curators to display optimally an
38 acceptable rank order for well-studied sets of paralogs. The ranking score considers several
39 factors, including alignment length, percent identity, and the number of paralogy methods that
40 identify the paralog. Additional Information for rank determination and alignment length are
41 available to the users via a clickable help icon located next to those column headers.

42
43 The paralog section was released with Alliance version 6.0.0. Forthcoming updates will include
44 the ability to sort and filter the table by column values and the availability of these data via our

1 bulk downloads page. The existing tables on the gene pages for Function, Disease, and
2 Expression all contain checkboxes for "Compare Ortholog Genes" that allow users to search
3 across species for these features. We will add the additional checkbox, "Compare Paralog
4 Genes" to provide similar functionality for paralogous genes in a future Alliance release.

5
6 **JBrowse sequence detail widget.** A recent Alliance 6.0.0 release includes a new "Sequence
7 Detail" section of all gene pages that uses JBrowse and javascript libraries to display an
8 interactive widget that allows users to download DNA and amino acid sequences of genes in
9 several possible configurations: genomic sequence highlighted with UTRs, coding and intronic
10 regions, CDS regions, and translated protein for example (Figure 3). In the next few releases,
11 we will extend the functionality of the widget variant detail pages, where both the wild-type and
12 variant sequences will be provided. When the variant occurs in the context of a protein coding
13 gene, changes to the coding sequence and resulting translated protein will also be displayed
14 and available for download.

15
16 **Model organism BLAST.** For more than two decades, some of the MOD members of the
17 Alliance have hosted their own custom BLAST interfaces (Altschul et al. 1990; e.g., FlyBase
18 Consortium. 1999), that have allowed users to search custom databases related to those model
19 organisms, e.g., subsets of related species or molecular clones and display BLAST hits in
20 Genome Browsers aligned with current gene models. We are now developing an updated and
21 integrated Alliance BLAST, powered by SequenceServer (Priyam et al. 2019), that optimizes
22 sequence analysis across model organisms. We have begun to update BLAST for the individual
23 MODs. The new WormBase BLAST is now available online, and can currently be accessed via
24 the tools menu on wormbase.org. The results are linked to Genome Browsers and Alliance
25 gene pages (**Figure 4**). This tight connection allows users to navigate seamlessly between their
26 BLAST results and the wealth of information available within the Alliance, enhancing the
27 efficiency and depth of genetic research. For example, users can retrieve BLAST results for a
28 sequence of interest and then easily navigate across Genome Browsers for different organisms,
29 with a comparison to different tracks revealing how that sequence aligns with gene models,
30 variants, and experimental tools (**Figure 5**). From a project perspective, developing Alliance
31 BLAST with a common cloud-optimized infrastructure will increase efficiency by reducing the
32 cost of compute overhead and eliminating the need to manage separate MOD systems, which
33 will then allow more focus on developing new functionality to support researchers. Our focus in
34 the upcoming year is directed toward enhancing the user interface, reflecting our commitment to
35 providing an intuitive platform for researchers in model organism genetics. We plan to produce
36 more analysis tools as part of the evolving Alliance portal, thereby broadening the range of
37 resources available for genetic research within the community.

38
39 **AllianceMine.** AllianceMine, a sophisticated, multifaceted search and retrieval tool that utilizes
40 the InterMine software (Smith et al. 2012), offers a unified view of harmonized data, enabling
41 advanced queries across multiple species. For instance, gene lists can be processed as input
42 and simultaneously query different annotations, such as 'Show me genes associated with a
43 (specific disease term)' (**Figure 6**). The results from queries can be combined for further
44 analysis, and saved or downloaded in customizable file formats. Queries themselves can be

1 customized by modifying predefined templates or by creating new templates to access a
2 combination of specific data types. Thus, this powerful tool can be used in multiple ways,
3 namely, for search, discovery, curation, and analysis.

4
5 AllianceMine currently showcases harmonized data encompassing genes, diseases, Gene
6 Ontology (GO), orthology, expression, alleles, variants, and FASTA formatted genome
7 sequences. The tool also offers predefined queries or "templates" for cross-species searching.
8 Continual optimization will ensure timely data synchronization with the main Alliance site, as
9 well as integration of newly harmonized data types. Another aspect of improvement will be the
10 addition of more templates, widgets, and pre-compiled lists, which can serve as logical input for
11 templated queries.

12
13 **SimpleMine.** We designed SimpleMine for biologists to get essential information for a list of
14 genes without any command-line or programming skill, or patience to learn the awesome power
15 of AllianceMine discussed above. Users can submit a list of gene names or IDs to access more
16 than 20 types of essential data with which they are associated. The results are one line per
17 gene with detailed information separated by four types of separators: tab, comma, bar, and
18 semicolon. Users can choose to display the output as HTML or to download a tab-delimited file.
19 Alliance SimpleMine contains ten species curated by the Alliance MODs. It provides easy gene
20 name/ID conversion among MOD ID, public name (the commonly used name for public
21 consumption), NCBI, PANTHER, Ensembl, and UniProtKB. Users can find summarized
22 anatomic and temporal expression patterns, variants, genetic and physical interactions. Other
23 essential gene information includes disease association and orthologs among all ten species.
24 The infrastructure of SimpleMine allows users to perform species-specific searches for lists of
25 genes that take about two seconds to return results, or mixed-species searches that take about
26 10 seconds to complete.

27
28 **Pathway displays with metabolites (GO Causal Activity Models; GO-CAMs).** We
29 implemented a pathway display on Alliance gene pages that presents both GO-CAM (Thomas
30 et al. 2019) and Reactome pathway (Milacic et al. 2024) models. The display queries both the
31 Reactome and GO Application Programming Interfaces (APIs) and shows the number of
32 pathways from each resource that contain the gene of interest. If a gene appears in multiple
33 pathways, users can select which pathway to display. For the GO-CAM models, the viewer has
34 been improved relative to previous releases of the Alliance website (**Figure 7**). First, the layout
35 has been improved to show clearly the overall causal flow through a pathway, from top to
36 bottom and branching as necessary. Second, the viewer displays not only the activities of
37 genes/proteins in a pathway, but also metabolites, which is particularly useful for visualizing
38 metabolic pathways. These metabolites may be either intermediates in a pathway, or regulators
39 of a protein activity. For signaling pathways, we distinguish between direct and indirect
40 regulation, and between positive, negative, or unknown effects.

41 **Harmonized Data Models**

42
43 The transition of data from individual MODs to the Alliance infrastructure requires data
44 harmonization so that existing analogous MOD data classes (types/tables) can be loaded into

1 Alliance databases using a consistent schema and language. The first step is for biocurators
2 from each Alliance knowledge center to agree on which data classes are analogous and can be
3 treated as a single, consolidated data class. The biocurators then align the properties (table
4 columns) of the consolidated data class, including identifiers, types of values, and whether
5 entity-property-value associations/triples require their own respective metadata and/or evidence
6 records. To enable this process the Linked Data Modeling Language (LinkML). We now have a
7 standard workflow and common data modeling patterns that have streamlined the process,
8 which we expect to complete in the next year. The LinkML specifications, authored in human-
9 readable files, are used to programmatically generate JavaScript Object Notation (JSON)
10 schema specifications, which allow Data Quartermasters (DQMs) to move data to the persistent
11 store. These specifications also inform curation software developers how to generate initial
12 backend (Java models and APIs) and front-end infrastructure (curation user interface data
13 tables and detail pages). Once DQMs have submitted their data files for a particular data class,
14 the data are loaded into the persistent store and validated (see persistent store architecture
15 description below) and thus automatically populated into data tables and the curation interface.
16 The data, having been harmonized, ingested, validated, and displayed to curators in the
17 curation software, can now flow through to the public site according to the data pipeline
18 described (see persistent store architecture description below).

19
20 Many Alliance data classes have completely (or nearly completely) harmonized data models in
21 LinkML (see https://github.com/alliance-genome/agr_curation_schema) including: disease
22 annotations, alleles, variants, expression annotations, and references. Although many other
23 data classes have partially harmonized models, ongoing and future harmonization efforts will
24 focus on completing harmonized models for the remaining curated data classes: genes,
25 transcripts, proteins, non-transcribed genome features, affected genomic models (AGMs;
26 strains, genotypes, fish), phenotype annotations, molecular and genetic interactions, gene
27 regulation annotations, high-throughput expression dataset metadata (including for RNA-seq,
28 single-cell RNA-seq, and proteomics datasets), species, reagents such as DNA clones and
29 antibodies, images, persons, laboratories, companies, and various entity set classes like gene
30 sets, which can be used for storing assay results and performing downstream analyses like
31 ontology term enrichment, alignments, and other entity set processing calculations.

32 33 **Persistent Store architecture**

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

44

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

19
20 To enable CRUD operations on persistent store data, curation APIs and a curation user
21 interface accessible with Okta authentication have been implemented (**Figure 9**). Curators can
22 now access data tables for the following data types: genes, alleles, variants, affected genomic
23 models (AGMs; e.g. strains, genotypes), publications (accessed via Alliance Bibliographic
24 Central (ABC) APIs), experimental conditions, constructs, disease annotations, molecules (not
25 already managed by Chemical Entities of Biological Interest (ChEBI)), ontology terms, and
26 controlled vocabularies and their terms. CRUD operations have been fully enabled for disease
27 annotations, experimental conditions, and controlled vocabularies, read-update operations have
28 been enabled for alleles and variants, and read operations are enabled for the remaining data
29 types. Work is underway to fully enable CRUD operations on all remaining data classes and
30 their attributes including new data tables for transcripts, proteins, other (non-gene) genome
31 features, expression annotations, phenotype annotations, molecular interactions, genetic
32 interactions, gene regulation annotations, antibodies, images, and more. In addition to data
33 tables presenting all entries of a particular data class, the curation tool also has individual entity
34 detail pages (for example, see an allele detail page
35 <https://curation.alliancegenome.org/#/allele/MGI:6446761>) for data entry and editing on a
36 dedicated web page for one particular entity. The curation tool also enables user-specific and
37 MOD-specific custom user settings and preferences to provide a user interface most compatible
38 with individual curators' workflows.

39
40 In the next year, the curation tool will include batch creation of data entities (e.g., annotations,
41 reagents), batch editing, data history inspection and auditing, undo and review of latest
42 changes, publication constraints (constrain data view and entry to publication currently being
43 curated), customizations and MOD default settings for new entity creation and detail pages,
44 incorporation of data entity and topic tagging information from the ABC literature store (see

1 below), and incorporation of AI/ML into the curation workflow.

2
3 For releases of persistent store data to the Alliance public website, Postgres database
4 snapshots are taken and sent to a separate Postgres instance that feeds the data via the
5 curation APIs (instantiated as a library) into the public site indexer where various data filtering
6 and transformations occur before making those processed data available to our public website
7 APIs via our Elasticsearch index. The Alliance public website user interface, using existing UI
8 infrastructure, is then modified or created to accommodate the data now flowing from the
9 persistent store database.

10
11 **Security, stability and backups.** All services and data provided by the Alliance to its
12 community are hosted on Amazon web services (AWS). This provides us with industry leading
13 availability of up to 99.99% on services like EC2, which we use to host our virtual servers. We
14 use additional AWS-managed services such as Elastic Beanstalk for application deployment,
15 AWS Relational Database Service for hosting our relational (Postgres) databases, and Amazon
16 OpenSearch Service for hosting our search indexes, which all provide automatic updates and
17 maintenance for increased reliability. All files hosted at the Alliance of Genome Resources are
18 stored in S3 buckets, which ensures industry leading durability and availability. Furthermore, we
19 make daily backups of our relational databases and have processes in place that enable easy
20 restore of those backups in case of failure or data corruption. All Search indexes are derived
21 from the persistent relational database and can be regenerated at any moment when required.

22
23 We make use of separated subnets between public-facing and private systems, and only
24 services requiring public access are given public IP addresses, ensuring that public-facing
25 services such as our curation interface can be accessed by our curators world-wide (through
26 Okta Authentication), although the supporting back-end services such as the supporting
27 databases can be kept private. Access to all services is furthermore restricted to allow access
28 only to the required ports and services through the use of AWS Security Groups to control the
29 allowed network traffic. AWS IAM users, groups, and roles are used to control the allowed AWS
30 operations and access among Alliance developers. In all cases, the principle of least privilege is
31 applied, so that the potential attack surface is reduced to a minimum (for example by not
32 granting blanket AWS admin permissions to developers who do not have an AWS admin
33 function). Access keys to any system can be revoked when misuse of those access keys is
34 detected. We also configured our github repositories to be scanned automatically for accidental
35 secret credential leakages through the use of GitGuardian software.

36 **Literature Acquisition**

37
38 We designed and are implementing a literature system, Alliance Bibliographic Central (ABC),
39 that will support curation, and in the future, end users. The ABC supports the tasks of reference
40 acquisition, triage, and curation workflow management. Specifically, the ABC is an ecosystem of
41 online tools and supporting Alliance databases that manage all references and related metadata
42 that are 'in corpus' for the member MODs.

43
44 Literature acquisition at the Alliance begins with automated, organism-specific PubMed queries

1 to retrieve candidate references for each MOD's corpus. References matching the search
2 criteria are then added to the ABC by assigning an Alliance reference identifier and importing
3 associated bibliographic information to the database. Subsequently, curators manually sort
4 references as either 'in' or 'out of corpus' based on the curation policies of the MOD and
5 eliminate any false positive results from the initial search. While many thousands of papers are
6 published each year, only some have information that is currently curated. For example, in
7 2022, the curatable literature size after triage was: 3181 for ZFIN, 3221 for SGD,; 2130 for
8 FlyBase, 1419 for WormBase, and 437 for Xenbase. Once references are sorted, they enter
9 MOD-specific curation workflows supported by task-specific ABC curator interfaces to, for
10 example, add reference files, manually tag references with specific entities (e.g., genes, alleles,
11 and data types) and topics (e.g., phenotypes, anatomic expression) using the Alliance Tags for
12 Papers (ATP) ontology, and merge duplicate references. In addition to adding reference files
13 manually, the full text of 'in corpus' references included in the PubMed Central (PMC) open
14 access set is also automatically downloaded. Curators may also use the ABC to add non-
15 PubMed references. An additional key feature of the ABC is a search interface that allows
16 curators to retrieve references based on various criteria including their in/out of corpus status,
17 bibliographic data, and publication data range, if desired. Reference acquisition functionality can
18 easily be extended to integrate additional MODs into the Alliance infrastructure.

19

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

26

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

43

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

6
7 The Alliance is committed to creating Textpresso instances tailored to the unique needs of each
8 member database, all of which will be managed within the Alliance software ecosystem and
9 connected to the ABC as a single reference data source. This will reduce the overhead of
10 managing Textpresso at individual MODs while also simplifying development and deployment of
11 new features. Users will benefit from simplified access to Textpresso from the Alliance website.
12 We also plan to integrate Textpresso searches further into specific Alliance web pages such as
13 gene or allele pages. Users will be able to obtain additional references to biological entities
14 through Textpresso searches, adding information from potentially non-curated literature to the
15 list of curated references currently linked on those pages. Textpresso will be available to
16 Alliance biocurators and to the general public through MOD-customized websites and via APIs
17 for programmatic access.

18
19 **Artificial Intelligence (AI).** The Alliance member MODs have a track record of implementing
20 ML tools to enhance literature triage and curation efficiency. Notable examples include RGD's
21 early adoption of standard software architectures such as UIMA (Unstructured Information
22 Management Architecture, an Apache.org project) and the development of the OntoMate
23 system (Liu et al. 2015) an ontology-driven literature search engine, as well as WormBase's
24 creation of Textpresso (Mueller et al. 2004) and document classifiers for paper triage.

25
26 The rise of Large Language Models (LLMs), such as BERT (Bidirectional Encoder
27 Representations from Transformers), and ChatGPT, has transformed the natural language
28 processing (NLP) landscape, but questions about their accuracy and "hallucinations" remain.
29 The Alliance is developing LLMs for tasks such as document classification, Named Entity
30 Recognition (NER), sentence classification, computationally assisted triage and curation and to
31 build a natural language query system to simplify access to its underlying structured data.

32
33 Alliance members have developed AI/ML classifiers for determining with high accuracy whether
34 papers returned from automated PubMed queries should be kept in their corpus or discarded
35 (Jiang et al. 2020) and classifiers that can determine whether specific data types relevant for
36 curation are present in a document (Fang et al. 2012). The Alliance is developing a central
37 solution by providing these types of classifiers to all members.

38
39 Efforts are also underway to improve existing species-specific entity extraction and classification
40 models, with a focus on incorporating human feedback in the loop and continuously training
41 models based on data validated by professional biocurators and community curators. A
42 centralized interface for "topic and entity tag" addition and validation during literature triage and
43 curation is under development as part of the ABC. The interface allows curators to associate
44 tags with publications and at the same time validate (or invalidate) results extracted from AI/ML

1 methods. This interface will streamline the collection of valuable training and testing sets and
2 will allow a more systematic approach to the creation and comparison of different AI/ML models.

3
4 Furthermore, the Alliance is adopting Evidence and Conclusion Ontology (ECO) terms to record
5 systematically the type of evidence, e.g. neural network method evidence, and assertion
6 method, e.g. automatic assertion, used for reference flagging and triage. This is especially
7 relevant for topic and entity tags. Using ECO terms aligns with FAIR data principles and offers
8 transparency in curation workflows.

9 10 **Application Programming Interfaces (APIs)**

11 Application Programming Interfaces (APIs) are a key component of Alliance Central's data
12 services infrastructure for rapid, modular software development. We currently support a dozen
13 APIs with hundreds of endpoints (**Figure 11**). New APIs will be added as data harmonization
14 and modeling of additional data entities are completed. We will expand public site APIs to
15 generate all data needed for SimpleMine, AllianceMine, etc. from single endpoints. Current APIs
16 include Public site APIs (agr_java_software in the GitHub repo) and APIs available from a public
17 Swagger UI page. Because the public APIs support only GET endpoints, they do not require
18 authentication. All APIs that support both GET and PUT/POST/DELETE endpoints do require
19 authentication. Some of the key API endpoints available at
20 <https://www.alliancegenome.org/swagger-ui/> are: gene-summary, gene-disease, gene-
21 interactions, homologs-species, allele-phenotypes, expression ribbon-summary, etc.

22 23 **Data preservation in external repositories**

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

35
36 Ontologies that Alliance members maintain are also available from long-term repositories
37 including the OBO Foundry (<https://obofoundry.org/>) and Zenodo (zenodo.org).

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

42 The external repositories used by the Alliance of Genome Resources include the *OBO Foundry*
43 that was established in the early 2000s as a community-based initiative for development and
44 maintenance of biological and biomedical ontologies using standardized practices. The Foundry

1 is the ontology repository of choice for the Alliance because it is widely recognized as an
2 authoritative source of well-maintained ontologies for biology and biomedical research. *Zenodo*
3 is a general purpose repository maintained by CERN (European Council for Nuclear Research)
4 for storing and sharing documents, data, and other digital research materials across many
5 disciplines. Zenodo is a repository of choice for the Alliance, in part, because of the commitment
6 by the European Commission to support Zenodo as long as CERN exists.

7 8 **Outreach and interactions**

9
10 **The Alliance Helpdesk.** We established a common help desk email address
11 (help@alliancegenome.org) that is featured prominently on the Alliance website header and
12 footer under “Contact Us”. All inquiries submitted using this email are logged as tickets in the
13 Alliance Jira software system. Members of the User Support Working Group respond to user
14 questions and inquiries in a timely manner, typically within 48 hours. Time to resolve user
15 inquiries depends on the nature of the question or request. The Jira system tracks open tickets,
16 forward tickets, tracks their active/resolved status, and classifies them by subject. We use the
17 information, in part, to evaluate the design and utility of our user interfaces. For example, if
18 particular questions or subjects arise frequently, we re-evaluate the design and wording of the
19 search form and/or results display that caused user confusion.

20
21 **Online documentation.** We provide extensive user documentation about using the Alliance
22 data resources under the Help menu on the homepage (<https://www.alliancegenome.org/help>).
23 The online documentation provides guidance on such topics as how to use the search functions,
24 defines acceptable field parameters, and provides explanations of the displayed results. The
25 User Support Working Group also works closely with the User Interface Working Group and the
26 Developers to craft text for tooltips displayed on user interfaces.

27
28 **Frequently Asked Question (FAQ) pages.** The FAQ/Known Issues page provides answers to
29 commonly asked questions about the Alliance and also describes any known issues associated
30 with a particular software release. The link to the FAQ page is featured prominently on the
31 Alliance home page under the Help menu.

32
33 **Illustrated tutorials and videos.** We maintain several types of tutorial options that are
34 accessible from the Help menu (<https://www.alliancegenome.org/tutorials>). The most requested
35 types of tutorials are illustrated guides with screenshots on how to use various features of the
36 Alliance web portal. When new functionality is released, we post to social media channels and
37 issue “Tweertorials.” Short video tutorials are disseminated through the Alliance YouTube
38 channel.

39
40 **Alliance User Community Forum.** The Alliance supports a centralized community discussion
41 board implemented in Discourse (<https://community.alliancegenome.org/categories>) (**Figure**
42 **12**). Each model organism represented in the Alliance is represented as its own Discourse
43 category with model organism specific threads for news, discussion, and reagent information.
44 The forum also includes categories for job postings, meeting announcements, and general

1 information about the Alliance of Genome Resources. Alliance members with existing on-line
2 community forums are migrating users to the Alliance Central forum.

3
4 Users are not required to register to access the forum but must register to post messages,
5 questions, and announcements. On average, ~1,000 users a day access the forum. Posts
6 include jobs open and sought, news, meeting announcements and discussion of research
7 approaches, reagents and interpretation.

8
9 **Social Media.** In addition to a News and Events header that links to software release notes and
10 other Alliance Central updates, the Alliance uses standard social media venues to engage with
11 the user community, including FaceBook (www.facebook.com/alliancegenome/), Twitter (now,
12 X) (twitter.com/alliancegenome), Mastodon (<https://genomic.social/@AllianceGenome>), and
13 Bluesky (<https://bsky.app/profile/alliancegenome.bsky.social>).

14 15 **Prospects and Challenges**

16
17 **The tail of not-yet harmonized data.** One challenge in the central Alliance infrastructure
18 providing support for the union of MOD and GO features is the many unique dataset displays
19 and tools that have evolved in the individual MODs over two decades. Among the 8 resources
20 this comprises 150 years of branch length! Although horizontal tool transfer has occurred, it is
21 not complete. We are taking a few approaches to this problem. In some cases, where the data
22 are stand-alone, we will simply move the data and code to the Alliance. In the short term we will
23 likely run tools off their existing servers. As tools age out, we will evaluate whether there is a
24 broader mandate for that feature, and if so, implement it in the context of the Alliance.

25
26 There are types or aspects of our data that can be harmonized but have not yet been so. We
27 adopted LinkML to help with harmonization because it provides a common language to
28 represent structured data. The use of this language has spread to the point where our progress
29 on harmonization is much more rapid.

30
31 **AI.** As discussed above, we are actively considering AI/ML applications throughout the project.
32 Our practical approach is driven by us being subject matter experts. Because we have relied on
33 human expert curation, we are in a unique position to evaluate and use the output of various
34 AIs. Future plans include development of tools for creating training sets and a model manager
35 for tracking ML models' performance. Integration with specialized biocuration tools such as
36 Ontomate and Textpresso is part of the strategy, with a vision of harmonizing AI/ML solutions
37 across member sites.

38
39 We will also explore the use of AI/ML in gene function summarization. Included on gene pages
40 at the Alliance are short textual gene summaries based on curated and structured data
41 annotations that provide users a quick overview of gene function. The current automated
42 system for generating gene summaries has produced more than 160,000 summaries (Alliance
43 version 6.0.0) for nine species, including humans (Kishore et al. 2020). However, to increase
44 the coverage of genes further, we will explore the use of LLMs. This is especially relevant for

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

7
8 We will use prompt engineering and finetuning of LLMs to improve accuracy of the generated
9 summaries. As part of a continual improvement process, we will ask professional biocurators to
10 evaluate summaries, and we will develop a scoring system based on several features such as
11 readability of summaries, inclusion of key gene data, and checking for inaccurate and false
12 data. To improve and keep gene summaries up to date, we plan to retrieve newly published
13 articles that contain gene data that were not available when the LLM was trained and add
14 extracted relevant text from the identified articles to the LLM prompt. To do so, we will use tools
15 such as Textpresso (Muller et al. 2004) and Ontomate (Liu et al. 2015).

16
17 **Community curation.** Some Alliance MODs employ community curation pipelines to engage
18 authors in curation of their papers. For example, FlyBase utilizes the Fast Track Your Paper
19 (FTYP) (Bunt et al 2012; Larkin et al. 2021) tool that allows users to curate scientific papers,
20 identify data types, and associate relevant genes with the reference. Authors using FTYP to
21 ensure their papers appear quickly on the FlyBase website, help highlight data needing manual
22 curation, and prioritize their publication for further curation.

23
24 Similarly, WormBase developed ACKnowledge (Author Curation to Knowledgebase; Arnaboldi
25 et al. 2020), a semi-automated curation tool that lets authors curate their publications with the
26 help of ML. Authors receive an email with a link to a form pre-populated by document-level
27 classifiers that identify data types and several NER pipelines that extract lists of entities. Authors
28 can correct and validate the extracted data using the form and submit curated information to
29 WormBase. We will continue to provide these services to our community and develop a unified
30 infrastructure which will help expand the service to other member communities.

31
32 Several Alliance members also collaborate with publishing groups, such as microPublication
33 Biology (<https://www.micropublication.org/>) or the Genetics Society of America (<https://genetics-gsa.org/publications/>), to streamline pre-publication data integrity verification and curation by
34 curators and authors, enabling MODs to quality-check and work with authors to correct data
35 reporting before publication and promptly incorporate it into Alliance Knowledgebases upon
36 article publication.

37
38
39 **Dealing with satellite genomes and genetic models.** In addition to the core genomes and
40 associated data, our resources store and present information about the genes and genomes of
41 relatively closely related organisms. For example, WormBase includes some genetically-studied
42 nematodes such as *Caenorhabditis briggsae* that benefit from the rich data models typical of *C.*
43 *elegans*. Genetic screens and positional cloning (Inoue et al. 2007; Sharanya et al. 2012),
44 CRISPR editing (Cohen and Sternberg, 2019; Cohen et al. 2022; Ivanova and Moss 2023), as

1 well as transcriptomic analyses (Jhaveri et al. 2022) are now routinely done in this species. For
2 the Alliance to take on this responsibility of WormBase, we need to support such satellite model
3 organisms. Our plan is to support community gene structure annotation (e.g., for *Drosophila*,
4 Sargent et al, 2020; for *C. elegans*, Moya et al. 2023) using the Apollo curaton system designed
5 specifically for such activity (Dunn et al. 2019).

6
7 **High Throughput expression data and single cell RNA-seq plans.** We harmonized high-
8 throughput expression metadata of mouse, rat, yeast, worm, fly, and zebrafish. Users can
9 browse them with species, assay type (microarray, RNA-seq, tiling array, and proteomics),
10 tissue, sex, and curated categories. We plan to add single-cell RNA-seq as a new assay type,
11 making such datasets easily identifiable within our collection, with links to other resources,
12 including Gene Expression Omnibus, EBI single-cell RNA-seq Expression Atlas, and CZI
13 CellxGene, to display the information above, we will implement a content-rich expression detail
14 page that will provide a unified way to access all expression data associated with a specific
15 gene, including link outs to other sources and MOD-specific single-cell RNA-seq gene
16 expression graphs (**Figure 13**).

17
18 **Disease Portal(s).** Providing users with ready and easy access to curated and harmonized
19 model organism disease data and tools is crucial to accelerate research related to the
20 pathogenesis of human disease. The Alliance has a wealth of disease-relevant data from eight
21 model organism species and human data, such as: genes, alleles and variants implicated in
22 disease, genotypes and strains that serve as disease models, and related data such as
23 modifiers (herbals, chemicals, small molecules, etc.) that ameliorate or exacerbate the disease
24 condition and may serve as candidates for potential drug development. To provide an easy
25 entry point for clinical researchers and human geneticists to access the consolidated data and
26 tools, we are in the process of designing and implementing a topic-specific resource--an
27 Alzheimer's disease (AD) portal that will serve as a paradigm for other disease portals (**Figure**
28 **14**). The AD portal will include: orthologous genes in animal model systems, models with a
29 mutation orthologous to one in a patient group, models with a specific set of phenotypes, and/or
30 modifiers that have been shown to alter the disease condition. Building on the experience and
31 pages developed for the AD portal, we will expand this paradigm to other disease portals. In
32 addition to the specific disease portals we also plan to provide "compare" functionalities across
33 diseases. Features planned for the disease portal with AD as an example include: a home page
34 with an overview of the data in the portal, an autocomplete search box, links to other AD
35 resources, and a list of the most recent papers from PubMed and/or from the ABC store (see
36 example portal page below). The pages in the portal will be modeled on existing pages at the
37 Alliance and will include gene summaries, alleles and variants, phenotypes, gene interactions,
38 pathways, biological processes (based on GO), gene expression, etc. We also plan to provide
39 visualizations of data analysis, for example, diseases that share genes and protein interactions
40 that may point to common underlying molecular mechanisms. Up-to-date data sets, e.g., genes,
41 strains, modifiers (drugs, chemicals, herbals, etc. shown to either ameliorate or exacerbate
42 phenotypes) will be available as downloadable files. Disease-specific data sets will also be
43 available for query from AllianceMine. We will also provide up-to-date links to disease-specific
44 literature, and search capabilities through literature search engines such as the Textpresso

1 instance dedicated to AD (<http://alzheimer.textpresscentral.org>; corpus size - 96,000 papers).
2 Not all papers are curatable by the MODs given their extensive but not comprehensive data
3 models, and thus literature search will remain important.
4

5 **The Alliance in the ecosystem of knowledgebases.** The Alliance has a unique and
6 complementary role relative to other informatics resources that support comparative biology. For
7 example, NCBI's new Comparative Genomics Resource (CGR; Bornstein et al 2023) focuses
8 on developing analysis tools and resources for *sequence-based* genome comparisons across a
9 large number of species, the Alliance focuses on standardized annotations, harmonized
10 biological concepts, and comparison of *biological knowledge*. The CGR supports comparative
11 sequence analysis for all eukaryotes whereas the Alliance is primarily focused on model
12 organisms used widely in biomedical research. These model organisms have a tremendous
13 amount of highly valuable genetic, transgenic, and phenotypic data generated with multiple
14 types of assays and are uniquely represented by the Alliance Knowledge Centers. The CGR
15 uses the standardized gene summaries from the Alliance and follows nomenclature and
16 ontology standards developed and maintained by Alliance members. For sequence analysis, the
17 Alliance leverages sequence-based analysis tools developed and maintained by the CGR.
18 Resource developers by and large appreciate the magnitude of the tasks we face in order to
19 provide researchers with the information they need and strive to fill in the many gaps in
20 services.
21
22

23 **Acknowledgements**

24 We thank our multiple communities for their patience and feedback about the prospect of the
25 Alliance and their love of their own MODs. We also thank the members of our Scientific
26 Advisory Board (Gary Bader, Alex Bateman, Helen Berman, Shawn Burgess, Andrew Chisholm,
27 Phil Hieter, Brian Oliver, Calum Macrae, Titus Brown, Abraham Palmer and Michelle Southard-
28 Smith) for cogent advice, and NHGRI Program Staff (Sandhya Xirasagar, Ajay Pillai, Valentina
29 di Francesco, Sarah Hutchison, and Helen Thompson) for guidance.
30

31 **Funding**

32 The core funding for the Alliance is from the National Human Genome Research Institute and
33 the National Heart, Lung and Blood Institute (U24HG010859). The curation of data and their
34 harmonization is supported by National Human Genome Research Institute grants
35 U24HG002659 (ZFIN), U24HG002223 (WormBase), U41HG000739 (FlyBase), U24HG001315
36 (SGD), U24HG000330 (MGD), P41HD064556 (Xenbase), U24HG011851 (Reactome + GO)
37 and U41HG012212 (GO Consortium), as well as grant R01HL064541 from the National Heart,
38 Lung and Blood Institute (RGD), P41HD062499 from the Eunice Kennedy Shriver National
39 Institute of Child Health and Human Development (GXD), and the Medical Research Council-
40 UK grant MR/L001020/1 (WormBase). Additional effort was supported by DOE DE-AC02-
41 05CH11231. Curation tools are supported in part by the National Library of Medicine NLM
42 R01LM013871.
43
44

1 **Competing Interests**

2 The authors declare no competing interests.

4 **References**

- 6 Alliance of Genome Resources, C., *Harmonizing model organism data in the Alliance of*
7 *Genome Resources*. Genetics, 2022. **220**(4).
- 8 Altenhoff AM, Train CM, Gilbert KJ, Mediratta I, Mendes de Farias T, Moi D, Nevers Y,
9 Radoykova HS, Rossier V, Warwick Vesztrocy A, Glover NM, Dessimoz C. OMA orthology in
10 2021: website overhaul, conserved isoforms, ancestral gene order and more. Nucleic Acids
11 Res. 2021 Jan 8;49(D1):D373-D379. doi: 10.1093/nar/gkaa1007. PMID: 33174605; PMCID:
12 PMC7779010.
- 13 Altschul SF, Gish W, Miller W, Eugene W. Myers, Lipman DJ. 1990. Basic local alignment
14 search tool. *J Mol Biol* **215**: 403-410.
- 15 Anderson, W.P. and G. Global Life Science Data Resources Working, *Data management: A*
16 *global coalition to sustain core data*. Nature, 2017. **543**(7644): p. 179.
- 17 Bornstein, K., et al., *The NIH Comparative Genomics Resource: addressing the promises and*
18 *challenges of comparative genomics on human health*. BMC Genomics, 2023. **24**(1): p. 575.
- 19 Bowes JB, Snyder KA, James-Zorn C, Ponferrada VG, Jarabek CJ, Burns KA, Bhattacharyya B,
20 Zorn AM, Vize PD. The Xenbase literature curation process. Database (Oxford). 2013 Jan
21 9;2013:bas046. doi: 10.1093/database/bas046. PMID: 23303299; PMCID: PMC3540419.
- 22 Bradford, Y.M., et al., *From multiallele fish to nonstandard environments, how ZFIN assigns*
23 *phenotypes, human disease models, and gene expression annotations to genes*. Genetics,
24 2023. **224**(1).
- 25 Bult, C.J. and P.W. Sternberg, *The alliance of genome resources: transforming comparative*
26 *genomics*. Mamm Genome, 2023.
- 27 Bunt SM, Grumbling GB, Field HI, Marygold SJ, Brown NH, Millburn GH; FlyBase Consortium.
28 Directly e-mailing authors of newly published papers encourages community curation. Database
29 (Oxford). 2012 May 2;2012:bas024. doi: 10.1093/database/bas024. PMID: 22554788; PMCID:
30 PMC3342516.
- 31 Carotenuto R, Pallotta MM, Tussellino M, Fogliano C. *Xenopus laevis* (Daudin, 1802) as a
32 Model Organism for Bioscience: A Historic Review and Perspective. Biology (Basel). 2023 Jun
33 20;12(6):890. doi: 10.3390/biology12060890. PMID: 37372174; PMCID: PMC10295250.
- 34 Cohen S, Sternberg P. Genome editing of *Caenorhabditis briggsae* using CRISPR/Cas9 co-
35 conversion marker *dpy-10*. MicroPubl Biol. 2019 Oct
36 11;2019:10.17912/micropub.biology.000171. doi: 10.17912/micropub.biology.000171. PMID:
37 32550401; PMCID: PMC7252229.
- 38 Cohen SM, Wrobel CJJ, Prakash SJ, Schroeder FC, Sternberg PW. Formation and function of
39 dauer ascarosides in the nematodes *Caenorhabditis briggsae* and *Caenorhabditis elegans*. G3
40 (Bethesda). 2022 Mar 4;12(3):jkac014. doi: 10.1093/g3journal/jkac014. PMID: 35094091;

1 PMCID: PMC8895998.

2 Cosentino S, Iwasaki W. SonicParanoid: fast, accurate and easy orthology inference.
3 Bioinformatics. 2019 Jan 1;35(1):149-151. doi: 10.1093/bioinformatics/bty631. PMID: 30032301;
4 PMCID: PMC6298048.

5 Davis, P., et al., *WormBase in 2022-data, processes, and tools for analyzing Caenorhabditis*
6 *elegans*. Genetics, 2022. **220**(4).

7 Dunn NA, Unni DR, Diesh C, Munoz-Torres M, Harris NL, Yao E, et al. Apollo: democratizing
8 genome annotation. PLoS Comput Biol. 2019;15:e1006790.

9 Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for comparative genomics.
10 Genome Biol. 2019 Nov 14;20(1):238. doi: 10.1186/s13059-019-1832-y. PMID: 31727128;
11 PMCID: PMC6857279.

12 Engel SR, Wong ED, Nash RS, Aleksander S, Alexander M, Douglass E, Karra K, Miyasato SR,
13 Simison M, Skrzypek MS, Weng S, Cherry JM. New data and collaborations at the
14 Saccharomyces Genome Database: updated reference genome, alleles, and the Alliance of
15 Genome Resources. Genetics. 2022 Apr 4;220(4):iyab224. doi: 10.1093/genetics/iyab224.
16 PMID: 34897464; PMCID: PMC9209811.

17 Fang R, Schindelman G, Van Auken K, Fernandes J, Chen W, Wang X, Davis P, Tuli MA,
18 Marygold SJ, Millburn G, Matthews B, Zhang H, Brown N, Gelbart WM, Sternberg PW.
19 Automatic categorization of diverse experimental information in the bioscience literature. BMC
20 Bioinformatics. 2012 Jan 26;13:16. doi: 10.1186/1471-2105-13-16. PMID: 22280404; PMCID:
21 PMC3305665.

22 Fisher M, James-Zorn C, Ponferrada V, Bell AJ, Sundararaj N, Segerdell E, Chaturvedi P,
23 Bayyari N, Chu S, Pells T, Lotay V, Agalakov S, Wang DZ, Arshinoff BI, Foley S, Karimi K, Vize
24 PD, Zorn AM. Xenbase: key features and resources of the Xenopus model organism
25 knowledgebase. Genetics. 2023 May 4;224(1):iyad018. doi: 10.1093/genetics/iyad018. PMID:
26 36755307; PMCID: PMC10158840.

27 FlyBase C. 1999. The FlyBase database of the Drosophila Genome Projects and community
28 literature. *Nucleic Acids Res* **27**: 85-88.

29 Fuentes D, Molina M, Chorostecki U, Capella-Gutiérrez S, Marcet-Houben M, Gabaldón T.
30 PhylomeDB V5: an expanding repository for genome-wide catalogues of annotated gene
31 phylogenies. *Nucleic Acids Res*. 2022 Jan 7;50(D1):D1062-D1068. doi: 10.1093/nar/gkab966.
32 PMID: 34718760; PMCID: PMC8728271.

33 Gene Ontology, Consortium., *The Gene Ontology knowledgebase in 2023*. Genetics, 2023.
34 **224**(1).

35 Gramates, L.S., et al., *FlyBase: a guided tour of highlighted features*. Genetics, 2022. **220**(4).

36 Howe, D.G., et al., *Model organism data evolving in support of translational medicine*. Lab Anim
37 (NY), 2018. **47**(10): p. 277-289.

38 Hu Y, Comjean A, Rodiger J, Liu Y, Gao Y, Chung V, Zirin J, Perrimon N, Mohr SE.
39 FlyRNAi.org-the database of the Drosophila RNAi screening center and transgenic RNAi

1 project: 2021 update. *Nucleic Acids Res.* 2021 Jan 8;49(D1):D908-D915. doi:
2 10.1093/nar/gkaa936. PMID: 33104800; PMCID: PMC7778949.

3 Hu Y, Flockhart I, Vinayagam A, Bergwitz C, Berger B, Perrimon N, Mohr SE. An integrative
4 approach to ortholog prediction for disease-focused and other functional studies. *BMC*
5 *Bioinformatics.* 2011 Aug 31;12:357. doi: 10.1186/1471-2105-12-357. PMID: 21880147; PMCID:
6 PMC3179972.

7 Inoue T, Ailion M, Poon S, Kim HK, Thomas JH, Sternberg PW. Genetic analysis of dauer
8 formation in *Caenorhabditis briggsae*. *Genetics.* 2007 Oct;177(2):809-18. doi:
9 10.1534/genetics.107.078857. Epub 2007 Jul 29. PMID: 17660533; PMCID: PMC2034645.

10 Ivanova M, Moss EG. Orthologs of the *C. elegans* heterochronic genes have divergent functions
11 in *C. briggsae*. *Genetics.* 2023 Oct 3:iyad177. doi: 10.1093/genetics/iyad177. Epub ahead of
12 print. PMID: 37788363.

13 Jhaveri N, van den Berg W, Hwang BJ, Muller HM, Sternberg PW, Gupta BP. Genome
14 annotation of *Caenorhabditis briggsae* by TEC-RED identifies new exons, paralogs, and
15 conserved and novel operons. *G3 (Bethesda).* 2022 Jul 6;12(7):jkac101. doi:
16 10.1093/g3journal/jkac101. PMID: 35485953; PMCID: PMC9258526.

17 Jiang X, Li P, Kadin J, Blake JA, Ringwald M, Shatkay H. Integrating image caption information
18 into biomedical document classification in support of biocuration. *Database (Oxford).* 2020 Jan
19 1;2020:baaa024. doi: 10.1093/database/baaa024. PMID: 32294192; PMCID: PMC7159034.

20 Kostiuk V, Khokha MK. *Xenopus* as a platform for discovery of genes relevant to human
21 disease. *Curr Top Dev Biol.* 2021;145:277-312. doi: 10.1016/bs.ctdb.2021.03.005. Epub 2021
22 Apr 23. PMID: 34074532; PMCID: PMC8734201.

23 Larkin A, Marygold SJ, Antonazzo G, Attrill H, Dos Santos G, Garapati PV, Goodman JL,
24 Gramates LS, Millburn G, Strelets VB, Tabone CJ, Thurmond J; FlyBase Consortium. FlyBase:
25 updates to the *Drosophila melanogaster* knowledge base. *Nucleic Acids Res.* 2021 Jan
26 8;49(D1):D899-D907. doi: 10.1093/nar/gkaa1026. PMID: 33219682; PMCID: PMC7779046.
27 (<https://pubmed.ncbi.nlm.nih.gov/33219682/>)

28 Liu W, Laulederkind SJ, Hayman GT, Wang SJ, Nigam R, Smith JR, De Pons J, Dwinell MR,
29 Shimoyama M. OntoMate: a text-mining tool aiding curation at the Rat Genome Database.
30 *Database (Oxford).* 2015 Jan 25;2015:bau129. doi: 10.1093/database/bau129. PMID:
31 25619558; PMCID: PMC4305386.

32 Milacic M, Rothfels K, Mathews L, Wright A, Jassal B, Shamovsky V, Trinh Q, Gillespie M,
33 Sevilla C, Tiwari K, Ragueneau E, Gong C, Stephan1 R, May B, Haw R, Weiser J, Beavers D,
34 Conley P, Hermjakob H, Stein LD, D'Eustachio P, Wu G (2024) The Reactome Pathway
35 Knowledgebase 2024. *Nucleic Acids Res.*, in press. PMID: 37941124

36 Mitros T, Lyons JB, Session AM, Jenkins J, Shu S, Kwon T, Lane M, Ng C, Grammer TC,
37 Khokha MK, Grimwood J, Schmutz J, Harland RM, Rokhsar DS. A chromosome-scale genome
38 assembly and dense genetic map for *Xenopus tropicalis*. *Dev Biol.* 2019 Aug 1;452(1):8-20. doi:
39 10.1016/j.ydbio.2019.03.015. PMID: 30980799.

40 Moya ND, Stevens L, Miller IR, Sokol CE, Galindo JL, Bardas AD, Koh ESH, Rozenich J, Yeo

1 C, Xu M, Andersen EC. Novel and improved *Caenorhabditis briggsae* gene models generated
2 by community curation. *BMC Genomics*. 2023 Aug 25;24(1):486. doi: 10.1186/s12864-023-
3 09582-0. PMID: 37626289; PMCID: PMC10463891.

4 Müller HM, Van Auken KM, Li Y, Sternberg PW. Textpresso Central: a customizable platform for
5 searching, text mining, viewing, and curating biomedical literature. *BMC Bioinformatics*. 2018
6 Mar 9;19(1):94. doi: 10.1186/s12859-018-2103-8. PMID: 29523070; PMCID: PMC5845379.

7 Nevers Y, Jones TEM, Jyothi D, Yates B, Ferret M, Portell-Silva L, Codo L, Cosentino S,
8 Marcet-Houben M, Vlasova A, Poidevin L, Kress A, Hickman M, Persson E, Piližota I, Guijarro-
9 Clarke C; OpenEBench team the Quest for Orthologs Consortium; Iwasaki W, Lecompte O,
10 Sonnhammer E, Roos DS, Gabaldón T, Thybert D, Thomas PD, Hu Y, Emms DM, Bruford E,
11 Capella-Gutierrez S, Martin MJ, Dessimoz C, Altenhoff A. The Quest for Orthologs orthology
12 benchmark service in 2022. *Nucleic Acids Res*. 2022 Jul 5;50(W1):W623-W632. doi:
13 10.1093/nar/gkac330. PMID: 35552456; PMCID: PMC9252809.

14 Nevers Y, Kress A, Defosset A, Ripp R, Linard B, Thompson JD, Poch O, Lecompte O.
15 OrtholInspector 3.0: open portal for comparative genomics. *Nucleic Acids Res*. 2019 Jan
16 8;47(D1):D411-D418. doi: 10.1093/nar/gky1068. PMID: 30380106; PMCID: PMC6323921.

17 Oliver, S.G., et al., *Model organism databases: essential resources that need the support of*
18 *both funders and users*. *BMC Biol*, 2016. **14**: p. 49.

19 Persson E, Sonnhammer ELL. InParanoid-DIAMOND: faster orthology analysis with the
20 InParanoid algorithm. *Bioinformatics*. 2022 May 13;38(10):2918-2919. doi:
21 10.1093/bioinformatics/btac194. PMID: 35561192; PMCID: PMC9113356.

22 Priyam A, Woodcroft BJ, Rai V, Moghul I, Munagala A, Ter F, Chowdhary H, Pieniak I, Maynard
23 LJ, Gibbins MA et al. 2019. Sequenceserver: A Modern Graphical User Interface for Custom
24 BLAST Databases. *Mol Biol Evol* **36**: 2922-2924.

25 Ringwald, M., et al., *Mouse Genome Informatics (MGI): latest news from MGD and GXD*.
26 *Mamm Genome*, 2022. **33**(1): p. 4-18.

27 Sargent L, Liu Y, Leung W, Mortimer NT, Lopatto D, Goecks J, Elgin SCR. G-OnRamp:
28 Generating genome browsers to facilitate undergraduate-driven collaborative genome
29 annotation. *PLoS Comput Biol*. 2020 Jun 4;16(6):e1007863. doi: 10.1371/journal.pcbi.1007863.
30 PMID: 32497138; PMCID: PMC7272004.

31 Session AM, Uno Y, Kwon T, Chapman JA, Toyoda A, Takahashi S, Fukui A, Hikosaka A,
32 Suzuki A, Kondo M, van Heeringen SJ, Quigley I, Heinz S, Ogino H, Ochi H, Hellsten U, Lyons
33 JB, Simakov O, Putnam N, Stites J, Kuroki Y, Tanaka T, Michiue T, Watanabe M, Bogdanovic
34 O, Lister R, Georgiou G, Paranjpe SS, van Kruijsbergen I, Shu S, Carlson J, Kinoshita T, Ohta
35 Y, Mawaribuchi S, Jenkins J, Grimwood J, Schmutz J, Mitros T, Mozaffari SV, Suzuki Y,
36 Haramoto Y, Yamamoto TS, Takagi C, Heald R, Miller K, Haudenschild C, Kitzman J,
37 Nakayama T, Izutsu Y, Robert J, Fortriede J, Burns K, Lotay V, Karimi K, Yasuoka Y, Dichmann
38 DS, Flajnik MF, Houston DW, Shendure J, DuPasquier L, Vize PD, Zorn AM, Ito M, Marcotte
39 EM, Wallingford JB, Ito Y, Asashima M, Ueno N, Matsuda Y, Veenstra GJ, Fujiyama A, Harland
40 RM, Taira M, Rokhsar DS. Genome evolution in the allotetraploid frog *Xenopus laevis*. *Nature*.

- 1 2016 Oct 20;538(7625):336-343. doi: 10.1038/nature19840. PMID: 27762356; PMCID:
2 PMC5313049.
- 3 Sharanya D, Thillainathan B, Marri S, Bojanala N, Taylor J, Flibotte S, Moerman DG, Waterston
4 RH, Gupta BP. Genetic control of vulval development in *Caenorhabditis briggsae*. *G3*
5 (Bethesda). 2012 Dec;2(12):1625-41. doi: 10.1534/g3.112.004598. Epub 2012 Dec 1. PMID:
6 23275885; PMCID: PMC3516484.
- 7 Smith RN, Aleksic J, Butano D, Carr A, Contrino S, Hu F, Lyne M, Lyne R, Kalderimis A,
8 Rutherford K, Stepan R, Sullivan J, Wakeling M, Watkins X, Micklem G. InterMine: a flexible
9 data warehouse system for the integration and analysis of heterogeneous biological data.
10 *Bioinformatics*. 2012 Dec 1;28(23):3163-5. doi: 10.1093/bioinformatics/bts577. Epub 2012 Sep
11 27. PMID: 23023984; PMCID: PMC3516146.
- 12 Sternberg et al., 2024. WormBase article, Genetics
- 13 Thomas PD, Ebert D, Muruganujan A, Mushayahama T, Albou LP, Mi H. PANTHER: Making
14 genome-scale phylogenetics accessible to all. *Protein Sci*. 2022 Jan;31(1):8-22. doi:
15 10.1002/pro.4218. Epub 2021 Nov 25. PMID: 34717010; PMCID: PMC8740835.
- 16 Thomas, P.D., et al., *Gene Ontology Causal Activity Modeling (GO-CAM) moves beyond GO*
17 *annotations to structured descriptions of biological functions and systems*. *Nat Genet*, 2019.
18 **51**(10): p. 1429-1433.
- 19 UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Res*.
20 2023 Jan 6;51(D1):D523-D531. doi: 10.1093/nar/gkac1052. PMID: 36408920; PMCID:
21 PMC9825514.
- 22
- 23 Van Auken K, Fey P, Berardini TZ, Dodson R, Cooper L, Li D, Chan J, Li Y, Basu S, Muller HM,
24 Chisholm R, Huala E, Sternberg PW; WormBase Consortium. Text mining in the biocuration
25 workflow: applications for literature curation at WormBase, dictyBase and TAIR. *Database*
26 (Oxford). 2012 Nov 17;2012:bas040. doi: 10.1093/database/bas040. PMID: 23160413; PMCID:
27 PMC3500519.
- 28 Vedi, M., et al., *2022 updates to the Rat Genome Database: a Findable, Accessible,*
29 *Interoperable, and Reusable (FAIR) resource*. *Genetics*, 2023. **224**(1).
- 30 Wood V, Sternberg PW, Lipshitz HD. Making biological knowledge useful for humans and
31 machines. *Genetics*. 2022 Apr 4;220(4):iyac001. doi: 10.1093/genetics/iyac001.

32
33
34
35
36
37
38
39
40

Legends

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

Figure 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

1 the similarity and identity percentages of genes paralogous to *hlh-25*. The methodology count,
2 indicating the number of algorithms supporting the paralogous relationship, is also included. In
3 this ranking, *hlh-27* is identified as the primary paralog due to its high similarity and identity
4 scores, despite being recognized by fewer methods than *hlh-28*.

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

8
9 **Figure 4. Screenshot of results from the Alliance SequenceServer BLAST tool.** The results
10 have been enhanced relative to the default Sequence Server results page by the addition of
11 links to Alliance JBrowse and to the corresponding gene page (in this case, *C. elegans* *abi-1*) at
12 the Alliance website for each BLAST hit.

13
14 **Figure 5. Output of a BLAST search** After a user clicks on the JBrowse link for a BLAST hit
15 they are directed to the web service where they will see a track for the BLAST hit and how the
16 hit aligns with other tracks.

17
18
19 **Figure 6. AllianceMine example.** Using a simple template, a disease ontology (DO) term, in
20 this case “autism,” is chosen, and all genes associated with this DO term are returned in a
21 downloadable table.

22
23 **Figure 7. Alliance Pathway Viewer.** The pathway widget displays gene products (light purple
24 rectangles), and chemicals (light blue rectangles) and the flow of information and material
25 between them (relations). These relations, shown in legend indicate direct or indirect regulation
26 that can be positive, negative or of unknown effect direction. For metabolites, grey-blue shows
27 that a metabolite mediates the information flow between gene products. In addition, blue lines
28 with circles indicate input to a reaction; pink indicates output of a reaction.

29
30 **Figure 8. Evolution of Data Flow.** Graphical summary showing the design of short term
31 infrastructure initially deployed to support rapid delivery of unified data to the community and the
32 planned production system. Red, data quartermasters at MODs; Yellow, data; Brown, database;
33 Green, transformations; Blue, user interface.

34
35 **Figure 9. Alliance Curation tool.** Screenshot of the Alliance curation tool interface showing an
36 example of curated annotations of Affected Genomic Models managed in the persistent store.

37
38 **Figure 10. Textpresso for SGD literature at the Alliance.** ([http://sgd-
39 textpresso.alliancegenome.org/tpc/search](http://sgd-textpresso.alliancegenome.org/tpc/search))

40
41 **Figure 11. Swagger interface for the Alliance APIs.**

42
43 **Figure 12. Alliance community forum home page.**

44

1 **Figure 13. Mockup of an Expression Detail page.** This example shows one of the current
2 features of WormBase – single cell data from two studies – displayed on what will be part of an
3 Alliance Gene Expression detail page.

4
5 **Figure 14. Mockup of the Alzheimer’s Disease Portal showing the Home page and the**
6 **Data access page.** These views illustrate the type of information that will be available with a
7 disease-focus.

8
9
10

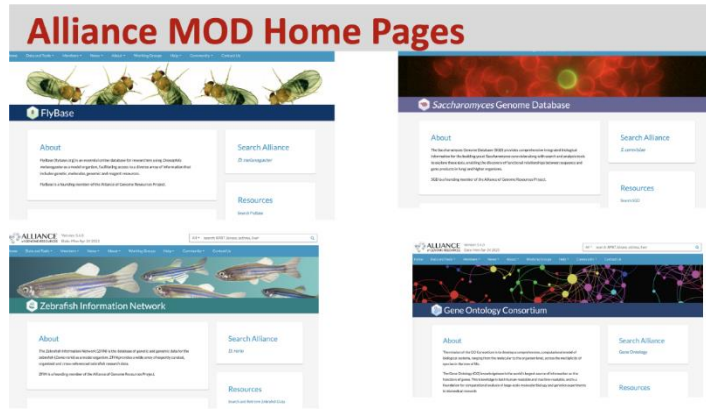


Figure 1
94x53 mm (x DPI)

1
2
3
4

Paralogy 

Gene symbol	Rank	Alignment Length (aa)	Similarity %	Identity %	Method Count	Method								
						Exonix Compare	HyPhy	tblastn	OMA	OrthoFinder	OrthoMCL	OrthoVenn	OrthoVenn2	OrthoVenn3
hnh-27	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 type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
hnh-28	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 type="checkbox"/>	<input type="checkbox"/>
hnh-29	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 type="checkbox"/>	<input type="checkbox"/>
hnh-26	4	274	48	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 type="checkbox"/>	<input type="checkbox"/>
ref-1	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 type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Figure 2
82x33 mm (x DPI)

1
2
3
4

Sequence Details ?

Transcript: Mode

Copy plain text

Copy with highlights

>NM_001031689.3-gene

```
AGGCTTGCTGTGTCACCTCGGCCCGCTCGGCGCGCCG  
GGCCGCCCTTACCTGCAGGCTTCTCCCGCCGCG  
CTGGGCACCGGGCGCCAGACAGACACTGGCCATGAC  
AGCTGGACGTACGGGGCTGGTGTGCTGCGCCTATC  
TGGGCCCGACAGGTGAGCGCTGGGAGTCGGGTTG  
TCCCTGTCAGTCTGccgtctctctctccccccagcCTTCCCTCTGCTCTCCGCTCCCTTCCAATTCTCAGACTATTAGAAC  
TCTGTAAGAAACATCGGGATTTAAGTGaaagagcacaggggtggggacCAAAGACCTGCCTTGTGTGCTAACTTGACCA  
CAGACGAGTCTCCTACCTTTTGGGCTTCAGTTTTGGACGATGATCTCCAAGTTCTTTAGACTTGAAAATTTACTGATT  
GCAGTTGCACCCCTCCGAAGTGAGGTAGTTTGAAGGCATCTGaaatgtcctcttttttttttttttttgcgaagaAGATGCTCTGT  
AGTCTtctgtaaaatttaatttttgaagactTTAGTTCTCAAAAATTGCACCTGGTGAATCCTCTTTTCGTCCAGGTAGAAAT  
TTAGTCTGTGGTCTGTCTGGTCTGAGGGAGTGAAACCTCTCGGATGTTTTGTTTCTGTGCATGTGCTGTTTCTTGAGGAGA  
AGCAGCATCCATTGCCTTCAAAGGATTTATCAGAAGGGTTCacgaacaaaaaaaaaagaagaaaagggttaGGAATCAGT  
CTGATCGAGTTCACGGTTCAGCCCTGATTTGGCTGGTGTAAACAGGATATTTAAGACCTAGAAGACAGATTGcagttcagag  
aaagaaaaattgaggttagttattttggtatttagtaGGTCTCCACTGCTAGAGATTTAGAAATTTGAGTCACCATCCATAA  
ATTCAGTTGATAACTGTTGAGTGCCTCCTTTTGGTGGGATACTGGAGAGGAATAAAGACAGACAAGTTGCCAGTGTTCCTG  
GAGCTTTCTTACAGTCTGATGGGgaagatagattaaaaacaagccaataaataattaaaaactggCATGTGAAGAAAACATA
```

- ✓ 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

- up/downstream
- TR
- oding
- intron
- Genomic (i.e., unprocessed)
- Amino acid

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

1
2
3
4

Figure 3
165x87 mm (x DPI)

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)

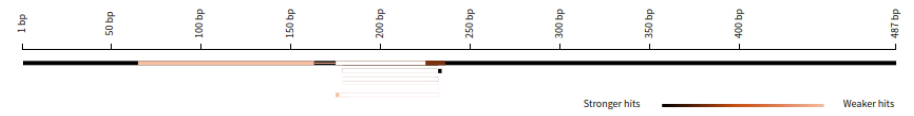
Parameters: task blastn, evaluate 1e-05, sc-match 2, sc-mismatch -3, gap-open 5, gap-extend 2, filter L;m;

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

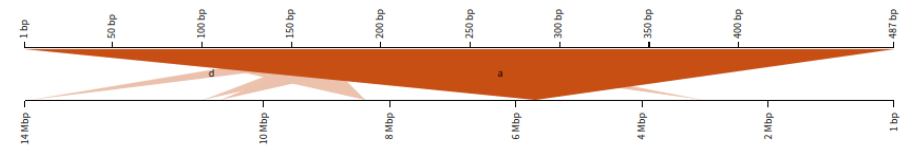
Similar sequences

	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%

III 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 CTGAAAATAATTTGCTTTTCGTGTTTTGACAAAACGTTTTCAAAAAAAAAAGGGAGCGAAAAATCTGACATAACTTATACAT 84
          |||
Subject 5691930 CTGAAAATAATTTGCTTTTCGTGTTTTGACAAAACGTTTTCAAAAAAAAAAGGGAGCGAAAAATCTGACATAACTTATACAT 5691847
    
```

1
2
3
4

Figure 4
165x107 mm (x DPI)

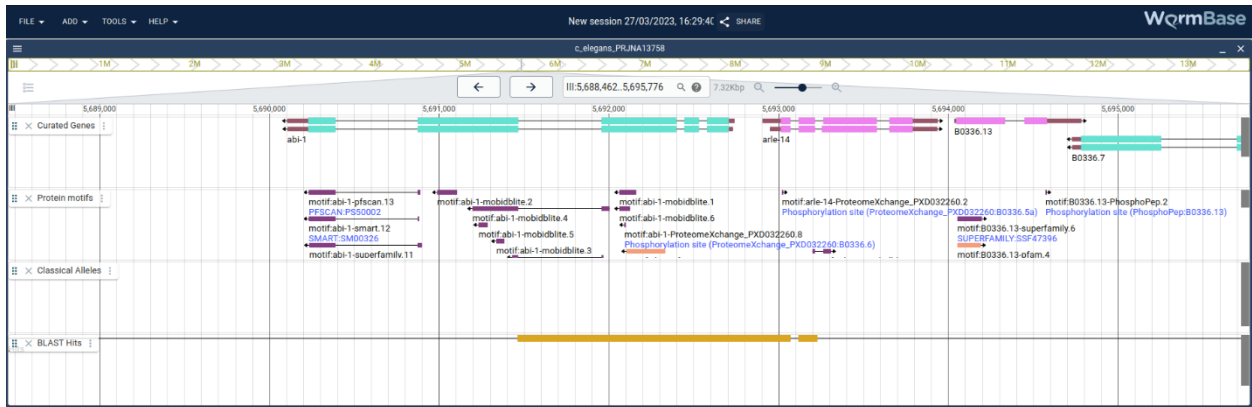


Figure 5
165x55 mm (x DPI)

1
2
3
4

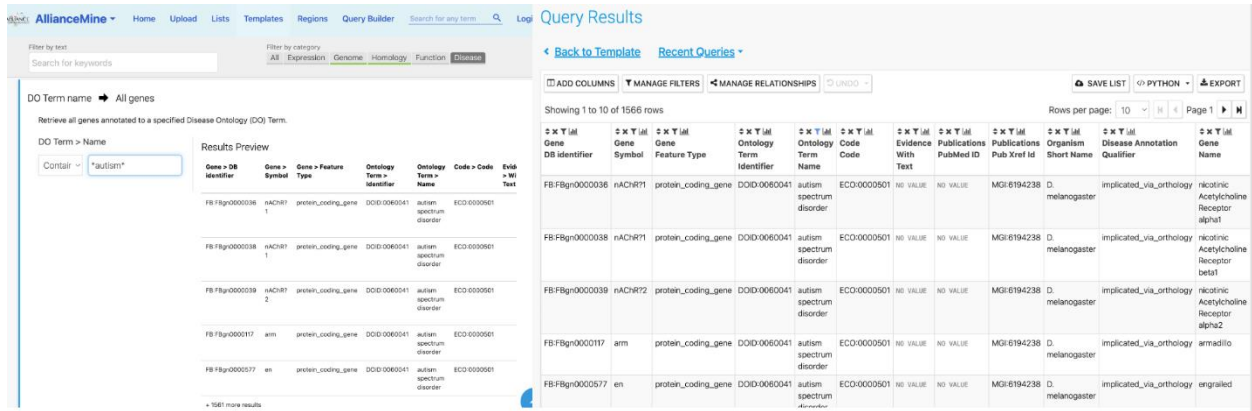


Figure 6
165x54 mm (x DPI)

1
2
3
4

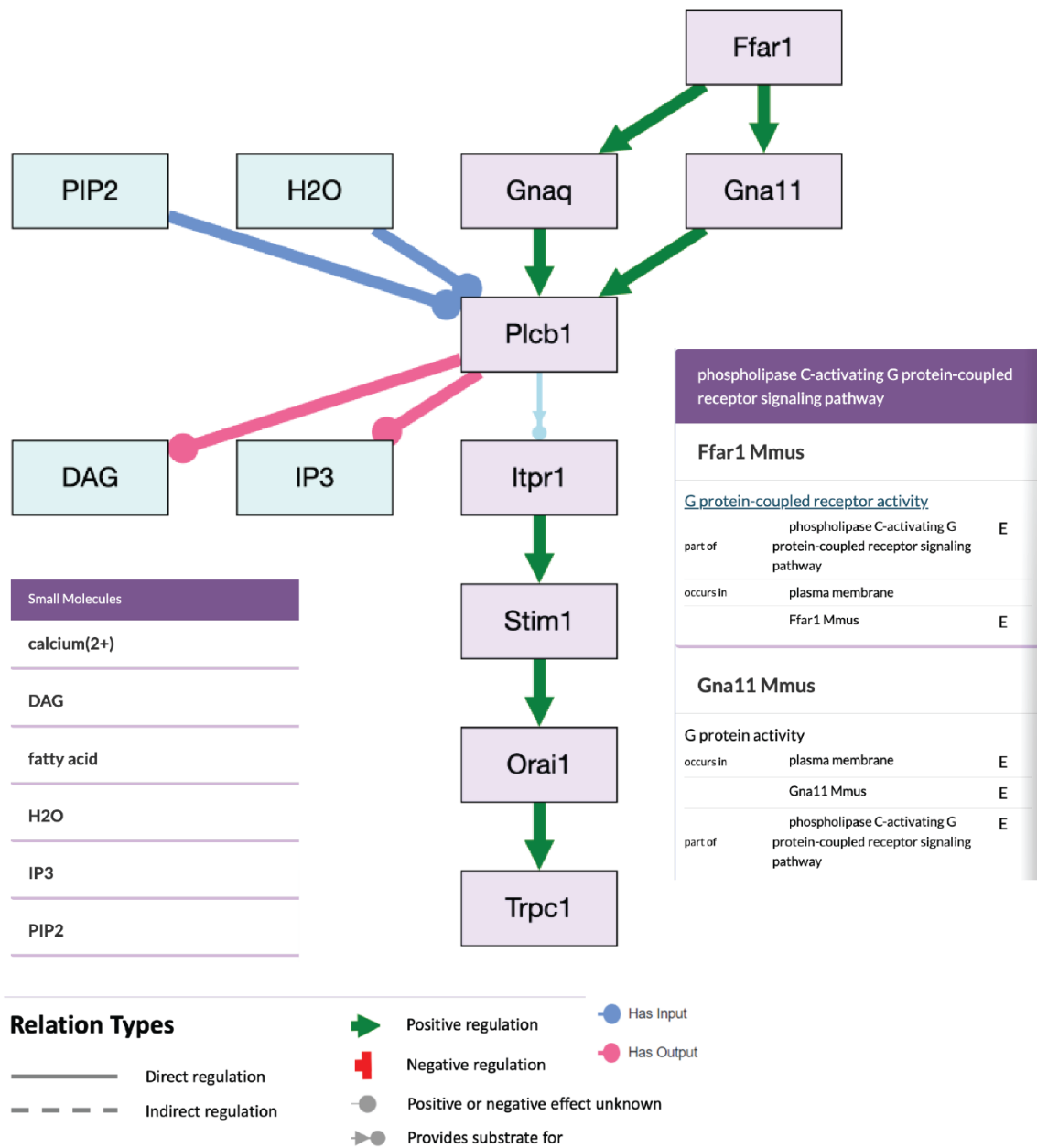


Figure 7
165x173 mm (x DPI)

1
2
3
4

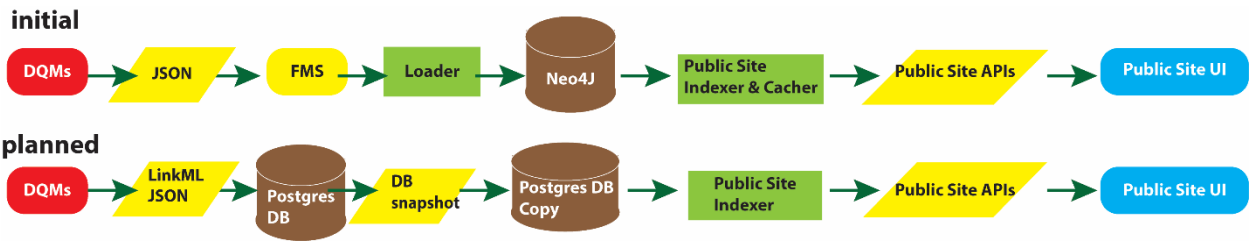


Figure 8
165x31 mm (x DPI)

1
2
3
4

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-WBStrain0001	FGP29	strain	Caenorhabdit...	WB		
WB-WBStrain0001	IE4314	strain	Caenorhabdit...	WB		
MGI:5008182	Akr1a1 ^{OST} ...	genotype	Mus musculu...	MGI		
MGI:6492714	Atg7 ^{em1} (MPC...	genotype	Mus musculu...	MGI		
ZFIN:ZDB-FISH-11	crb2a ^{#289/m...}	fish	Danio rerio (...)	ZFIN		

Figure 9
93x51 mm (x DPI)

1
2
3
4

1
2
3
4

SEARCH SCOPE:
DOCUMENT

SEARCH LOCATION:
DOCUMENT

Available Literature Info
Current site contains 1 literature:
S. cerevisiae (89212 papers)

Pick Category from Tree

Include children of selected categories

Type in category

or pick from tree below:

- root
- ChEBI (Tp:0000100)
- Gene Ontology (Tp:0000103)
- SGD Curation (Tp:0000017)
- Sequence Ontology (Tp:0000101)
- disease (DOID:4)

SELECT LITERATURE Current selection: S. cerevisiae

Figure 10
61x34 mm (x DPI)

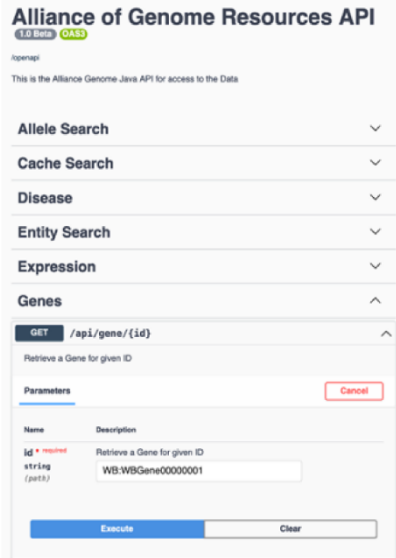


Figure 11
54x75 mm (x DPI)

1
2
3
4

The screenshot shows the Alliance of Genome Resources website. The top navigation bar includes the logo, 'Sign Up', and 'Log In' buttons. The main content is organized into three columns: 'Category', 'Topics', and 'Latest'. The 'Category' column lists sections like 'Alliance of Genome Resources', 'Job Postings', 'Positions Wanted', 'Meeting Announcements', 'Model Organism: Flies', and 'Model Organism: Frogs'. The 'Topics' column shows the number of items in each category. The 'Latest' column displays a list of recent posts with their titles, icons, and dates.

Category	Topics	Latest
Alliance of Genome Resources ■ News and Announcements ■ Site Feedback ■ Data Discussion ■ General Discussion	29	Welcome to Discourse Worms Nov '20
Job Postings Open positions and job announcements. ■ Flies ■ Frogs ■ Mice ■ Rats ■ Worms ■ Yeast ■ Zebrafish ■ Other	1.1k	MMRRC Newly Available Strains July 2023 & MMRRC Newly Accepted Strains July 2023 Stocks
Positions Wanted Are you a graduate student, postdoc, or young faculty member looking for a position? Post your details and requirements here. ■ Flies ■ Frogs ■ Mice ■ Rats ■ Worms ■ Yeast ■ Zebrafish	11	Drug-induced shrinkage of nematodes Scientific Discussion
Meeting Announcements Announcements and discussions about upcoming meetings ■ Flies ■ Frogs ■ Mammals/Human ■ Worms ■ Yeast ■ Zebrafish	132	How to enter data in Kaplan Meier graph? Methods & Reagents
Model Organism: Flies Discussion related to <i>Drosophila melanogaster</i> . ■ Reagents ■ FlyBase	8	Multi-purpose embryo extracts-Freon Free protocol Methods & Reagents
Model Organism: Frogs ■ News and Announcements ■ Scientific Discussion ■ Stocks	4	Xenopus Developmental Biology 1-week course Sept 11-15, 2023 Frogs
		Project Manager, Rare Disease Translational Center at JAX Job Postings

Figure 12
94x72 mm (x DPI)

1
2
3
4

Summary ribbon

Associated Reagents

Associated Phenotypes

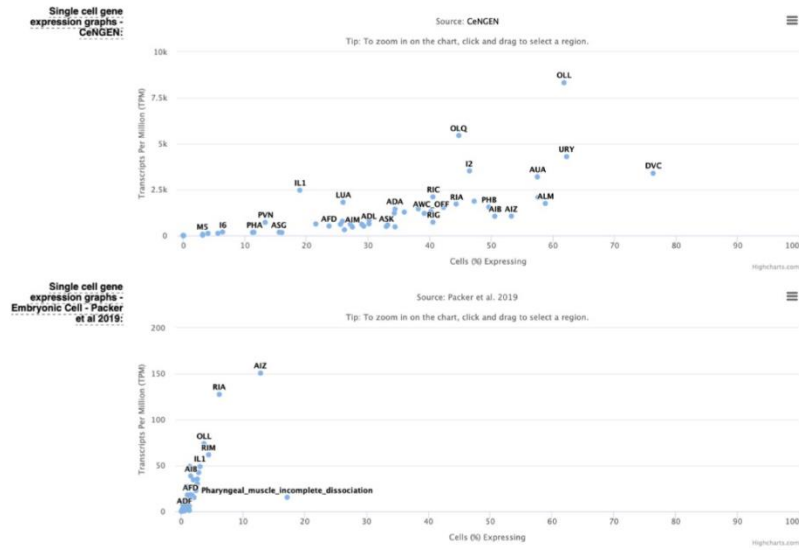
Associated Images

Associated Movies

sc-RNA Seq graphs

Other sources

EXPRESSION
 eat-4 expression



- 1
- 2
- 3
- 4

Figure 13
 142x91 mm (x DPI)

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 JAN;4(2):1421-1471. PMID: 35292937

Ihara et al., Treatment of Alzheimer's disease with framework nucleic acids. *Cell Prolif* 2020 Apr;53(4):1770. PMID: 32162173

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

Hogg et al., Use of Zebrafish Genetic Models to Study Pathology of the Amyloid Beta and Neurofibrillary Tangle Pathways in Alzheimer's Disease. *Cell Neuropharmacol* 2022 Mar 4;20(3):524-539. doi: PMID: 34335617

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	Species	Model	Allele	Publication
APP	U08006	Homo sapiens	Rat	APP	...
PTEN	U08006	Homo sapiens	Rat	PTEN	...

Alleles

Allele	Accession	Species	Model	Gene	Publication
APP	U08006	Homo sapiens	Rat	APP	...
PTEN	U08006	Homo sapiens	Rat	PTEN	...

Models

Model	Accession	Species	Gene	Allele	Publication
APP	U08006	Homo sapiens	APP	APP	...
PTEN	U08006	Homo sapiens	PTEN	PTEN	...

Publications

Pub	Author	Title	Year	PMID	Category
Yang et al.	Yang et al.	Long term exercise pre-training attenuates Alzheimer's disease-related pathology in a transgenic rat model of Alzheimer's disease.	2022	35292937	Genes
Ihara et al.	Ihara et al.	Treatment of Alzheimer's disease with framework nucleic acids.	2020	32162173	Genes

1
2
3

Figure 14
124x70 mm (x DPI)

Downloaded from https://academic.oup.com/genetics/advance-article/doi/10.1093/genetics/iyae049/7637331 by guest on 22 April 2024

39