

VEuPathDB User Documentation

Contract Name: VEuPathDB

Contract #: 75N93019C00077

Contractor: University of Pennsylvania

Reporting Period: September 2020

Prepared By: Jessica Kissinger and Omar Harb

Reviewed By: Brian Brunk, Jeremy DeBarry, other staff
as appropriate

Submitted By: Jeremy DeBarry

Table of Contents

INTRODUCTION	4
ABOUT VEUPATHDB.ORG	4
CURRENT FUNDING.....	5
DATA ACCESS POLICY	5
CITING VEUPATHDB IN PUBLICATIONS AND PRESENTATIONS	5
COMMUNITY INTERACTIONS AND DATA SUBMISSION POLICIES	6
HOW TO SUBMIT DATA TO VEUPATHDB	7
How to submit data for integration in VEuPathDB	8
<i>Genome Sequence and/or Annotation</i>	8
<i>High Throughput or Next Generation Sequencing</i>	8
<i>Microarray</i>	9
<i>Proteomics</i>	9
<i>Quantitative Proteomics</i>	9
<i>ChIP-chip</i>	10
<i>Isolates typed by sequencing limited genetic loci</i>	10
<i>Isolates or Strains typed by High Throughput Sequencing</i>	10
<i>General Data Submission</i>	11
Data submission and release on VEuPathDB databases.....	11
<i>General principles:</i>	11
<i>Why submit my data to VEuPathDB?</i>	11
<i>How do I submit data to VEuPathDB?</i>	12
<i>What data types are supported by VEuPathDB?</i>	12
THE DATA PRODUCTION CYCLE	13
Data Management SOPs (Standard Operating Procedures) for VEuPathDB databases.....	14
HOW TO USE OUR SITES	15
Online instructional material.....	15
Print-based instructional material	16
<i>Site Search</i>	16
<i>Search Strategies</i>	22
<i>Exploring the Gene Page</i>	28
<i>JBrowse Basics</i>	33
<i>Gene Ontology (GO) Enrichment</i>	46
<i>Advanced Search Strategies</i>	52
<i>Public Strategies</i>	61
<i>Regular Expressions & Genomic Colocation</i>	62
<i>Variant calling in VEuPathDB galaxy (Part 1)</i>	71
RELATED SITES OF INTEREST TO OUR COMMUNITIES.....	76
VEUPATHDB PUBLICATIONS AND CITATIONS.....	76

PUBLICATIONS THAT USE OUR RESOURCE	77
RELEASE NOTES.....	78
EXAMPLE RELEASE NOTES - VECTORBASE 48 RELEASED.....	78
ANALYSES METHODS:.....	81
GENOME ANALYSES.....	81
SUPPLEMENTS TO THE EBI PIPELINES.....	86
IN-HOUSE GENOME ANALYSES IN LIEU OF THE EBI PIPELINE	87
PROTEOMICS	89
RNA-SEQUENCE.....	89
CHIP-SEQUENCE	89
COPY NUMBER VARIATION.....	89
GENETIC VARIATION AND SNP CALLING	89
PROTEIN ARRAY DATA.....	90
METABOLIC PATHWAYS	90
DATASET DESCRIPTIONS:.....	90
TECHNICAL INFRASTRUCTURE AND SOFTWARE DOCUMENTATION:	90
BROWSER COMPATIBILITY STATEMENT	91
DATA LOADING AND DATABASE SCHEMA	91
CODE AVAILABILITY	92
WEB PRESENTATION SYSTEM AND USER INTERFACES.....	92
SOFTWARE CODE REPOSITORY	92
SYSTEM HARDWARE AND THIRD-PARTY SOFTWARE	92
OVERVIEW OF THE VEUPATHDB DATA PRODUCTION WORKFLOW AND ARCHITECTURE.....	93
VEUPATHDB WEBSITE PRIVACY POLICY.....	95
INTRODUCTION.....	95
<i>Information Automatically Collected</i>	95
<i>Information You Directly Provide</i>	96
<i>“Contact Us” Form.</i>	96
<i>How VEuPathDB Uses Cookies</i>	97
<i>Google Analytics</i>	97
<i>Third-Party Websites and Applications</i>	97
<i>Your Rights based on the General Data Protection Regulation (GDPR)</i>	99
VEUPATHDB ACCESSIBILITY CONFORMANCE	100
VEUPATHDB PERSONNEL.....	100
VEUPATHDB MANAGEMENT.....	100
CURRENT VEUPATHDB TEAM MEMBERS	100
VEUPATHDB ACKNOWLEDGEMENTS	102
<i>VEuPathDB Community Representatives</i>	102
<i>Previous Scientific Working Group</i>	103
WEBSITE USAGE STATISTICS	104

<i>Website usage links:</i>	104
<i>Sample awstats report from FungiDB.org</i>	104
VEUPATHDB GLOSSARY	107

Introduction

This material is made available to our users and NIH administrators to help them use VEuPathDB resources and understand the underlying tools, experiments and analyses provided by this Bioinformatics Resource Center funded in part by the US National Institute of Allergy and Infectious Diseases (Contract HHSN75N93019C00077). Note that the content of this document is also provided through the website, often in a context dependent manner. The web links are provided in the appropriate places throughout.

This report summarizes the categories of user documentation and provides a link to the landing page for each. The documentation can be divided into four broad categories: Site usage; Analysis methods; Dataset descriptors and Technical infrastructure and software documentation.

This report is available from all VEuPathDB sites, e.g. <https://veupathdb.org/>, from the Help menu.

About VEuPathDB.org

The Eukaryotic Pathogen, Vector and Host Informatics Resource (VEuPathDB) is one of two [Bioinformatics Resource Centers \(BRCs\)](#) funded by the US National Institute of Allergy and Infectious Diseases (NIAID), with additional support from the Wellcome Trust (UK).

These resources stem from support initially provided for the Plasmodium Genome Database by the Burroughs Wellcome Fund (2000-2) and a research grant from NIAID (2002-6). The BRC program was initiated in 2004 to provide public access to computational platforms and analysis tools enabling collection, management, integration and mining of genomic information and other large-scale datasets relevant to infectious disease pathogens including their interaction with mammalian hosts and invertebrate vectors of disease. Two BRCs are currently funded:

- VEuPathDB focuses on eukaryotic pathogens and invertebrate vectors of infectious diseases, encompassing data from prior BRCs devoted to parasitic species (EuPathDB), fungi (FungiDB) and vector species (VectorBase).
- BV-BRC - [PATRIC](#) & [VIRP](#) focus on bacterial and viral pathogens.

VEuPathDB has also received funding from the Wellcome Trust (UK), Bill & Melinda Gates Foundation, and US Department of Agriculture to support informatics efforts focusing on kinetoplastida and fungal organisms with special emphasis on improving functional annotation for select genome sequences and families of genes.

VEuPathDB provides access to diverse genomic and other large scale datasets related to eukaryotic pathogens and invertebrate vectors of disease (see [Data Summary](#)). Organisms supported by this resource include (but are not limited to) the NIAID list of [emerging and re-emerging infectious diseases](#).

Component web sites are constructed using a common infrastructure and standard data analysis and loading procedures, allowing the use of [VEuPathDB](#) as a single point of entry for each (or all) of these community resources, and the opportunity to leverage orthology for searches across taxa.

[Gene Metrics Table](#) summarizes the number of annotated genes for the organisms currently supported by this website.

- [Genome Summary](#) provides a list of all organisms available in this website.
- [Data Sets](#) provides a list of all information in this website integrated into VEuPathDB, with relevant references.

Current Funding

The Eukaryotic Pathogen, Vector and Host Informatics Resources (VEuPathDB) is funded by the National Institute of Allergy and Infectious Diseases (NIH/DHHS) under Contract No. NIH HHS 75N93019C00077.

VEuPathDB also receives funding from the Wellcome Trust (UK) to support informatics efforts focusing on kinetoplastida and fungal organisms with special emphasis on improving functional annotation of genomes. Grant numbers: 212929/Z/18/Z and 218288/Z/19/Z.

Data Access Policy

All data on these websites are provided freely for public use, through the contributions of many researchers involved in generating genome sequences, functional genomics datasets, and additional information. These data often derive from ongoing research and are not guaranteed to be accurate. When using data obtained from VEuPathDB, it is important to cite the original publications and contributors. Please see [Citing VEuPathDB](#).

Citing VEuPathDB in Publications and Presentations

If you use a VEuPathDB resource, we invite you to please cite the most relevant publication. This [PubMed filter](#) provides a list of the most recent VEuPathDB publications.

Please note that much of the data in VEuPathDB is provided by independent researchers, please cite them if you use their data. See [Data Sets](#) for a list of all information integrated into VEuPathDB, and related publications.

For acknowledgements in presentations, you may wish to use one or more of the following logos (right/control click to copy):



Additional resources leveraging the same infrastructure: [MicrobiomeDB](#) and [ClinEpiDB](#).

Community Interactions and Data Submission Policies

VEuPathDB serves a global scientific community that demands direct active support and community involvement. VEuPathDB outreach activities include:

- Organizing and running hands on training workshops and webinars ([Google Map](#)).
- Developing educational material in the form of [exercises](#) and [online tutorials](#).
- Responding to support emails for users who contact us directly by clicking the "Contact Us" links in the header or footer of any VEuPathDB webpage (average response time is 48 hours).
- Holding open community meetings/forums with our diverse user base. These meetings are held in person at scientific conferences or using an online conferencing platform.
- Attending national and international meetings with active participation in the form of posters, presentations or help desks.
- Authoring [peer reviewed manuscripts](#).
- Maintaining active social media presence in the form of a [FaceBook page](#) and [Twitter feed](#).
- Providing a clear [data handling and release policy](#) to investigators to encourage submission of data prepublication.

How to Submit Data to VEuPathDB

The Eukaryotic Pathogen, Vector & Host Informatics Resources (<https://VEuPathDB.org>) is a Bioinformatics Resource Center (BRC) operated under contract from the US National Institute of Allergy and Infectious Diseases (NIAID) and the Wellcome Trust. VEuPathDB is charged with ensuring that genomic (and other large-scale) datasets pertaining to supported pathogens are conveniently accessible to the worldwide community of biomedical researchers. This document summarizes policies associated with releasing datasets on VEuPathDB. Our goal is to help the communities that we serve ensure that their data are FAIR, Findable, Accessible, Interoperable and Reusable.

VEuPathDB welcomes submissions of genomic-scale data concerning eukaryotic microbes, fungi, vectors of human disease, and host-pathogen interactions. The VEuPathDB contract from NIAID provides support for biosecurity pathogens, including *Babesia*, *Cryptosporidium*, *Entamoeba*, *Giardia*, *Microsporidia* (various genera), *Toxoplasma*, *Plasmodium*, and related taxa (*Acanthamoeba*, *Gregarina*, *Neospora*, *Theileria*) and also arthropod vectors (ticks, mites, mosquitoes, kissing bugs, tsetse flies, sand flies, lice, etc.) of human disease, as well as a snail that serves as an intermediate host, and comparator species. Support for kinetoplastid parasites (*Crithidia*, *Endotrypanum*, *Leishmania*, *Trypanosoma*) is provided by The Wellcome Trust. The FungiDB project encompasses a large (and growing) number of species supported by both NIAID and the Wellcome Trust. Please [contact us](#) if you have data from other species that should be incorporated into VEuPathDB! Please review our [Data Submission Policy](#).

Our most common data types include transcriptomics, proteomics, metabolomics, epigenomics, population-level and isolate information. In one form or another, VEuPathDB currently represents datasets in the following categories:

- Sequence (genomic [nuclear and organellar])
- ESTs and RNA-seq, generated on various platforms)
- Host-response data
- Comparative genomic information
- DNA polymorphism and population genetics data
- Sequences and metadata pertaining to field and clinical isolates and collections (with geo-spatiotemporal and other metadata)
- Chromatin modification data (ChIP-chip and ChIP-seq)
- Manually curated and automatically generated gene models and other annotation (GO terms, InterPro domains, etc.)
- Transcript and proteomic profiling
- Host response data sets (multiple platforms)
- Interactome data
- Protein structural information
- Metabolic pathways and metabolomics data
- Phenotype information, reagents (clones, antibodies, etc.)
- Publication references
- Image data, etc.

We also accept other genomic-scale data and are open to suggestions. Use the [Contact Us](#) link to make suggestions. We look forward to working with you!

How to submit data for integration in VEuPathDB

To submit your data for integration, fill out the appropriate VEuPathDB Dataset nomination form listed below. If your data cannot be submitted via our forms, use the [Contact Us](#) link to send a brief description (two or three sentences) of your data.

Genomes & high throughput sequencing data (e.g. RNA-Seq, ChIP-Seq, isolates typed by Whole Genome Sequencing or by sequencing limited genetic loci) must be available in The International Nucleotide Sequence Database Collaboration (INSDC) such as NCBI GenBank, EMBL-EBI ENA or DDBJ.

Once the dataset is prioritized for loading, we will export the data directly from INSDC. Note: while genome sequences must be available in INSDC, functional annotation (e.g. gene names, GO terms, etc.) can be submitted directly to VEuPathDB.

Tell us about your data as early as possible, to allow ample time for scheduling into VEuPathDB release cycles. Depending on the dataset type, we can provide instructions on how to transfer your data to us (e.g. formats of proteomics datasets may differ depending on the nature and scale of the data to be transferred), or we may be able to facilitate data submission to a repository (e.g. GenBank, GEO/ArrayExpress, etc.).

VectorBase resource: Gene manual annotations (change of exon-intron boundaries, creation of new genes) and metadata (gene names/symbols and functional description) can be submitted via Apollo (Coming soon). If you submit a gene annotation before you submit a manuscript for publication, we can generate GeneIDs that can be linked out to within the publication. Gene deletions are not handled via Apollo, please [Contact Us](#) with supporting evidence. Metadata can also be submitted by sending a spreadsheet file for batch submissions, follow this link for information about the 12 columns heading that are required. To add publications to a gene send us the corresponding PubMed link.

Genome Sequence and/or Annotation

Genomes must be available in The International Nucleotide Sequence Database Collaboration (INSDC) such as NCBI GenBank, EMBL-EBI ENA or DDBJ.

- If your genome **IS** uploaded to a repository, complete the [Genome Sequence and/or Annotation Description Form](#) making sure to include the accession numbers of your data when prompted. We will download your data from the repository.
- If the annotation file **is not** uploaded to a repository, use the [Contact Us](#) form to send us the genome annotation file only (e.g. gff file format).

High Throughput or Next Generation Sequencing – RNA, DNA or ChIP Sequencing

We prefer to download the raw read data in FASTQ format from a sequence read archive. We integrate the data into the database using the raw reads and use the raw reads during future database releases to remap or update our analyses when necessary.

How to transfer a copy of your data to VEuPathDB:

- Upload your data to a sequence read archive such as the European Nucleotide Archive or NCBI's Sequence Read Archive and provide us the accession numbers. We will export the data directly from INSDC
- Complete the appropriate data description form making sure to enter your data archive accession numbers when prompted

[RNA-Seq Data Description Form](#)

[DNA-Seq Data Description Form](#)

[ChIP-Sequencing Data Description Form](#)

Microarray

Transfer a copy of your data to VEuPathDB using one of these options:

- Upload your data to a repository such as GeneExpression Omnibus. Complete the data description form linked below making sure to enter your data archive accession numbers when prompted. We will export the data directly from a repository.
- Send your data as an attachment to an email. Use the [Contact Us](#) form to send us an email. Files (e.g. CEL, CSV) must include expression levels and probe set information. Pay special attention to clearly indicate the identity of columns in the data files you transferred to VEuPathDB.

[Microarray Data Description Form](#)

Proteomics

Excel or tab delimited text files are preferred. We can accommodate xml file format. Required columns include gene IDs, peptide sequences, peptide counts and scores.

How to transfer a copy of your data to VEuPathDB:

- Send your data as an attachment to an email. Use the [Contact Us](#) form to send us an email.
- Complete the [Proteomics Data Description Form](#) making sure to clearly indicate the content of each column in your file.

Quantitative Proteomics

Excel or tab delimited files are preferred. We can accommodate xml file format. Required columns include gene IDs and scores.

How to transfer a copy of your data to VEuPathDB:

- Send your data as an attachment to an email. Use the [Contact Us](#) form to send us an email.
- Complete the [Quantitative Proteomics Data Description Form](#) making sure to include a description of data columns, for example, time course units and arrangement if not apparent from column headers.

ChIP-chip

Your data files should include expression levels and probe set information.

Transfer a copy of your data to VEuPathDB using one of these options:

- Upload your data to a repository such as Gene Expression Omnibus and complete the [ChIP-chip Data Description Form](#) making sure to enter the archive accession numbers (if any) for your data when prompted.
- Send your data as an attachment to an email. Use the [Contact Us](#) form to send us an email.

Isolates typed by sequencing limited genetic loci

If your data **IS** uploaded to Genbank, use the [Contact Us](#) form to tell us about your data. Note: Genebank Isolate records and the associated metadata are automatically updated with each VEuPathDB release. There is no need to complete our Isolate Submission Form.

If your data **IS NOT** uploaded to Genbank, we can facilitate this upload. Complete the Isolate Submission Form linked below and we will use the information to generate a Genbank submission for your isolates. The new isolate records will be downloaded to VEuPathDB with the release. Use the [Contact Us](#) form to send us instructions for retrieving your data.

[Isolate Submission Form](#)

[Help for submitting Isolate Data](#)

Isolates or Strains typed by High Throughput Sequencing

We prefer to receive the raw read data in FASTQ or FASTA file format. We integrate your data into the database using the raw reads. We also use the raw reads during future database releases to remap your data when the reference genome is reloaded and to update our analyses when needed.

How to transfer a copy of your data to VEuPathDB:

- Upload your data to a sequence read archive such as DNA Data Bank of Japan, the European Nucleotide Archive or NCBI's Sequence Read Archive. We will retrieve your data using the read archive's accession numbers for your data set.
- Complete our [DNA Seq Data Description Form](#) making sure to enter the read archive accession numbers for your data when prompted. We also ask that you complete an abbreviated [Abbreviated Isolate Submission Form](#) to describe meta data associated with your isolates.

General Data Submission – Use the [Contact Us](#) form to tell us about data that does not fit any of the above categories

Data submission and release on VEuPathDB databases

issued February 2010, most recent revision 02 April 2020



The Eukaryotic Pathogen, Vector & Host Informatics Resources (<http://VEuPathDB.org>) is a Bioinformatics Resource Center (BRC) operated under contract from the US National Institute of Allergy and Infectious Diseases (NIAID) and the Wellcome Trust. VEuPathDB is charged with ensuring that genomic (and other large-scale) datasets pertaining to supported pathogens are conveniently accessible to the worldwide community of biomedical researchers. This document summarizes policies associated with releasing datasets on VEuPathDB and affiliated knowledgebases. Our goal is to help the communities that we serve ensure that their data are FAIR, Findable, Accessible, Interoperable and Reusable.

General principles:

- ***Data providers define the schedule for data release (in consultation with funders, publishers, etc).*** While there is no point in providing VEuPathDB with data that will never become public, deposition does not in itself authorize immediate release. Data become accessible to the public only when the data providers and VEuPathDB staff agree that it is accurately represented and ready to go live. Note that knowledgebase staff are not active research scientists; they are distinct from researchers in the groups responsible for VEuPathDB, who see new data only when it becomes accessible to the general public.
- ***Data providers know their data best.*** We expect to work with those who generated the underlying data to determine how best to analyze and represent new data types. This typically means taking in relatively raw data – often earlier, and in a more unprocessed form than the published dataset – and building an in-house analysis pipeline to ensure that all comparable datasets are handled similarly.
- ***The earlier we learn about new datasets, the easier it is to schedule timely release.*** The nature of our knowledgebase production, and competing demands from the many communities we support, means that several months' notice are often required to prepare for release. Note that it is often possible to use a preliminary dataset for planning, which can be swapped for the final version before public release.
- ***Experience has shown that data not deposited prior to publication often fails to emerge at all!*** After publication, it may be difficult to focus on tracking down the raw data, associated metadata, analysis methods, etc. It is never too early to discuss planned datasets with the VEuPathDB team!
- ***While not required, pre-publication data release often results in favorable attention from scientific colleagues (including journal editors and grant reviewers).*** Note that all major scientific journals now agree that early release of genomic-scale datasets does not compromise publication.

Why submit my data to VEuPathDB?

- Inclusion in VEuPathDB facilitates your own analysis of the data, in the context of

other genomic-scale experiments already available from researchers around the world.

- Electronic access permits others to analyze your data in greater depth than possible in print (even in advance of publication, if you wish to allow this).
- Availability within VEuPathDB keeps your data alive on a highly visible genomics knowledgebase resource: VEuPathDB is accessed by ~13,000 unique users each month.

How do I submit data to VEuPathDB?

- Fill out the appropriate form to indicate the data availability
- Contact the VEuPathDB by clicking the 'Contact Us' link on any VEuPathDB page or emailing us at help@VEuPathDB.org.
- Tell us about your data as early as possible, to allow ample time for scheduling into VEuPathDB release cycles.
- Once you tell us about your data, we will provide instructions on how to transfer your data to us (formats may differ depending on the nature and scale of the data to be transferred).
- In order to avoid any confusion and ensure accuracy, we adhere to strict Standard Operating Procedures (SOPs), as outlined below.

What data types are supported by VEuPathDB?

In one form or another, VEuPathDB currently represents sequence (genomic [nuclear and organellar], ESTs and RNA-seq, generated on various platforms), host-response data, comparative genomic information, DNA polymorphism and population genomics data, sequences and metadata pertaining to field and clinical isolates and collections (with geo-spatiotemporal and other metadata), chromatin modification data (ChIP-chip and ChIP-seq), manually curated and automatically generated gene models and other annotation (GO terms, InterPro domains, etc.), transcript and proteomic profiling, host response data sets (multiple platforms), interactome data, protein structural information, metabolic pathways and metabolomics data, phenotype information, reagents (clones, antibodies, etc.), publication references, image data, etc. We can support additional data types as needed.

Please let us know if you have data to provide that is not currently supported! What species are supported by VEuPathDB?

The VEuPathDB contract from NIAID provides support for biosecurity pathogens, including *Babesia*, *Cryptosporidium*, *Entamoeba*, *Giardia*, *Microsporidia* (various genera), *Toxoplasma*, *Plasmodium*, and related taxa (*Acanthamoeba*, *Gregarina*, *Neospora*, *Theileria*) and also arthropod vectors (ticks, mites, mosquitoes, kissing bugs, tsetse flies, sand flies, lice, etc.) of human disease, as well as a sail that serves as an intermediate host, and comparator species. Support for kinetoplastid parasites (*Crithidia*, *Endotrypanum*, *Leishmania*, *Trypanosoma*) is provided by The Wellcome Trust. The FungiDB project encompasses a large (and growing) number of species supported by both NIAID and the Wellcome Trust.

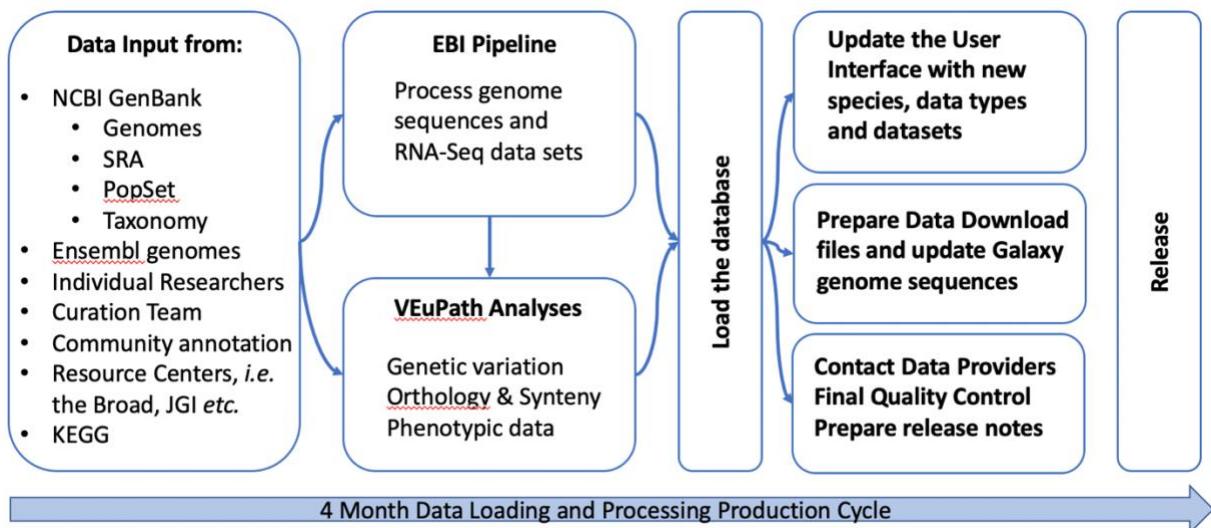
Please contact us if you have data from other species that should be incorporated into VEuPathDB!

The Data Production Cycle

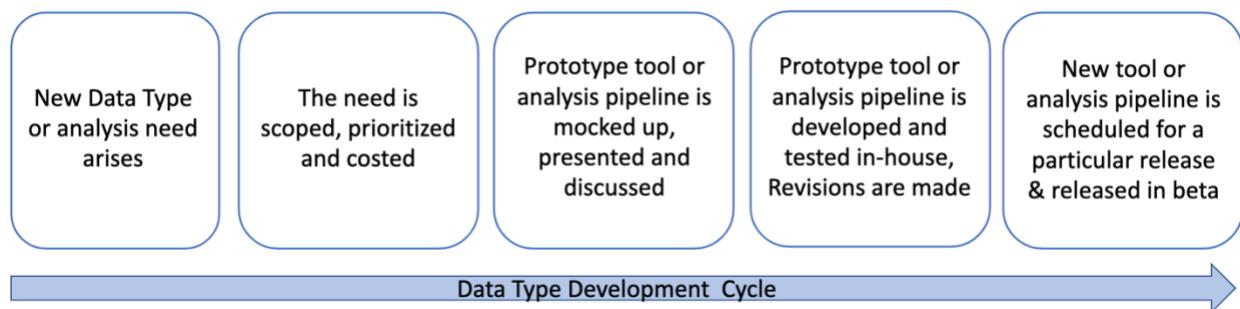
Guiding Principle:

- We balance the need to process as many datasets as possible from known data types in our production pipeline with the need to scale in capacity and containerize analyses for these same data types and the need to continually build new tools and infrastructure in our data type development to accept new data types and facilitate emerging community analysis needs as they emerge.

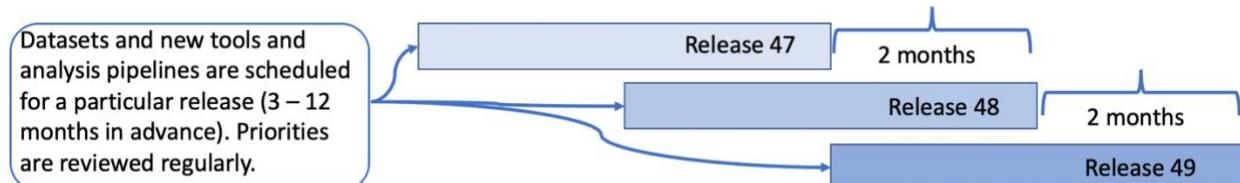
(A)



(B)



(C)



High level views of the VEuPathDB Production Cycle (A), Datatype Development Cycle (B), and Release Timeline (C; overlapping 4 month cycles of production and data type development with bimonthly releases).

Data Management SOPs (Standard Operating Procedures) for VEuPathDB databases

VEuPathDB routinely handles datasets provided prior to publication, in addition to those already in the public domain. In order to ensure timely and accurate data integration we strictly adhere to the following Standard Operating Procedures (SOPs):

1. Datasets come to our attention in several ways, including
 - Direct contact from researchers generating the data (during the earliest stages of project design, as
 - data is being produced, in the course of data analysis, or in the context of manuscript preparation).
 - Information provided by members of the VEuPathDB or larger pathogen community.
 - Information obtained by VEuPathDB staff at meetings and conferences.
 - Publicly available information from the scientific literature, genomic dataset repositories, etc.

* Note that VEuPathDB can often facilitate data deposition in the appropriate archival repositories (GenBank, SRA, GEO/ArrayExpress, etc.)

- 2. Decisions to include a dataset in VEuPathDB are based on value to the research community.** When prioritizing data for integration, we rely heavily on discussions with active researchers, including the scientific advisory committees established for each of the taxonomic groups supported by VEuPathDB and the objectives and priorities of our funders. **Please contact us if you are interested in participating in these discussions.**
- 3. Regardless of how we first learn about a given dataset, communication is established with the original data producer** through email, teleconference, and/or face-to-face meetings to discuss the desirability and feasibility of integration into VEuPathDB. In the course of these discussions, we consider what data are likely to be available, data formats and transfer protocols, questions the community may wish to ask of this data, and ways to represent or display such information. We also collect appropriate metadata (regarding samples, experimental protocols, etc.), and information on data sources, data providers, appropriate citation, etc.
- 4. Data provided to VEuPathDB is housed on secure servers and never shared outside of VEuPathDB staff without prior consent of the data provider.** Note that database staff are not active researchers; they are distinct from students and postdocs in the groups responsible for VEuPathDB, who see new datasets only when they become accessible to the general public.
- Datasets are assigned a provisional release date, in consultation with the data provider. **Scheduling a dataset does not mean that it will be released without final examination and approval by the data provider!** We operate on the assumption that those who generate the data are best placed to evaluate its proper integration and representation in the knowledgebase. Note that this 'golden rule' applies to both published and unpublished data.

6. Three to four months before the scheduled release date, the data provider is contacted by the **Outreach** team, to ensure that we have the most up-to-date version of the data, along with appropriate metadata and information on data sources and citations. The **Data Loading** team then processes and integrates this data into our internal knowledgebases.
7. After data loading is complete, the **Data Development** team begins to analyze and develop searches against the data. At this point we will likely communicate with the data provider, if questions arise.
8. **Once data development is underway, the data provider is given access to a password-protected version of the VEuPathDB website containing their data.** This development site is similar to the current production knowledgebase, except that it also includes new data from the provider. We also provide instruction on how to search and view these new data, including sample searches integrating new data with relevant information already available in the knowledgebase. Important questions to consider include:
 - Does the knowledgebase accurately represent your data?
 - Are the values and/or graphical displays provided appropriate?
 - Are the questions that one can ask of your data appropriate?
 - Are there additional questions that you would like to see implemented?
 - Are the data appropriately described, including relevant metadata and reference / citation details?
9. **A series of exchanges typically ensues**, in which we work iteratively with data providers to address any concerns, with changes reviewed on the password protected site so that providers can view and interrogate their data in the context of the rest of the database.
10. **Public release is only considered after everyone is satisfied with how the data is represented.** If the provider is not yet ready to authorize data public release, data is rescheduled for a future release, and removed from the development site before it goes live.
11. Once data are approved for public release, a description is included in the ‘News’ accompanying the next release, **highlighting new datasets and functionality, and acknowledging all data providers.**
12. **Post-release quality assurance** provides the opportunity to modify displays and develop new queries if/as appropriate.

How to Use Our Sites

We provide our users with a variety of mechanisms to learn about how the VEuPathDB site works and can facilitate their research. In addition to the instructions provided here, users should know that they can learn about VEuPathDB in person, via recorded instructional material or via help located throughout the website, detailed in the list below.

Online instructional material

- a. Users can request that a member of the VEuPathDB outreach staff attend a lab meeting or institutional or regional workshop to ask questions and have training in a specific topic
- b. User can participate in or view previously recorded webinars (online recordings of training sessions conducted by VEuPathDB staff) that are provided as YouTube recordings.
- c. Users can browse and download “quick start” materials and tutorials to follow at their own pace, each with color pictures and well explained prompts to help users explore the VEuPathDB site or one of its projects, e.g. VectorBase, FungiDB or PlasmoDB. All project use the same underlying architecture and software, only the local environment is customized so user can learn from a tutorial on any organism.
- d. VEuPathDB host several workshops annual and all training materials (step by step tutorials with exercises and answers) are provided for all to use and can be easily downloaded.
- e. Context sensitive help is available via the many help icons located on the web site, particularly on search pages.
- f. Textual descriptions are associated with search pages to provide information about the data type and how to use the search.

To quickly access these online resources, please visit:

<https://beta.veupathdb.org/veupathdb.beta/app/static-content/landing.html>

Print-based instructional material

Site Search

Note: *this exercise uses VectorBase as an example database, but the same functionality is available on all VEuPathDB resources.*

Learning objectives:

- Use keywords in site search
- Explore site search results
- Filter site search results by categories
- Filter site search results by organisms
- Filter site search results by category fields
- Export results to a search strategy
- Find a specific gene using its ID in site search

- Enter the word *kinase* in the site search window (top center of the page, arrow in the image below). Then click enter on your keyboard or click on the search icon (square in the image below).

- How many results with the word kinase did you get? Are all the results genes? Explore the filter panel on the left side of the webpage. Filter the results so that you

Filter results	
Genome	43,494
Genes	
Metabolism	
Metabolic pathways	293
Compounds	85
Data access	
Data sets	1
Searches	3

Filter fields	
Select a result filter above	

Filter organisms	
select all clear all expand all collapse all	
<input type="checkbox"/> Arthropoda	42,156
<input type="checkbox"/> Mollusca	1,338

only view gene results (hint: click on the word *genes* in the *Filter results* section; arrow in image below).

3. How many of the genes included the word kinase in their product descriptions? Notice that once you filter the result by genes (click on the *Genes* filter), the fields section expands to reveal additional filtering options. Once you select the *Product descriptions* field you are provided the option to *apply* this filter or cancel it (box middle panel below). Once a filter is applied it can be cleared by clicking on *Clear filter* (box left panel below).

Filter results

Hide zero counts

Genome Genes **Clear filter** 43,494

Filter Gene fields

[select all](#) | [clear all](#)

<input type="checkbox"/> EC descriptions and numbers	24,315
<input type="checkbox"/> GO terms	16,978
<input type="checkbox"/> Orthologs	25,638
<input type="checkbox"/> PDB chains	17,900
<input type="checkbox"/> Product descriptions	8,085
<input type="checkbox"/> PubMed	3

Filter organisms

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#)

<input type="checkbox"/> Arthropoda	42,156
<input type="checkbox"/> Mollusca	1,338

Filter results

Hide zero counts

Genome Genes **Clear filter** 43,494

Filter Gene fields **Apply** **X**

[select all](#) | [clear all](#)

<input type="checkbox"/> EC descriptions and numbers	24,315
<input type="checkbox"/> GO terms	16,978
<input type="checkbox"/> Orthologs	25,638
<input type="checkbox"/> PDB chains	17,900
<input checked="" type="checkbox"/> Product descriptions	8,085
<input type="checkbox"/> PubMed	3

Filter organisms

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#)

<input type="checkbox"/> Arthropoda	42,156
<input type="checkbox"/> Mollusca	1,338

Filter results

Hide zero counts

Genome Genes **Clear filter** 8,085

Filter Gene fields **Clear filter**

[select all](#) | [clear all](#)

<input type="checkbox"/> EC descriptions and numbers	24,315
<input type="checkbox"/> GO terms	16,978
<input type="checkbox"/> Orthologs	25,638
<input type="checkbox"/> PDB chains	17,900
<input checked="" type="checkbox"/> Product descriptions	8,085
<input type="checkbox"/> PubMed	3

Filter organisms

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#)

<input type="checkbox"/> Arthropoda	8,003
<input type="checkbox"/> Mollusca	82

4. How many of the above genes are found in *Anopheles gambiae* str. PEST? How did you find this number? (hint: explore the *Filter organisms* section of the results filter). Select the correct organism and apply the filter.
5. Export the results to a search strategy. (hint: to achieve this click the *Export as a search strategy* button at the top right-hand corner of the page).

Filter organisms **Apply** **X**

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#)

- Arthropoda 8,003
 - Arachnida 1,072
 - Insecta 6,931
 - Diptera 6,542
 - Culicidae 4,698
 - Aedes 386
 - Anopheles 3,923
 - Anopheles albimanus 173
 - Anopheles arabiensis 221
 - Anopheles atroparvus 203
 - Anopheles christyi 157
 - Anopheles coluzzii 192
 - Anopheles culicifacies 201
 - Anopheles darlingi 243
 - Anopheles dirus 184
 - Anopheles epiroticus 168
 - Anopheles farauti 212
 - Anopheles funestus 199
 - Anopheles gambiae 244
 - str. PEST 244
 - Anopheles maculatus 130
 - Anopheles melas 201
 - Anopheles merus 212
 - Anopheles minimus 177
 - Anopheles quadriannulatus 200
 - Anopheles sinensis 402
 - Anopheles stephensi 204
 - Culex 389

My Search Strategies

Opened (1) All (1) Public (3) Help

Unnamed Search Strategy *

Text 244 Genes Step 1

244 Genes (219 ortholog groups)

Gene Results		Genome View		Analyze Results	
Genes: 244 Transcripts: 310 <input type="checkbox"/> Show Only One Transcript Per Gene					
<input type="button" value="4"/> <input type="button" value="1"/> <input type="button" value="2"/> <input type="button" value="3"/> <input type="button" value="16"/> Rows per page: 20					
Gene ID	Transcript ID	Organism	Genomic Location (Gene)	Product Description	
AGAP004699	AGAP004699-RA	Anopheles gambiae str: PEST	AgamP4_2L:1,973,601..1,976,987(+)	RAF proto-oncogene serine/protein kinase [Source:VB C Annotations]	

6. Return to the site search results page. How did you do this? (hint: you can achieve this in two ways: 1. Click on your browser's back arrow. 2. Click on the back to results arrow in the site search window. Notice that your previous results and filter settings were preserved.

7. Clear all filters. How did you do this? (hint: you can achieve this in two ways: 1. You can click on each of the clear filter options in the filter results panel on the left (boxes below). 2. You can click on the single *clear filters* option in the site search window.

Filter results Hide zero counts

Genome Genes

Filter Gene fields

select all | clear all

<input type="checkbox"/> EC descriptions and numbers	410
<input type="checkbox"/> GO terms	436
<input type="checkbox"/> Orthologs	515
<input type="checkbox"/> PDB chains	396
<input checked="" type="checkbox"/> Product descriptions	244
<input type="checkbox"/> PubMed	1

Filter organisms

select all | clear all | expand all | collapse all

<input checked="" type="checkbox"/> Arthropoda	8,003
<input type="checkbox"/> Mollusca	82

8. Try the *Hide zero counts* check box in the *Filter results* panel. What does this do?

Category	Value (Left: Hide zero counts checked)	Value (Right: Hide zero counts unchecked)
Genome	Genes: 43,494	Genomic sequences: 43,494
Metabolism	Metabolic pathways: 293 Compounds: 85	
Data access	Data sets: 1 Searches: 3	
Filter fields	Select a result filter above	
Filter organisms	select all clear all expand all collapse all Arthropoda: 42,156 Mollusca: 1,338	
Instructional	Tutorials Workshop exercises	
About	News General info pages	

9. Try running a search with a wild card. The wild card is denoted by an asterisk *. The wild card can be used alone to retrieve all results available to the site search or combined with a word such as *kinase to retrieve compound words ending with the word kinase like phosphofructokinase. As usual results can then be explored using the filters in the *Results filter* on the left side of the website.

Result Type	Count	Description
Compound	4,457,608	Vismine D, Vismadin, Vismagin, ribostamycin sulfate salt, nalidixic acid, Voacamine, vobasine, vobtusine, volemitol, vomeritin, vomitoxin, voriconazole

All results matching *

Export as a Search Strategy
to download or data mine ►

1 - 20 of 4,457,608

Compound - CHEBI:10000 Vismine D

Compound - CHEBI:10001 Vismadin

Compound - CHEBI:10002 Vismagin

Compound - CHEBI:10003 ribostamycin sulfate salt

Definition: An aminoglycoside sulfate salt resulting from the reaction of ribostamycin with sulfuric acid.

Compound - CHEBI:10014 nalidixic acid

Definition: A monocarboxylic acid comprising 1,8-naphthyridin-4-one substituted by carboxylic acid, ethyl and methyl groups at positions 3, 1, and 7, respectively.

Compound - CHEBI:10014 Voacamine

Compound - CHEBI:10015 vobasine

Definition: An indole alkaloid that is vobasanine in which the bridgehead methyl group is substituted by a methoxycarbonyl group and an additional oxo substituent is present in the 3-position.

Compound - CHEBI:10016 vobtusine

Compound - CHEBI:10017 volemitol

Definition: A heptitol that is heptane-1,2,3,4,5,6,7-heptol that has R-configuration at positions 2, 3, 5 and 6.

Compound - CHEBI:10018 vomeritin

Definition: A cyanogenic glycoside that is (4R)-4-hydroxycyclopent-2-ene-1-carbonitrile attached to a beta-D-glucopyranosyloxy at position 1.

Compound - CHEBI:10019 vomicine

Compound - CHEBI:10022 Vomitoxin

Compound - CHEBI:10023 voriconazole

The screenshot shows the VectorBase search results for the query '*kinase'. The results are filtered by organism to Aedes aegypti LVP_AGWG. There are 45,121 results in total, with the first 20 shown. Each result includes the gene ID, name, source, and organism information. A sidebar on the left provides filtering options for genome, metabolic pathways, data access, filter fields, and filter organisms.

Category	Value
Genome	Genes: 44,659
Metabolic pathways	Compounds: 367
Data access	Data sets: 1
	Searches: 3
Filter fields	Select a result filter above
Filter organisms	Arthropoda: 43,291 Mollusca: 1,368

10. Try searching for a specific gene ID. Enter the gene ID below in the site search window:

AAEL007018

The screenshot shows the search results for the gene ID AAEL007018. The results are filtered by organism to Aedes aegypti LVP_AGWG. There is 1 result shown, which is AAEL007018, udp-glucose 4-epimerase. A sidebar on the left provides filtering options for genome, filter gene fields, and filter organisms.

Category	Value
Genome	Genes: 1
Filter Gene fields	Gene ID: 1 Transcripts: 1
Filter organisms	Arthropoda: 1 Insecta: 1

Notice that the gene of interest appears at the top for easy access. You can click on the Gene ID to go the gene page.

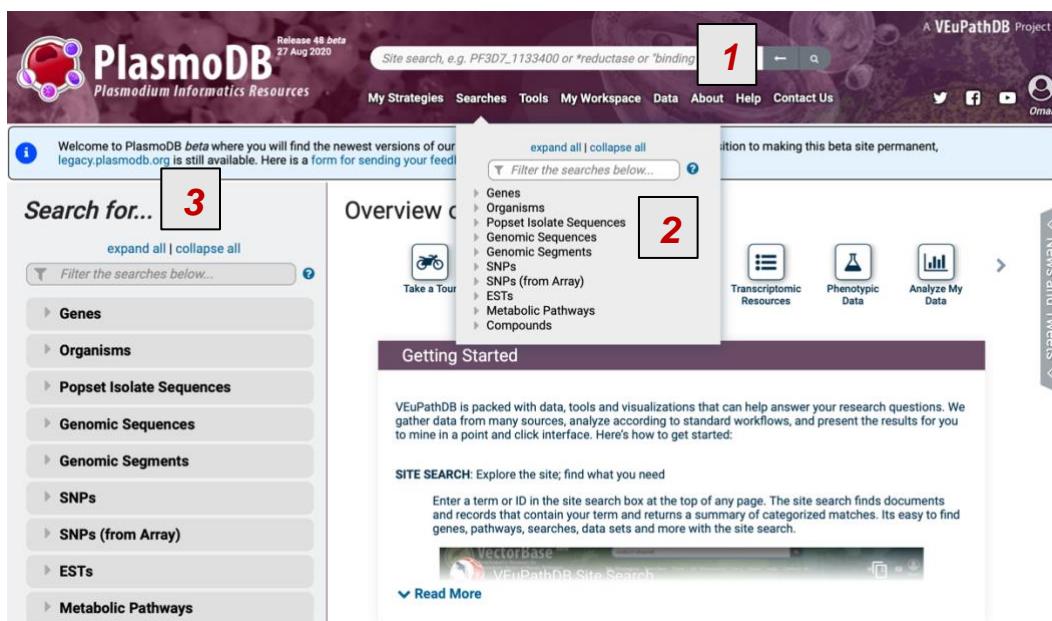
Search Strategies

Note: this exercise uses *PlasmoDB* (<https://PlasmoDB.org>) as an example database, but the same functionality is available on all VEuPathDB resources.

Learning objectives:

- Running a search to start a search strategy
- Adding steps in a search strategy
- Adding and sorting results
- Revising steps

There are three options to start a Search Strategy. 1) From the “Site Search” box ---> Export as Search Strategy, 2) In the site header from the “Searches” menu and 3) In the home page (left hand side) from the “Search for ...” section.



1. Go to the home page and in the “Search for ...” section on the left, filter the searches by typing the word transmembrane to find the **Transmembrane Domain Count search** in the filtered results.
2. Click on the transmembrane (TM) domain count search to get to the search page.

Search for...

expand all | collapse all

Filter the searches below...

- ▶ Genes
- ▶ Organisms
- ▶ Popset Isolate Sequences
- ▶ Genomic Sequences
- ▶ Genomic Segments
- ▶ SNPs
- ▶ SNPs (from Array)
- ▶ ESTs
- ▶ Metabolic Pathways
- ▶ Compounds

Search for...

transm

Genes

Protein targeting and localization

Transmembrane Domain Count

Configure this search to find all genes from *Plasmodium vivax* P01 that have at least 6 TM domains and at most 8 TM domains. See image below if you need help with the configuration.

Identify Genes based on Transmembrane Domain Count

1 SELECTED, OUT OF 45

[add these](#) | [clear these](#) | [select only these](#)
[select all](#) | [clear all](#)

vivax [?](#)

Plasmodium vivax
 Plasmodium vivax P01
 Plasmodium vivax Sal-1
 Plasmodium vivax-like sp.
 Plasmodium vivax-like PvI01

[add these](#) | [clear these](#) | [select only these](#)
[select all](#) | [clear all](#)

Minimum Number of Transmembrane Domains

[?](#)

Maximum Number of Transmembrane Domains

[?](#)

[Get Answer](#)

3. How many genes did you obtain? (hint: look at the number results in the strategy step in yellow, or the number right above the results and below the search strategy).

My Search Strategies

The screenshot shows the VEuPathDB search interface. At the top, it says "Opened (1) All (404) Public (42) Help". Below that is the "Unnamed Search Strategy *". A red arrow points to the "Transmb Dom 101 Genes" button. Another red arrow points to the "101 Genes (97 ortholog groups)" link. The main area shows a table of gene results with columns: Gene ID, Transcript ID, Organism, Genomic Location (Transcript), and # TM Domains. The table lists 101 entries, each corresponding to a Plasmodium vivax P01 gene.

Gene ID	Transcript ID	Organism	Genomic Location (Transcript)	# TM Domains
PVP01_0104300	PVP01_0104300.1	Plasmodium vivax P01	PvP01_01_v1:213970..224298(+)	6
PVP01_0113600	PVP01_0113600.1	Plasmodium vivax P01	PvP01_01_v1:607892..610837(+)	6
PVP01_0114900	PVP01_0114900.1	Plasmodium vivax P01	PvP01_01_v1:664719..665660(-)	6
PVP01_0317200	PVP01_0317200.1	Plasmodium vivax P01	PvP01_03_v1:744775..747000(+)	6
PVP01_0606500	PVP01_0606500.1	Plasmodium vivax P01	PvP01_06_v1:261479..262988(-)	6
PVP01_0703300	PVP01_0703300.1	Plasmodium vivax P01	PvP01_07_v1:187292..193820(-)	6
PVP01_0706600	PVP01_0706600.1	Plasmodium vivax P01	PvP01_07_v1:352530..354459(-)	6

4. Explore the results table. Try the following things:

- Sort the #TM domain column to show genes with 8 TM domains first.
- Add a column for transcript length (Click on add columns and find the transcript length column, then click on update columns).

The screenshot shows the "Select Columns" dialog box and the main results table. The dialog box has a red box around the "Update Columns" button. A red arrow points from the dialog to the "# TM Domains" column in the results table. The results table lists genes with their transcript IDs, organisms, genomic locations, and transcript lengths. The transcript length column is highlighted with a red box.

Gene ID	Transcript ID	Organism	Genomic Location (Transcript)	# TM Domains	Transcript Length
PVP01_0208500	PVP01_0208500.1	Plasmodium vivax P01	PvP01_02_v1:352862..356879(+)	8	1551091..1552320(+)
PVP01_0412900	PVP01_0412900.1	Plasmodium vivax P01	PvP01_04_v1:530187..531415(-)	8	1011300..1011300(-)
PVP01_0509600	PVP01_0509600.1	Plasmodium vivax P01	PvP01_05_v1:429961..434135(-)	8	102700..102700(-)
PVP01_0702700	PVP01_0702700.1	Plasmodium vivax P01	PvP01_07_v1:158919..167549(-)	8	12700..12700(-)
PVP01_0817900	PVP01_0817900.1	Plasmodium vivax P01	PvP01_08_v1:778581..780290(+)	8	12900..12900(+)
PVP01_0914100	PVP01_0914100.1	Plasmodium vivax P01	PvP01_09_v1:654670..655704(+)	8	12960..12960(+)
PVP01_0936100	PVP01_0936100.1	Plasmodium vivax P01	PvP01_09_v1:1551091..1552320(-)	8	128700..128700(-)
PVP01_1011300	PVP01_1011300.1	Plasmodium vivax P01	PvP01_10_v1:501278..502747(-)	8	1300..1300(-)
PVP01_1028700	PVP01_1028700.1	Plasmodium vivax P01	PvP01_10_v1:1224465..1226132(-)	8	13100..13100(-)

The screenshot shows the results table with the transcript length column added. A red arrow points to the transcript length column header. The table lists genes with their transcript IDs, organisms, genomic locations, and transcript lengths. The transcript length column is highlighted with a red box.

Gene ID	Transcript ID	Organism	Genomic Location (Transcript)	# TM Domains	Transcript Length
PVP01_0208500	PVP01_0208500.1	Plasmodium vivax P01	PvP01_02_v1:352862..356879(+)	8	1551091..1552320(+)
PVP01_0412900	PVP01_0412900.1	Plasmodium vivax P01	PvP01_04_v1:530187..531415(-)	8	1011300..1011300(-)
PVP01_0509600	PVP01_0509600.1	Plasmodium vivax P01	PvP01_05_v1:429961..434135(-)	8	102700..102700(-)
PVP01_0702700	PVP01_0702700.1	Plasmodium vivax P01	PvP01_07_v1:158919..167549(-)	8	12700..12700(-)
PVP01_0817900	PVP01_0817900.1	Plasmodium vivax P01	PvP01_08_v1:778581..780290(+)	8	12900..12900(+)
PVP01_0914100	PVP01_0914100.1	Plasmodium vivax P01	PvP01_09_v1:654670..655704(+)	8	12960..12960(+)
PVP01_0936100	PVP01_0936100.1	Plasmodium vivax P01	PvP01_09_v1:1551091..1552320(-)	8	128700..128700(-)
PVP01_1011300	PVP01_1011300.1	Plasmodium vivax P01	PvP01_10_v1:501278..502747(-)	8	1300..1300(-)
PVP01_1028700	PVP01_1028700.1	Plasmodium vivax P01	PvP01_10_v1:1224465..1226132(-)	8	13100..13100(-)

5. Add a step to your strategy. Click on the add step button then find the search for genes with GO Terms. When you find it select it and configure the search to find all genes with the GO term “Transporter activity”. See screen shots below if you need help.

The image consists of two vertically stacked screenshots of the PlasmoDR interface, specifically the 'Add a step to your search strategy' dialog.

Screenshot 1: Main Interface Overview

- The top part shows the main search interface with a sidebar titled 'My Search Strategy' containing an 'Opened (1)' entry for 'All (404) Public (42)'.
- A red arrow points from the 'Add a step' button in the sidebar to the 'Add a step to your search strategy' dialog.
- The dialog has three tabs: 'Combine with other Genes', 'Transform into related records', and 'Use Genomic Colocation to combine with other features'. The first tab is selected.
- Under 'Combine with other Genes', there are four options: '1 INTERSECT 2' (selected), '1 UNION 2', '1 MINUS 2', and '2 MINUS 1'. A red arrow points to the '1 INTERSECT 2' option.
- Below the tabs, there are three radio button options: 'A new search', 'An existing strategy', and 'My basket'. 'A new search' is selected.
- The search bar contains the text 'go:' followed by a red arrow pointing to the search input field.
- At the bottom right of the dialog, there is a note: 'The results will be [radio button] intersected with [dropdown menu] the results of Step 1.'

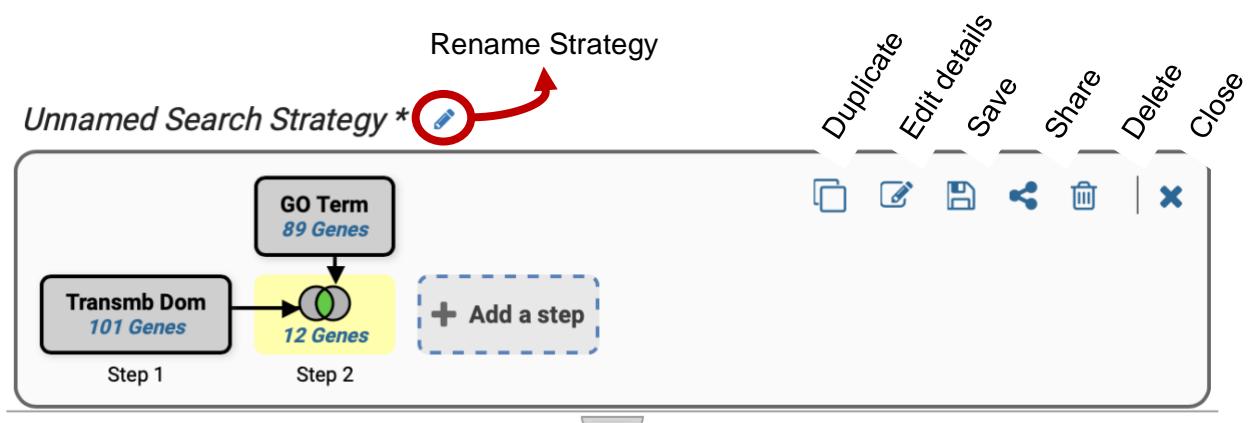
Screenshot 2: Detailed Search Configuration

- This screenshot shows the detailed configuration of the search step.
- Organism:** The 'vivax' checkbox is selected, indicated by a red arrow.
- Evidence:** 'Curated' and 'Computed' evidence types are selected.
- Limit to GO Slim terms:** The 'No' radio button is selected.
- GO Term or GO ID:** The input field contains 'GO:0005215 : transporter activity : 2' (preceded by a red arrow).

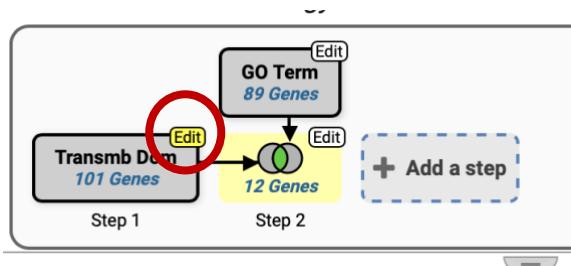
Notice that you have different options on how to combine results from searches in your strategy. What do each of the operations do?

Operator	Combined Result will contain
<input checked="" type="radio"/>  2 INTERSECT 3	IDs in common between the two lists
<input type="radio"/>  2 UNION 3	IDs from list 2 and list 3
<input type="radio"/>  2 MINUS 3	IDs unique to 2
<input type="radio"/>  3 MINUS 2	IDs unique to 3
	IDs whose features are near each other (collocated) in the genome

6. You can rename, duplicate, delete, save and share strategies (saving and sharing strategies requires creating an account and logging in). You can also rename each individual step in a strategy.



7. Revising a step in a strategy. You can revise any step in a strategy by moving your mouse over the step you want to revise until you see the edit button appear on the step.



8. Revise the first step in your strategy and change the TM domain parameter to include genes with a minimum of 5 TMs and a maximum of 12 TMs. How does this change your final results?

The screenshot shows the VEuPathDB interface. At the top, there are links for 'Opened (1)', 'All (404)', 'Public (42)', and 'Help'. Below this is the title 'Unnamed Search Strategy *' with an edit icon.

The main area displays a search strategy with two steps:

- Step 1:** 'Transmb Dom' with '101 Genes' and an 'Edit' button.
- Step 2:** 'GO Term' with '89 Genes' and an 'Edit' button.

A yellow box highlights 'Step 2' and its 'Edit' button. A red arrow points from this highlighted area to the 'Revise' button in the top right corner of the detailed step editor window.

The detailed step editor window for 'Transmb Dom' shows the following information:

- Details for step Transmb Dom** (with an edit icon)
- 101 Genes**
- Organism:** Plasmodium vivax P01
- Minimum Number of Transmembrane Domains:** 6
- Maximum Number of Transmembrane Domains:** 8
- Give this search a weight:** (with a dropdown menu)

Below the main search area, there is an 'Organism Filter' with options 'select all', 'clear all', 'expand all', and 'collapse all'. To the right, there is a 'Rows per page:' dropdown set to 50, a 'Download' button, and an 'Ad' button.

At the bottom, a modal window titled 'Revised Step' contains two input fields with arrows pointing to them:

- Minimum Number of Transmembrane Domains:** A field containing '5' with a red arrow pointing to it.
- Maximum Number of Transmembrane Domains:** A field containing '12' with a red arrow pointing to it.

The 'Revise' button at the bottom right of this modal window is also circled in red.

Exploring the Gene Page

Note: this exercise uses *TriTrypDB* (<https://TriTrypdb.org>) as an example database, but the same functionality is available on all VEuPathDB resources.

Learning objectives

Gene pages:

- Become familiar with the information in gene pages
- Navigate to and from the gene pages

1. Navigation to the Gene pages

For this exercise visit the gene page for Tb927.10.13780 (Glycogen synthase kinase 3). How did you get to this gene? (hint: copy and paste the ID in the site search, then click on the gene ID in the results.

The screenshot shows the TriTrypDB gene page for Tb927.10.13780. On the left, there are filter sections for 'Genome Genes', 'Filter Gene fields' (with options like Cellular localization, External links, Gene ID, GO terms, Transcripts), and 'Filter organisms' (with options like Trypanosomatidae and Trypanosoma). The main content area displays the gene details: 'Gene - Tb927.10.13780 Glycogen synthase kinase 3 short'. It shows the gene name or symbol (GSK3s), organism (Trypanosoma brucei brucei TREU927), and fields matched (Cellular localization, External links, Gene ID, GO terms, Transcripts). A red arrow points to the gene name 'Glycogen synthase kinase 3 short'.

2. Explore the top section of the gene page

What information is in this section? Can you easily find which chromosome this gene is located on? Does this gene have any user comments?

[Add to basket](#) [Add to favorites](#) [Download Gene](#)

Tb927.10.13780 Glycogen synthase kinase 3 short

Name: GSK3s
Type: protein coding
Chromosome: 10
Location: Tb927_10_v5.1.3361,774..3,366,257(-)

This genome is actively curated at GeneDB. User comments added to this gene will be reviewed and incorporated into the official annotation if appropriate.

Shortcuts

Species: Trypanosoma brucei
Strain: brucei TREU27
Status: Curated Reference Strain

Synteny BLAT Alignments Phenotype SNPs Transcriptomics Protein Features Proteomics

View this gene at GeneDB 
View 3 user comments, or add a comment 

Also see Tb927.10.13780 in the [Genome Browser](#) or [Protein Browser](#)

3. Explore the gene model section.

Scroll down to the gene model section of the gene page. What direction is the transcript relative to the chromosome? Does the gene have UTRs?

1 Gene models

Exons in Gene 1

Transcripts 1

Gene Models

[View in JBrowse genome browser](#)

Annotations:

- Annotated Transcripts (UTRs in White when available)
- Scroll and zoom ?

Annotations visible in the JBrowse interface:

- Tb927.10.13720 mRNA
- Tb927.10.13730
- Tb927.10.13740
- Tb927.10.13750
- Tb927.10.13760
- Tb927.10.13770
- Tb927.10.13780
- Tb927.10.13790
- Tb927.10.13800
- Tb927.10.13820
- Tb927.10.13830
- Tb927.10.13840
- Tb927.10.13850
- Tb927.10.13860
- Tb927.10.13870
- Tb927.10.13880
- Tb927.10.13890
- Tb927.10.13900
- Tb927.10.13910
- Tb927.10.13920
- Tb927.10.13930
- Tb927.10.13940
- Tb927.10.13950
- Tb927.10.13960
- Tb927.10.13970
- Tb927.10.13980
- Tb927.10.13990
- Tb927.10.14000
- Tb10.NT.148
- Tb10.NT.149
- Tb10.NT.149A
- Tb927.10.13700
- Tb927.10.13710
- Tb927.10.13720
- Tb927.10.13730
- Tb927.10.13740
- Tb927.10.13750
- Tb927.10.13760
- Tb927.10.13770
- Tb927.10.13780
- Tb927.10.13790
- Tb927.10.13800
- Tb927.10.13820
- Tb927.10.13830
- Tb927.10.13840
- Tb927.10.13850
- Tb927.10.13860
- Tb927.10.13870
- Tb927.10.13880
- Tb927.10.13890
- Tb927.10.13900
- Tb927.10.13910
- Tb927.10.13920
- Tb927.10.13930
- Tb927.10.13940
- Tb927.10.13950
- Tb927.10.13960
- Tb927.10.13970
- Tb927.10.13980
- Tb927.10.13990
- Tb927.10.14000

[View in JBrowse genome browser](#)

How long is the transcript? You can determine transcript length by expanding the Transcripts section.

[VIEW IN GBROWSE](#) [GENOME BROWSER](#)

▼ Transcripts [Download](#) [Data sets](#)

Transcript	# exons	Transcript length	Protein length
Tb927.10.13780:mRNA	1	4484	352

4. Content navigation.

How do you find/navigate to the different sections of the page? Use the “Contents” menu on the left side, type a keyword and click on the menu, click on the work to navigate to it

on the page. In the example below the word “synteny” is used. You can also click on the images in the Shortcuts section in the top of the page.

5. Running an alignment of selected sequences

- Expand the “Orthologs and Paralogs in TriTrypDB” section.
- Select a few genes from the table using the checkbox.
- Scroll to the bottom of the table and click on the Run Clustal Omega button.

6. Exploring the genetic variation section

<input checked="" type="checkbox"/>	TcYC6_0115420	Trypanosoma cruzi Y C6	protein kinase
<input type="checkbox"/>	Tc_MARK_4866	Trypanosoma cruzi marinkellei strain B7	glycogen synt alpha, putative
<input type="checkbox"/>	TevSTIB805.10.14480	Trypanosoma evansi strain STIB 805	glycogen synt
<input type="checkbox"/>	DQ04_00191000	Trypanosoma grayi ANR4	putative glyco kinase-3 alpha
<input checked="" type="checkbox"/>	TM35_000033680	Trypanosoma theileri isolate Edinburgh	putative glyco kinase-3 alpha
<input type="checkbox"/>	TvY486_1013940	Trypanosoma vivax Y486	protein kinase

Select sequence type for Clustal Omega multiple sequence alignment:

Please note: selecting a large flanking region or a large number of sequences will take several minutes.

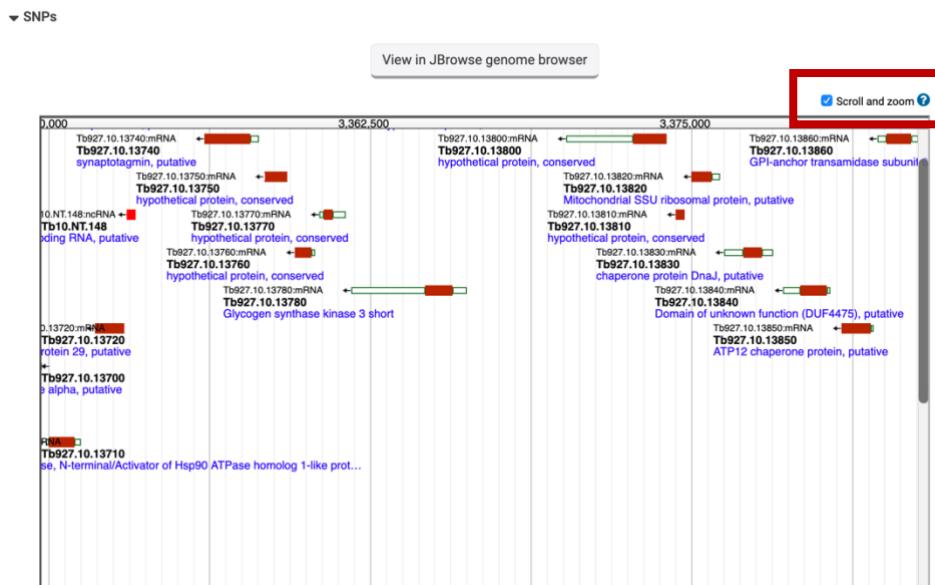
Protein CDS (spliced) Genomic

Output format:

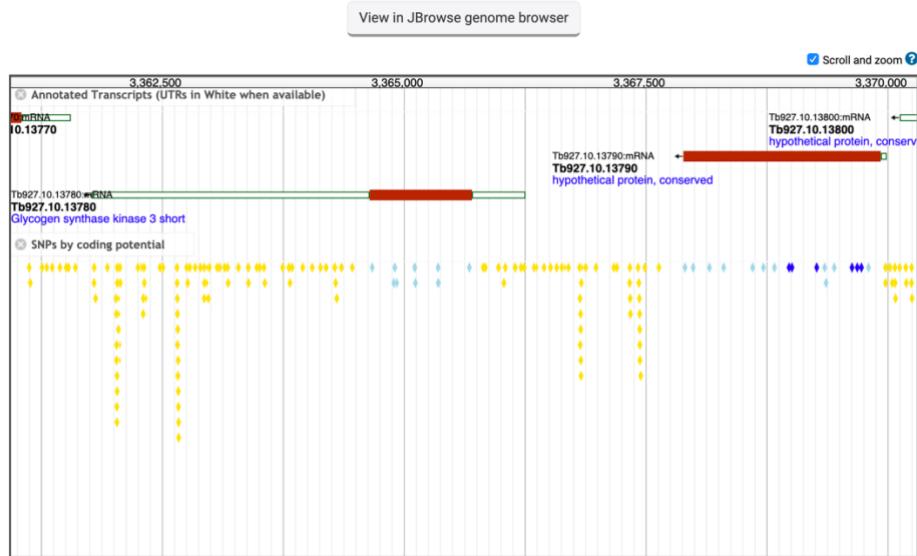
Go to the Genetic variation section of the gene page and expand the SNP section.

Notice that by default you cannot scroll within the embedded browser window. To enable scrolling, select the “Scroll and Zoom” check box in the upper right-hand side of the browser window. To scroll down within the browser window, you click and drag or use two-finger scrolling. You can also double click in an area to zoom in.

SNP color code: Dark blue (non-synonymous), light blue (synonymous), Yellow (non-



coding), Red (nonsense). What kind of SNPs are in this gene? Can you see any non-synonymous SNPs? How does this compare to the neighboring genes?



7. Explore other sections of the gene page.

Feel free to scroll around the gene page and ask questions for clarification. Here are some questions you may want to ask about this gene:

- Is there evidence that this protein is phosphorylated? (hint: go to the proteomics section and expand the Post Translational Modification section).

- b. Where is the protein localized? (hint: go to the Protein Targeting and Localization section and expand the cellular localization section).
- c. Is there any phenotypic data available for this gene? (hint: go to the Phenotype section and expand its subsections).
- d. Is there any RNA-Seq data available for this gene? (hint: go to the Transcriptomics section and expand the subsections called RNA-Seq transcription summary and Transcript Expression).

JBrowse Basics

Note: this exercise uses *TriTrypDB* as an example database, but the same functionality is available on all VEuPathDB resources.

Learning objectives:

- Navigate to the genome browser
- Use the menu and navigation bars
- Run searches
- Add pre-loaded data tracks
- Upload your own data tracks
- Configure tracks
- Download track data

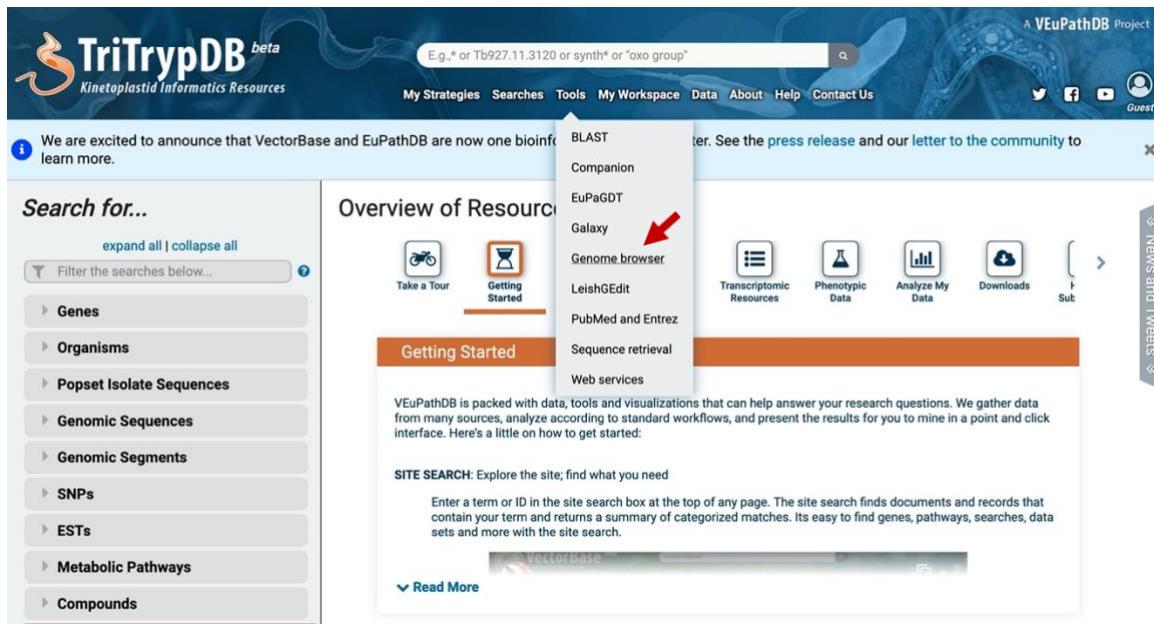
Navigating to the Genome Browser (JBrowse)

JBrowse is a fast and full-featured genome browser built with JavaScript and HTML5. You can read more about JBrowse and its features here:

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4830012/>

Links to the genome browser are available from multiple locations:

- a. The tools menu in the banner of any page.



The screenshot shows the TriTrypDB homepage. At the top, there's a search bar and a menu bar with links like 'My Strategies', 'Searches', 'Tools', 'My Workspace', 'Data', 'About', 'Help', and 'Contact Us'. A 'Guest' button is also present. A prominent red arrow points to the 'Genome browser' link under the 'Tools' menu. The main content area has sections for 'Overview of Resources' and 'Getting Started', along with links to 'BLAST', 'EuPaGDT', 'Galaxy', 'LeishGEedit', 'PubMed and Entrez', 'Transcriptomic Resources', 'Phenotypic Data', 'Analyze My Data', 'Downloads', and 'Web services'.

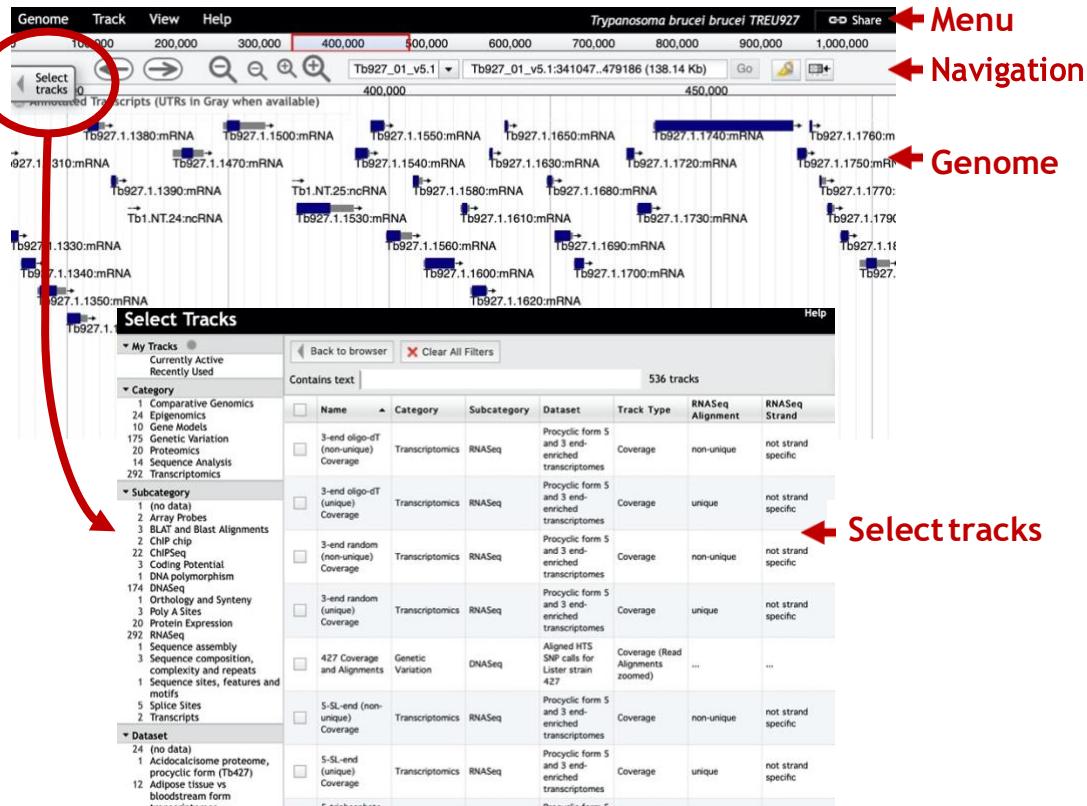
- b. From record pages such as gene, SNP or genomic sequence pages – these links are usually to a specific JBrowse configuration that includes data relevant to the section on that record page. For example, a JBrowse link from an RNAseq dataset on the gene page would display the gene of interest

along with the RNAseq data as density or coverage plots. These links are usually indicated by “View in JBrowse genome browser” button.

View in JBrowse genome browser

Getting around JBrowse.

- Use any of the above described JBrowse linking strategies to get to the genome browser.
- Once in JBrowse examine the following features:
 - The **menu bar**: located at the top of the JBrowse frame. This includes the Genome menu, Track menu, View menu, Help menu and the Sharing link. What do each of these do/provide?
 - The **navigation bar**: located below the menu bar. This contains zooming (magnifying glass icons), panning (left/right arrows) and highlighting (yellow highlighter) buttons, reference sequence selector (drop down with sequences from the selected genome sorted by length), a text box to search for features such as gene IDs and overview bar which shows the location of the region in view.
 - The **genome view**: this is where the data tracks are displayed.



- Selecting tracks: click on the “select track” button (top left). You can search/filter for tracks and then select them for display by checking the

check box next to the track name.

Navigating to a specific gene in JBrowse.

The goal of this step is to navigate to the sedoheptulose-1,7-bisphosphatase (SBPase) gene of *T. brucei* 927.

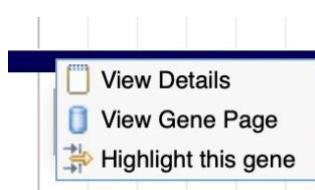
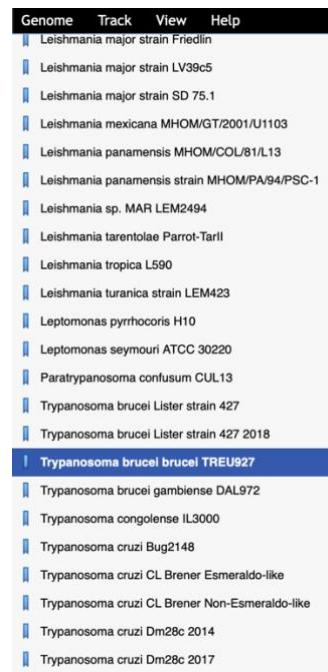
- f. Make sure the *Trypanosoma brucei* brucei TREU927 genome is selected from the genome menu.
- g. Start typing the word sedoheptulose in the search box. After a few seconds you should see the result of the search (do not hit enter). Select the gene from the search dropdown. This will take you to Tb927.2.5800.



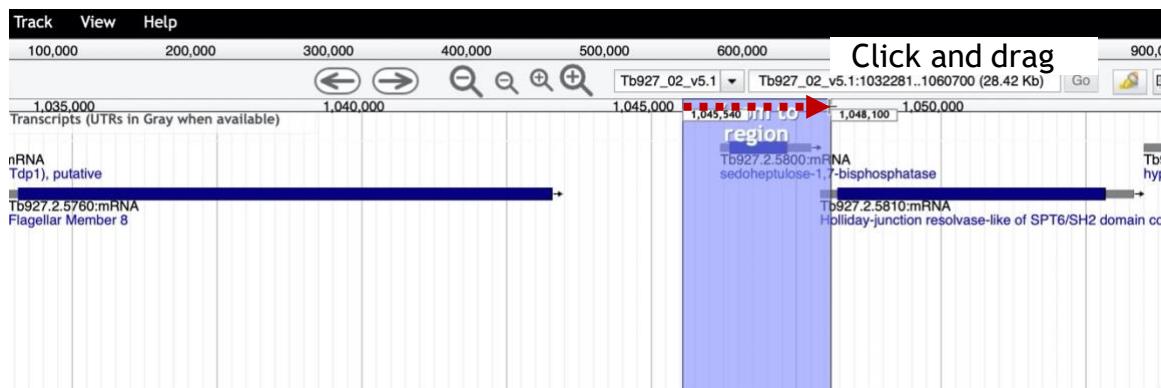
- h. You can get information about any feature in the genome view window by clicking on it. Click on the gene feature. What information is available in the popup?

Species:	Trypanosoma brucei brucei TREU927
ID:	Tb927.2.5800:mRNA
Gene ID:	Tb927.2.5800
Gene Type:	Protein Coding
Description:	sedoheptulose-1,7-bisphosphatase
5' UTR:	1046195..1046356
CDS:	1046357..1047355
3' UTR:	1047356..1047765
Download:	CDS protein
OrthoMCL:	OG5_134853
Links:	JBrowse Gene Page

- i. You can also right click (or control click) on a feature to display the context menu which provides quick links to highlight a feature, go to the feature page (like the gene page) or get the info popup (the same one you get when you click on the feature).
- j. What genes are immediately upstream and downstream of SBP? (Hint: use the zoom out button in the navigation bar). What is the difference between the small and large zoom buttons? (*Tip*: another way to zoom in and out is by clicking on shift and the up or down arrows. What happens if you click shift

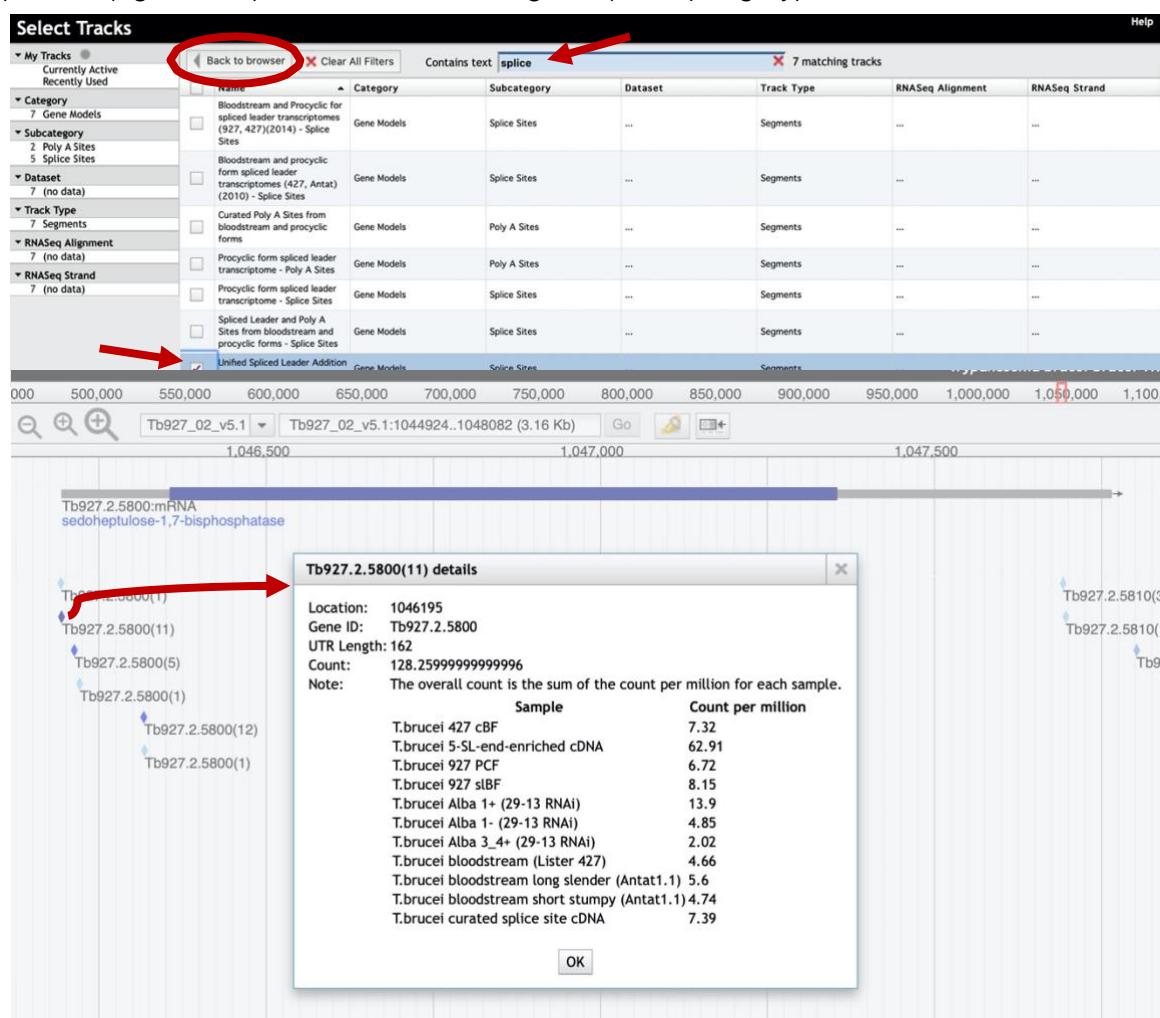


and left or right arrows? *Tip2:* you can also zoom in by clicking and dragging your cursor in the location ruler in the navigation bar).



Exploring transcription start sites.

Are you confident about the gene transcription start? (Note: gene features are in blue (left to right) or red (right to left) with untranslated regions (UTRs) in grey).



What additional data track would be useful for you to assess this? (hint: Click on the “Select Tracks” button to reveal all available tracks. Now type the word “splice” in the “contains text” box. This will filter all tracks that contain the word splice. Find the one called “Unified Splice Leader Addition Sites” and select it.

Click on the “Back to browser” button). What do the different diamond colors mean? Click on them and see if you can figure this out from the popups? Which color provides the most evidence for a splice junction?

Exploring synteny between genomes.

Synteny helps define conservation of homologous genes and gene order between genomes.

- Go to the “Select Tracks” tab on the left of the page and turn on the track called “Syntenic Sequences and Genes”. How did you find this track? One option is to click on the “Comparative Genomics” category on the left side to filter the

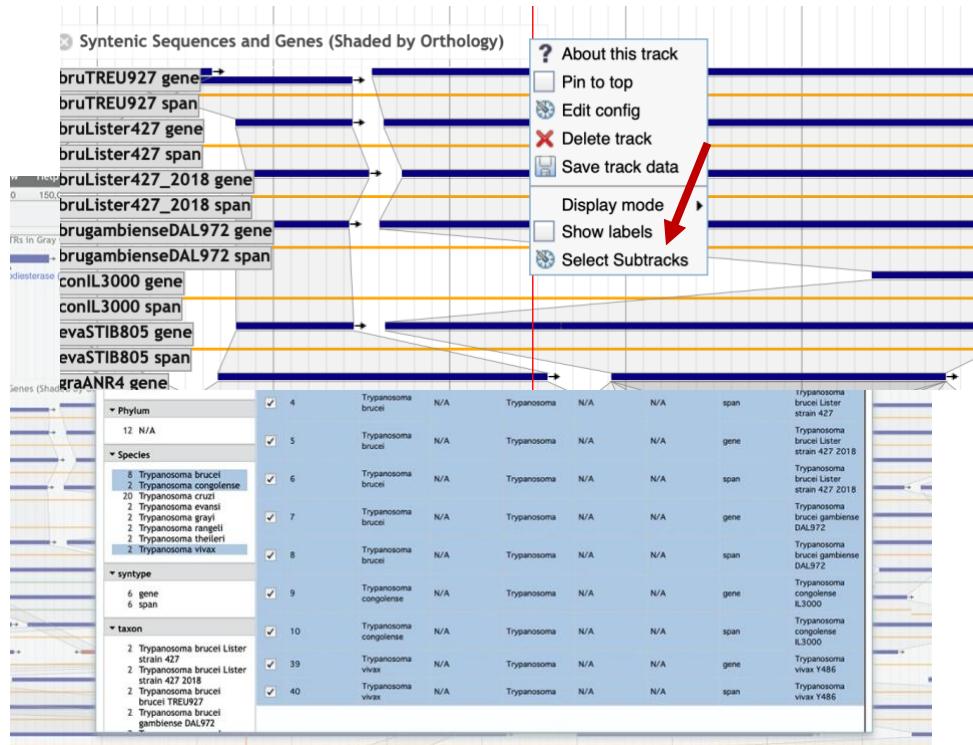
The screenshot shows the 'Select Tracks' interface. On the left, there's a sidebar with 'My Tracks' (Currently Active, Recently Used) and a 'Category' section expanded to show 'Comparative Genomics' (1 Comparative Genomics, 24 Epigenomics, 10 Gene Models, 175 Genetic Variation, 20 Proteomics, 14 Sequence Analysis, 292 Transcriptomics). The main area has a search bar with 'Contains text' and a table with columns: Name, Category, Subcategory, Dataset, Track Type, RNASEq Alignment, and RNASEq Strand. A single row is selected: 'Syntenic Sequences and Genes (Shaded by Orthology)' under 'Comparative Genomics / Orthology and Synteny'. The table shows 1 matching track.

tracks.

- Return to the browser by clicking “Back to Browser” and zoom out so you can see a couple of genes on either side of SBP (does not have to be exact)
- Configure the synteny track to include the following species subtracks: *Trypanosoma brucei* 927, *T. brucei* 427, *T. brucei gambiense*, *T. congolense*, *T. evansi*, *T. grayi*, *T. theileri* and *T. vivax*.
 - To configure the subtracks:
 - Click on the down arrow in the track name

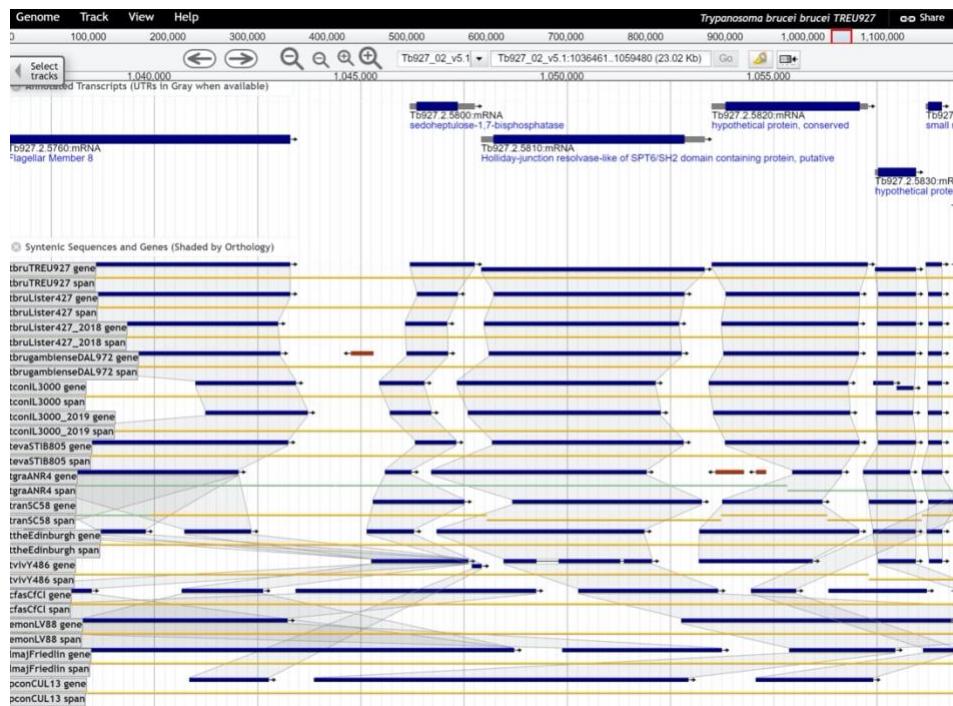


- Select the option called “Select Subtracks” from the menu

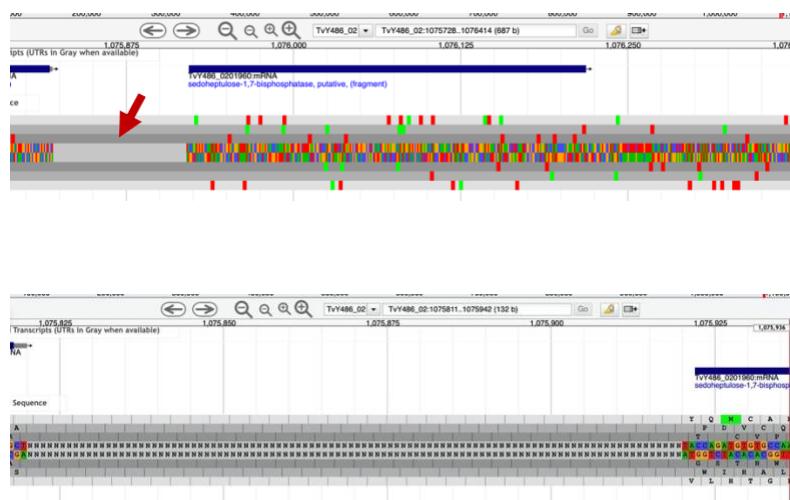


- In the next popup first uncheck all organisms, second use the filters on the left to select *Trypanosoma*, third select the species of interest (note that you should select both the gene and span subtracks for each species), fourth click on the save button at the bottom of the popup.

- What does the synteny track in this region look like? Feel free to zoom out some more. Are genes (in general) similarly organized between these species? What does the shading between genes mean?
- What direction is the SBPase gene relative to the chromosome?
- What genes are upstream and downstream of the SBPase? Are these genes syntenic?
- What does synteny look like if you add more distantly related species? Does
- SBPase appear to have orthologs in *Leishmania*? *Endotrypanum*? *Critidilia*?



- Examine the gene corresponding to the *T. vivax* SBPase in the synteny track. Hover over the gene image to find the gene name in the popup. Does this gene appear to be a fragment? What could be some possible reasons for this?
- Do you think all the genomes in the database are fully sequenced? Is it possible that gaps in sequence exist in the available genomes? Let's find out if there is a gap next to the SBPase gene in *T. vivax*:
 - Select *T. vivax* from the list of genomes in the menu bar.
 - Turn on the **annotated transcripts** and the **Reference sequence** tracks.
 - Search for the SBPase gene by typing "sedoheptulose" in the search box then select the gene.
 - Zoom to about 600bps. Do you see something missing on the left side of the gene?
 - Zoom in to this area (click and drag). What do you see? What do all of these Ns mean?



Exploring other data tracks in Jbrowse.

For this example, we will view *T. brucei* data, so the data tracks you turn on will display data only if the data is aligned to the *T. brucei* genome. Return to the SBPase gene in *T. brucei* by searching for the gene ID in the (Tb927.2.5800) in ‘Landmark or Region’ to redirect the browser. Then zoom to the area between 0.7M and the end of the chromosome.

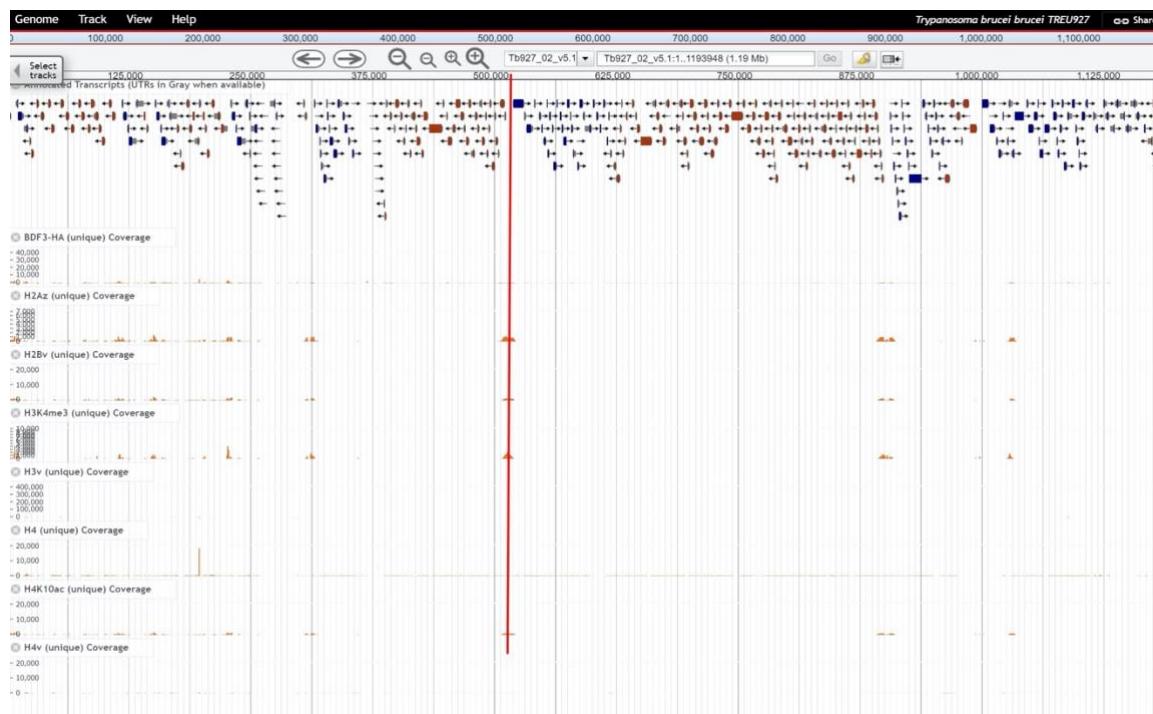
Turn on the ChIP-seq coverage plots and turn off the syntenic gene and region tracks. The data tracks are from an experiment called: **ChIP-Seq - Four histone Variants ChIP-Seq Coverage aligned to T brucei TREU927 (Cross) (linear plot)**. For this experiment, chromatin was immunoprecipitated using several different histone antibodies. The DNA that precipitated with the histone was sequenced and aligned to the *T. brucei* TREU927 genome. Peaks in the sequence coverage plots represent areas of histone binding. Different histone variants can be associated with start and termination sites for transcription (<http://www.ncbi.nlm.nih.gov/pubmed/19369410>)

- You may need to adjust the y-axis scaling to bring the tracks into proper view (try

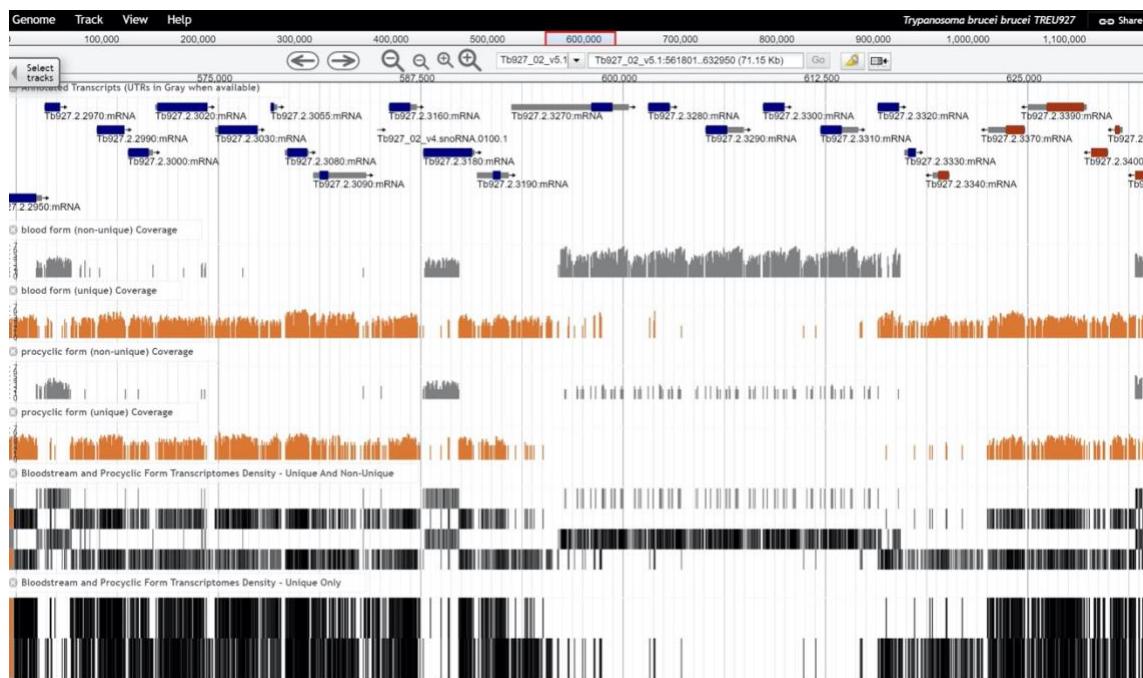
Name	Category	Subcategory	Dataset	Track Type	RNASeq Alignment	RNASeq Strand
BDF3+H4 (unique) Coverage	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Coverage
ChIP-Seq - Four histone Variants Density - Unique And Non-Unique	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Multi-Density
H2A (unique) Coverage	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Coverage
H2Bv (unique) Coverage	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Coverage
H3K4me3 (unique) Coverage	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Coverage
H3V (unique) Coverage	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Coverage
H4 (unique) Coverage	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Coverage
H4K10ac (unique) Coverage	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Coverage
H4v (unique) Coverage	Epigenomics	ChIPSeq	ChIP-Seq - Four histone Variants	Coverage

setting the score range to “global” by mousing over the track name, clicking the dropdown arrow and selecting “Change Score Range”).

- What does this data show you?
- Roughly how many polycistronic units does this chromosome have? Zoom out to the entire chromosome.

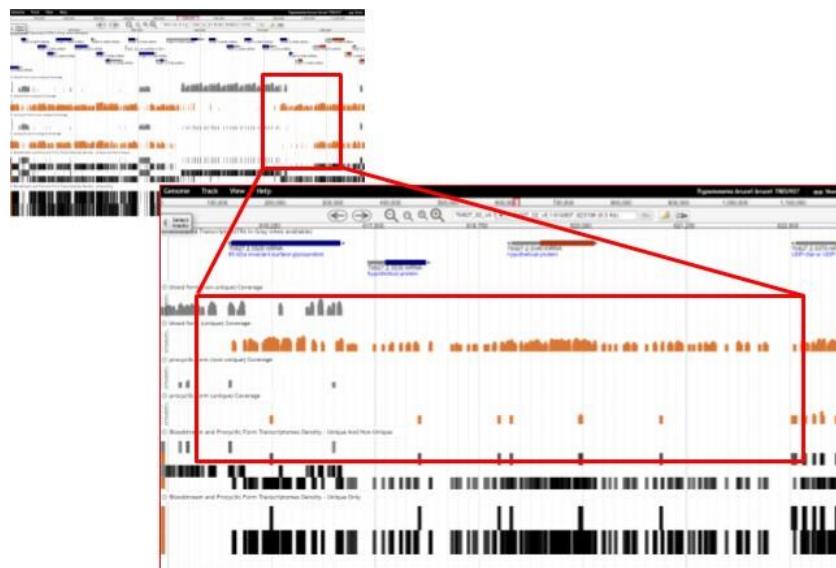


- Do the ChIP-seq peaks correlate with the direction of gene transcription (blue vs. red)?
- Now zoom back to around 50Kb. Turn off the ChIP-Seq tracks and turn on the RNASeq Coverage track called: **Bloodstream and Procyclic Form Transcriptomes mRNASeq Coverage aligned to T brucei TREU927**.



- Move to the **region around 0.6Mbs of the chromosome** (you should be on chromosome 2) and turn on all four subtracks. Take note of the orange and grey bars in the coverage plots. What do you think the grey bars indicate?

- Now zoom out to 100Kb – do you see a difference between the blood and procyclic forms?



- Zoom in to a gene that looks like it is differentially expressed. What are your conclusions? Are the reads supported by unique or non-unique reads?
- Can you turn on additional tracks that may give some more support to your conclusions?

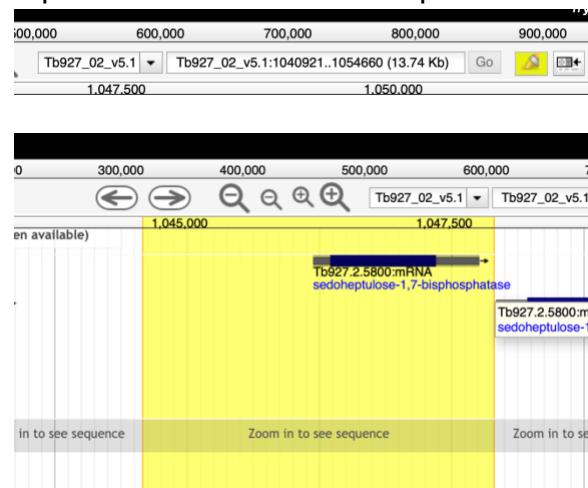
Hint: turn on the EST and *T. brucei* protein expression evidence tracks.

- Is there any proteomics evidence for this region?
- How about EST evidence? Click on an EST graphic (glyph) to get additional information.
- Turn off the RNA-seq graphs and make sure the *T. brucei* protein expression evidence tracks are on. **Zoom out to 500Kb**. Explore the evidence for gene expression based on mapped peptides from proteomics experiments – which gene in this view has the highest number of peptide hits? Try looking at the “All MS/MS peptides (feature density)” track for an overview.



Retrieving data from and uploading your own tracks to JBrowse

- k. Downloading sequence in FASTA format from a region of interest:
 - i. Make sure the “annotated transcripts” and the “reference sequence” tracks are turned on.
 - ii. Click on the “highlight a region” button in the navigation bar. It should turn yellow when activated.
 - iii. Click and drag in the genome view region and select the area you would like to highlight.
 - iv. Click on the down arrow on the reference sequence track and select “Save track data”.
 - v. In the next popup window you can keep everything as the default and either save or view the sequence.





i. Uploading data to JBrowse:

JBrowse can accept several standard-format data files by direct upload or through a URL if the data is stored remotely. Some file formats like BAM and VCF require indexing before uploading. In this exercise we will download a bigwig file from GEO and then upload it to JBrowse:

i. Go to this GEO sample record:

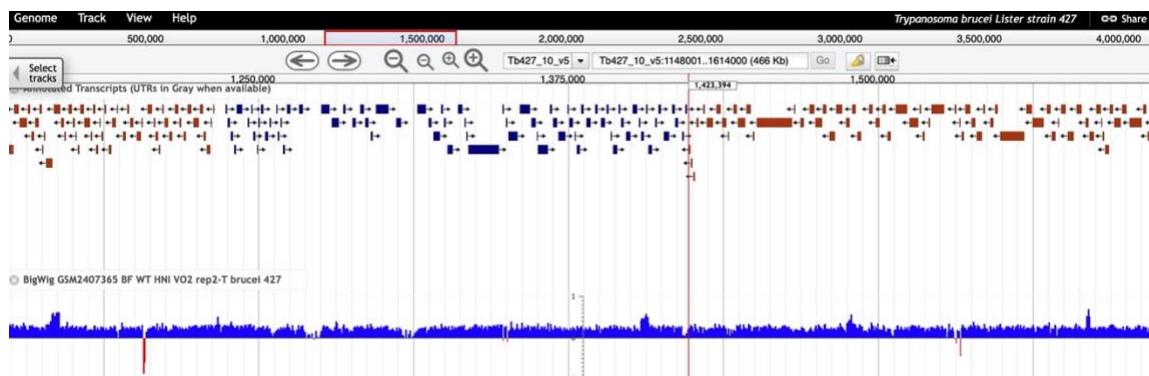
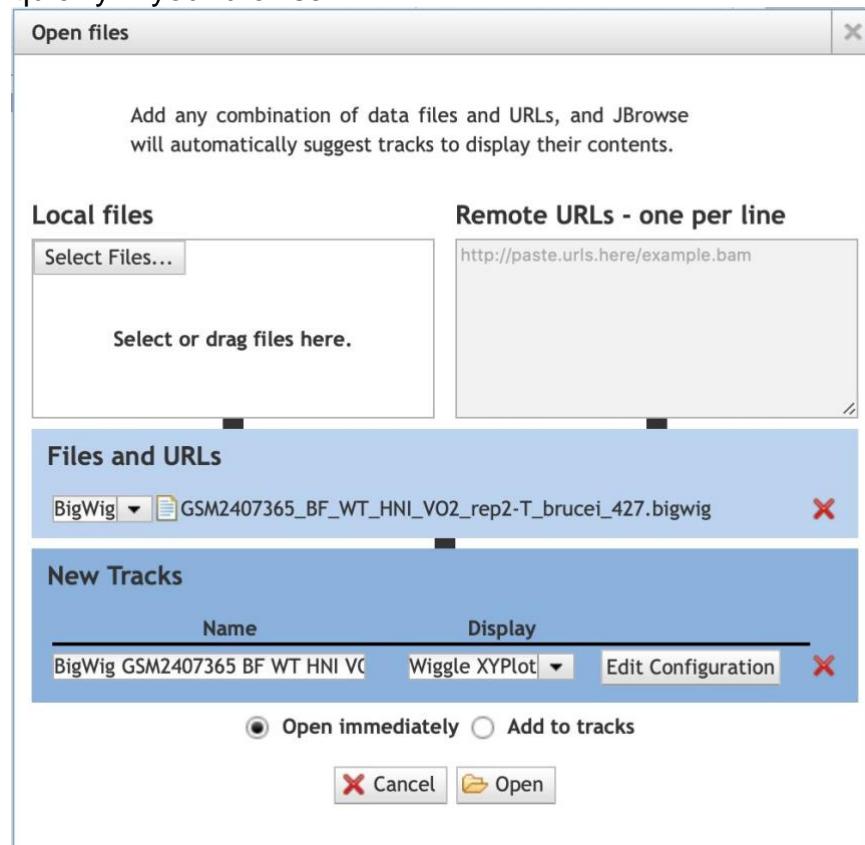
<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM2407365>

ii. Scroll down to the bottom of the page and download the bigwig file with the http link.

Supplementary file	Size	Download	File type/resource
GSM2407365_BF_WT_HNI_VO2_rep2-T_brucei_427.bigwig	12.4 Mb	(ftp)(http)	BIGWIG

- iii. Once the file is downloaded go to JBrowse and select *Trypanosoma brucei brucei* Lister 427 as the reference genome (hint: use the Genome link in the menu panel, top left).
- iv. Turn on the track for annotated transcripts if it is not on already.

- v. Click on the Tracks menu item and select “Open track file or URL”.
- vi. In the popup click on select file then select the file you just downloaded. JBrowse should automatically recognize that the file is in bigwig format.
- vii. Click on the Open button. The bigwig output should appear very quickly in your browser.



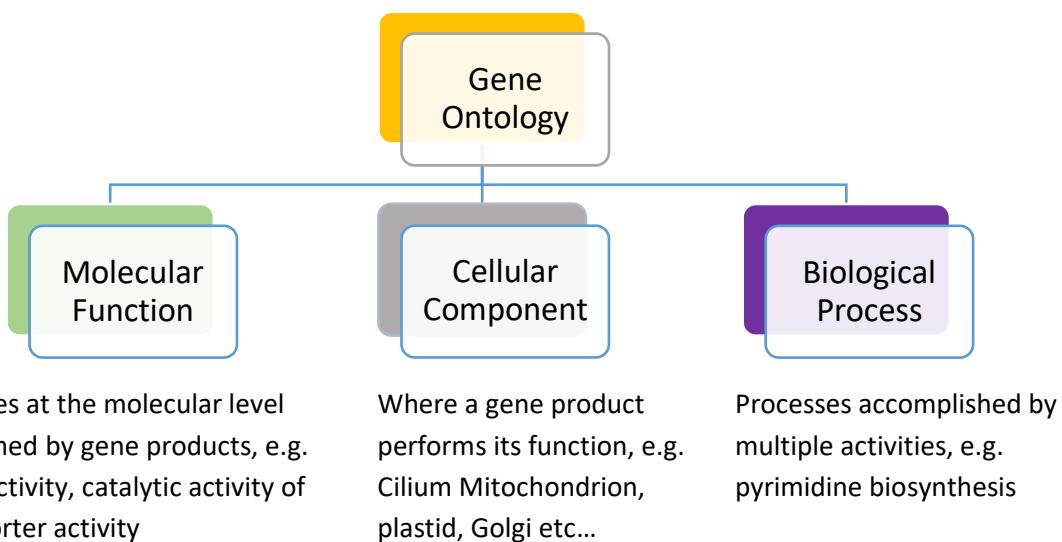
Gene Ontology (GO) Enrichment

Learning objectives:

- Run a GO enrichment analysis
- Explore GO enrichment results

Background:

The gene ontology describes the knowledge of biological sciences and divides this knowledge into three broad categories: Molecular function, cellular component and biological process.



To learn more about Gene ontology please visit: <http://geneontology.org/docs/ontology-documentation/>

Genes can be assigned a GO term either manually or computationally based on transfer by similarity, by domain association or by many other computational methods. GO terms can be used in enrichment analysis!

For example: Does my list of genes have an over-representation of specific GO terms compared to the rest of the genome?

A standard enrichment method is Fisher's exact test which is a statistical test used when analyzing contingency tables. Typically used when you have a small sample size. But when you are doing enrichment analysis on a list of genes with the background being the whole genome, your sample size is not small. As a result, the P- value you get from a Fisher's exact test might be misleading.

With a small sample size, a P-value of less than 0.05 is considered significant (5% chance of being wrong/random). But if you are doing an enrichment analysis with all genes in the genome

then each gene can be considered a test, so the chances of a type one error becomes higher. As a result, you should correct for this which can be done in different ways including Benjamini-Hochberg false discovery rate (FDR) or Bonferroni adjusted p-value

1. In order to run a GO enrichment analysis, we need a list of genes to test. This can be a list of gene IDs from your results that you can upload using the ID search or a gene list resulting from a search you conducted in the database. For this example, in ToxoDB, we will identify genes that are differentially regulated over time.

- a. Navigate to the RNA-Seq searches and find the data set called “**Oocyst Time Series (M4)**” from Fritz et al. A fast way of getting to the RNA-Seq searches is type ‘ma’ in the filter box on the left of the home page then click on the RNA-Seq Evidence link. See image below.

- b. The RNA-Seq evidence page include a list of all the data sets that are loaded in the database. To quickly find a dataset you can start typing key words in the “Filter Data Sets” box. For example, start typing the word “oocyst”.
- c. Once you find the data set of interest click on the fold change option. This will

Identify Genes based on RNA-Seq Evidence

make available to you all the parameters that you can manipulate to search this data set. For this exercise identify genes that are upregulated by 20-fold in the day 4 and day 10 time points compared to the day 0 time point. Parameters to set:

1. Up-regulated
2. 20-fold

3. Maximum
4. Day 0
5. Minimum
6. Day 4 and 10

Identify Genes based on *T. gondii* ME49 Oocyst Time Series (M4) RNA-Seq (fold change)

For the Experiment
Oocyst Time Series (M4) - Sense
return protein coding genes
that are up-regulated ?
with a Fold change >= 20 ?
between each gene's maximum expression value ?
(or a Floor of 10 reads ?)
in the following Reference Samples ?

day 0
 day 4
 day 10

[select all](#) | [clear all](#)

and its minimum expression value ?
(or the Floor selected above)
in the following Comparison Samples ?

day 0
 day 4
 day 10

[select all](#) | [clear all](#)

Example showing one gene that would meet search criteria
(Dots represent this gene's expression values for selected samples)

Up-regulated

Expression

Reference Samples Comparison Samples

Minimum Expression Value Comparison

Expression Value Reference

20 fold

[Get Answer](#)

- d. Once you have set the parameters you can click on the “Get Answer” button at the bottom of the search. This will return a one-step search strategy. How many genes did you get?
2. To run a GO enrichment analysis on these results, do the following:
 - a. Click on the Analyze Results tab right above the list of genes (arrow in image below).

My Search Strategies

Opened (1) All (1) Public (17) Help

Unnamed Search Strategy *

TgM4 Oocyst RNA-Seq (fc)
1,029 Genes

+ Add a step

Step 1

1,029 Genes (970 ortholog groups) [Revise this search](#)

[Gene Results](#) [Genome View](#) **Analyze Results**

Organism Filter
[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#)
 Hide zero counts

Search organisms...

Eimeriidae

0

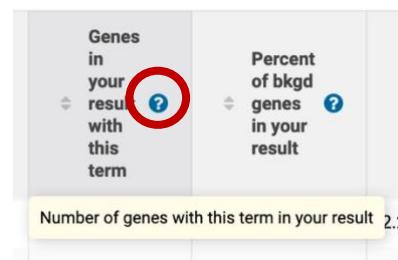
Gene Results Rows per page: 20

1 2 3 ... 52 Next

Download Add to Basket Add Columns

Page 48 of 121

- b. Clicking on the “Analyze Results” tab will reveal the different analyses that you can run on your results. Besides GO enrichment what other analyses are available?
- c. Click on the GO enrichment option. This will reveal the parameters that you can modify. For the purpose of this exercise, keep all the defaults and click on “Submit”.
- d. What is the top enriched GO term from this analysis?
- e. What do each of the columns in the analysis table represent? (hint: move your mouse over the question mark next to each column header to get more information)



A screenshot of the VEuPathDB Gene Results page. At the top, there are three tabs: 'Gene Results' (selected), 'Genome View', and 'New Analysis'. On the left, there is a vertical sidebar with a 'Hide Organism Filter' button. The main content area displays the message 'Analyze your Gene results with a tool below.' and three analysis tools:

- Gene Ontology Enrichment:** Shows a network diagram centered around the 'GO' term.
- Metabolic Pathway Enrichment:** Shows a network diagram of metabolic pathways.
- Word Enrichment:** Displays a list of enriched terms: kinase, phosphatase, exported, membrane.

- f. Try rerunning the GO enrichment analysis but this time select the Molecular Function ontology. What is the top enriched GO term?

Gene Results | Genome View | Gene Ontology Enrichment  | Gene Ontology Enrichment*  | Analyze Results

[Rename This Analysis | Duplicate]

Gene Ontology Enrichment

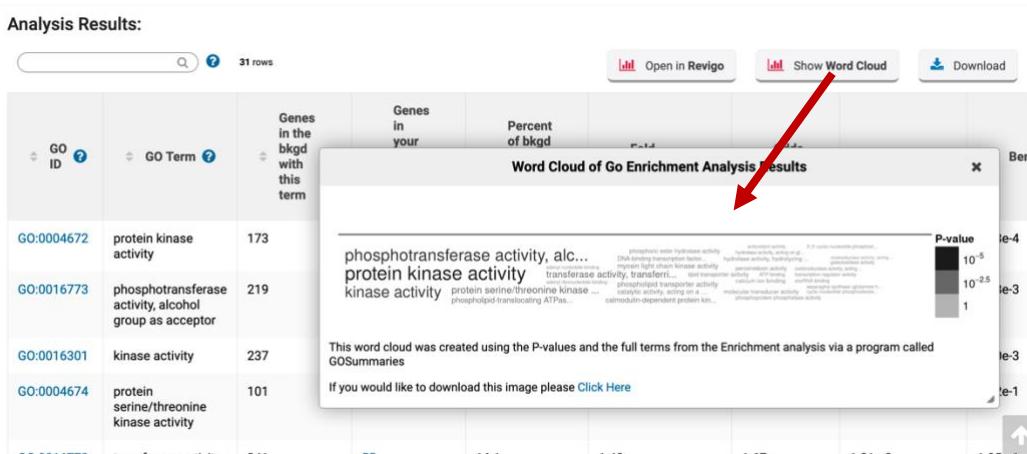
Find Gene Ontology terms that are enriched in your gene result. [Read More](#)

Parameters

Organism 	<input type="text" value="Toxoplasma gondii ME49"/> 
Ontology 	<input type="radio"/> Cellular Component <input checked="" type="radio"/> Molecular Function  <input type="radio"/> Biological Process
Evidence 	<input checked="" type="checkbox"/> Computed <input checked="" type="checkbox"/> Curated select all clear all
Limit to GO Slim terms 	<input checked="" type="radio"/> No <input type="radio"/> Yes
P-Value cutoff 	<input type="text" value="0.05"/> (0 - 1)

Submit

- g. Click on the “Word Cloud” button above the analysis results. What does this do? (See image below).



Additional resources:

Gene Ontology:

<http://geneontology.org/docs/ontology-documentation/>

Enzyme Commission numbers:

<https://www.qmul.ac.uk/sbcs/iubmb/enzyme/>

More info on Fischer's exact test:

<http://www.biostathandbook.com/fishers.html>

Fisher's Exact Test and the Hypergeometric Distribution (the M&M example):

<https://youtu.be/udyAvvaMjfM>

Some more info about Odds ratios:

<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2938757/>

False discovery rates and P value correction:

<http://brainder.org/2011/09/05/fdr-corrected-fdr-adjusted-p-values/>

GO Slim:

<http://www-legacy.geneontology.org/GO.slims.shtml>

REVIGO:

<https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0021800>

Advanced Search Strategies

Note: this exercise uses *PlasmoDB.org* as an example database, but the same functionality is available on all VEuPathDB resources.

Learning objectives:

- Integrate diverse datatypes in a search strategy
- Leverage orthology and phylogenetic profile searches

This exercise walks you through the process of building a multi-step strategy, integrating different datatypes to. The final search strategy identifies *Plasmodium* genes that are likely secreted, or membrane bound, highly polymorphic, “essential” for parasite survival, not conserved in mammals and expressed in liver stages of the *Plasmodium* life cycle. There are many ways to build these strategies and order the steps to reach a similar answer.

1. Identify all genes in PlasmoDB that are predicted to have a secretory signal peptide as defined by SignalP. An easy way to identify a search type is to filter the searches on the left of the home page. Start typing a word to identify the search type. For example, start typing the word “secreted”, you should see the searches being filtered even before you finish typing the complete word.

The screenshot shows the PlasmoDB beta homepage. On the left, there is a sidebar titled "Search for..." with various filters like "expand all" and "collapse all". A red arrow points from the "expand all" button to a dropdown menu labeled "Filter the searches below...". Below this, a list of categories includes "Genes", "Organisms", "Popset Isolate Sequences", "Genomic Sequences", "Genomic Segments", "SNPs", "SNPs (from Array)", "ESTs", and "Metabolic Pathways". In the center, there is a large search bar with the placeholder "Search for...". A red arrow points from the search bar to the input field where "secre" is typed. Below the search bar, there is a section titled "Protein targeting and localization" with a sub-section "Predicted Signal Peptide". Another red arrow points from the "Predicted Signal Peptide" link to the text "finds documents and records that contain your term and returns a summary of categorized matches. Its easy to find genes, pathways, searches, data sets and more with the site search." At the top of the page, there is a navigation bar with links for "My Strategies", "Searches", "Tools", "My Workspace", "Data", "About", "Help", and "Contact Us". The top right corner features a user profile for "Omar" and social media icons for Twitter, Facebook, and YouTube. A sidebar on the right contains links for "Analyze My Data", "Downloads", and "How to Submit Data".

2. Click on the search for genes by predicted signal peptide. On the next page select all organisms and click on the get answer button at the bottom of the page.
3. The next step is to combine the signal peptide results with results of genes that are predicted to have at least one transmembrane domain (TM). Click on the add step

Identify Genes based on Predicted Signal Peptide

Organism

*Note: You must select at least 1 values for this parameter.
45 selected, out of 45*

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#)

Filter list below... ?

- ▶ Plasmodium adleri
- ▶ Plasmodium berghei
- ▶ Plasmodium billcollinsi
- ▶ Plasmodium blacklocki
- ▶ Plasmodium chabaudi
- ▶ Plasmodium coatneyi
- ▶ Plasmodium cynomolgi
- ▶ Plasmodium falciparum
- ▶ Plasmodium fragile
- ▶ Plasmodium gaboni
- ▶ Plasmodium gallinaceum
- ▶ Plasmodium inui
- ▶ Plasmodium knowlesi
- ▶ Plasmodium malariae
- ▶ Plasmodium ovale curtisi
- ▶ Plasmodium praefalciparum
- ▶ Plasmodium reichenowi
- ▶ Plasmodium relictum
- ▶ Plasmodium vinckei
- ▶ Plasmodium vivax
- ▶ Plasmodium vivax-like sp.
- ▶ Plasmodium yoelii

[select all](#) | [clear all](#) | [expand all](#) | [collapse all](#)

Advanced Parameters

[Get Answer](#)

button in the search strategy panel.

My Search Strategies

[Opened \(1\)](#) All (415) Public (42) Help

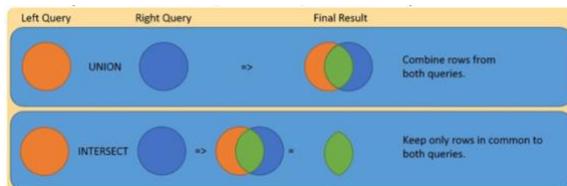
Unnamed Search Strategy *

Step 1

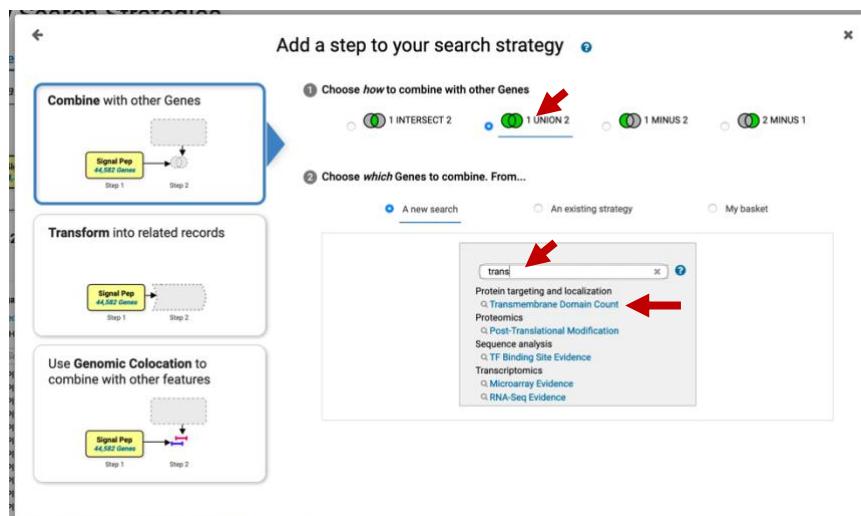
Signal Pep
44,582 Genes

+ Add a step

The popup window offers you option to add additional steps and ways to combine the searches (intersect, union, minus). For this exercise we are interested in finding genes that a signal peptide or a TM domain or both. What operation will you use to combine the searches – Union or Intersect?



Once you select the option for combining the searches, find the search for transmembrane domain count. Notice that you can use the same query filtering mechanism as before. Start typing transmembrane to find this search. Once you find it click on to open the search parameters.

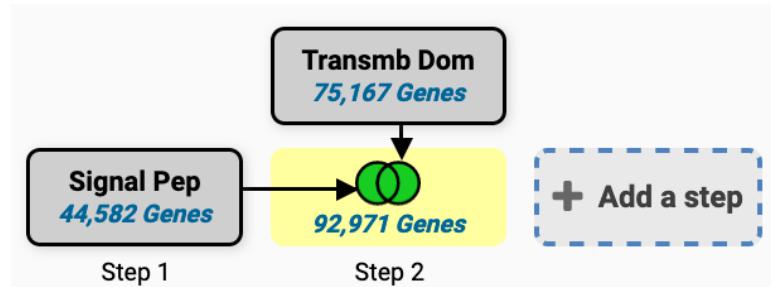


- For the TM search, again select all organisms, use the default parameters and click on the get answer button.

This screenshot shows the search parameters dialog box for "Transmembrane Domains".

- Organisms:** A list of Plasmodium species, all of which are selected (indicated by checked checkboxes).
- Minimum Number of Transmembrane Domains:** A text input field containing the value "1".
- Maximum Number of Transmembrane Domains:** A text input field containing the value "99".
- Run Step:** A button at the bottom right.

5. How many genes did you get? Since you used a union the number of results should be more than each of the individual steps that were combined.



6. Next, identify genes from step 2 that contain at least 5 non-synonymous SNPs (non-synonymous SNPs are single nucleotide polymorphisms that result in an amino acid change). Were you able to find the SNP search by clicking on add step and filtering the searches with a keyword? Which operation will you select to combine the searches?
 7. On the Genes by SNP characteristics search popup, select Plasmodium falciparum from the drop down and select all available isolates by selecting the checkbox at the top of

The screenshot shows the "Add a step to your search strategy" interface. It includes three main sections: "Combine with other Genes", "Transform into related records", and "Use Genomic Colocation to combine with other features". The "Combine with other Genes" section is active, showing a diagram of the search flow and a dropdown menu for combination methods. A red arrow points to the "Choose how to combine with other Genes" dropdown, which is set to "2 INTERSECT 3". Another red arrow points to the "Choose which Genes to combine. From..." dropdown, which is set to "A new search" and shows a search bar with "snp" and a dropdown menu with "SNP Characteristics" selected.

the filter panel (See image below).

The screenshot shows the "Search for Genes by SNP Characteristics" filter panel. It includes a dropdown for "Organism" set to "Plasmodium falciparum 307" and a "Set of Samples" section with a table showing sample types and their distribution. Red arrows point to the "Organism" dropdown and the "Sample type" table.

Sample type	Remaining Set of Samples	Set of Samples	Distribution
Blood	12 (8%)	201 (100%)	(100%)
Specimen from organism	189 (94%)	189 (94%)	(94%)

8. Next scroll down and select the following parameters. SNP class = Non-synonymous. Number of SNPs of above class ≥ 5 . After you select these parameters, scroll down to the bottom and click on Run Step.

[←](#) Add a step to your search strategy [?](#)

[expand all](#) | [collapse all](#)

② Read frequency threshold

80%

③ Minor allele frequency \geq

0

④ Percent isolates with a base call \geq

20

⑤ SNP Class

Non-Synonymous ←

⑥ Number of SNPs of above class \geq

5 ←

⑦ Number of SNPs of above class \leq

What do the results look like? What species are represented in the results? Is this surprising? Remember that your last search only queried *P. falciparum* data.

My Search Strategies

Opened (1) All (415) Public (42) Help

Unnamed Search Strategy * [🔗](#)

The screenshot shows a search strategy flowchart at the top. It starts with 'Transmb Dom' (75,167 Genes), which leads to 'SNPs' (3,806 Genes). From 'SNPs', it branches into two paths: 'Signal Pep' (44,582 Genes) leading to '92,971 Genes', and a direct path leading to '1,578 Genes'. The '1,578 Genes' path is highlighted with a dashed blue border and has a '+ Add a step' button next to it. Below the flowchart, the text '1,578 Genes (6,987 ortholog groups)' is displayed. A note says 'Some Genes in your combined result have Transcripts that were not returned by one or both of the two input searches. [Explore](#)'. The main results table has tabs for 'Gene Results', 'Genome View', and 'Analyze Results'. The 'Gene Results' tab is active. It shows 'Genes: 1,578' and 'Transcripts: 1,597'. There is a checkbox for 'Show Only One Transcript Per Gene'. The table includes columns for Gene ID, Transcript ID, Genomic Location (Gene), Product Description, and Ortholog Group. The results list includes entries for rifin (OG6_100719), erythrocyte membrane protein 1 (PfEMP1), and others. On the left, there is an 'Organism Filter' sidebar with a search bar and a list of Plasmodium species.

9. Determine how many of these genes are also differentially expressed in liver stages.
 Click on add step then search for the RNA-seq search. Type RNA in the search filter in the popup.
10. On the next page find data that queries liver stages. You can filter the data by typing the

Combine with other Genes

Transform into related records

Use Genomic Colocation to combine with other features

Add a step to your search strategy

① Choose how to combine with other Genes
 3 INTERSECT 4 3 UNION 4 3 MINUS 4 4 MINUS 3

② Choose which Genes to combine. From...
 A new search An existing strategy My basket

RNA

Transcriptomics
RNA-Seq Evidence

word liver in the filter box at the top of the page. This should yield two datasets from *P. cynomolgi* and *P. vivax*. For this exercise, select the fold change query for the *P. cynomolgi* dataset: Liver stage hypnozoite vs schizont transcriptomes (primary culture) (Voorverg-van der Wel et al.).

Search for Genes by RNA-Seq Evidence

The results will be intersected with | the results of Step 3.

Filter Data Sets:

Legend: DE Differential Expression FC Fold Change P Percentile SA SenseAntisense

Add a step to your search strategy

Search for Genes by RNA-Seq Evidence

The results will be intersected with | the results of Step 3.

Filter Data Sets: liver

Legend: DE Differential Expression FC Fold Change P Percentile SA SenseAntisense

Choose a Search

DE FC P
FC P SA
DE FC P
FC P SA

Organism	Data Set
<i>Plasmodium berghei</i> ANKA	5 asexual stages
<i>Plasmodium berghei</i> ANKA	P. berghei
<i>Plasmodium berghei</i> ANKA	Female gametocyte infections
<i>Plasmodium chabaudi</i> chabaudi	Trophozoite infections
<i>Plasmodium chabaudi</i> chabaudi	Trophozoites
<i>Plasmodium cynomolgi</i> strain M	Transcriptomes
<i>Plasmodium cynomolgi</i> strain M	Liver stage hypnozoite vs schizont transcriptomes (primary culture) (Voorverg-van der Wel et al.)
<i>Plasmodium cynomolgi</i> strain M	Hypnozoite, schizont and blood stage transcriptomes (laser microdissection) (Cubi et al.)
<i>Plasmodium falciparum</i> 3D7	Gamete Transcriptomes (Lasonder et al.)
<i>Plasmodium falciparum</i> 3D7	Mosquito or cultured sporozoites and blood stage transcriptome (NF54) (Hoffmann et al.)

11. Configure the RNA-Seq search to identify genes that are differentially regulated by at least 2-fold between all the hypozoite stages and the sporozoite stages. For example, select the hypozoite stages in the reference selection box and the sporozoite samples in the comparator selection box, then click on run step.

12. How many results did you get? Why did you get 0 results? How can you change this?

For the Experiment
Liver stage hypozoite vs schizont transcriptomes (primary culture) unstranded

return protein coding Genes
that are up or down regulated

with a Fold change ≥ 2
between each gene's average expression value
(or a Floor of 10 reads)

in the following Reference Samples

sporozoite 6-7 days pi
 sporozoite 9 days pi
 sporozoite 10 days pi
 hypozoite 6-7 days pi
 hypozoite 9 days pi

select all | clear all

and its average expression value
(or the Floor selected above)

in the following Comparison Samples

sporozoite 6-7 days pi
 sporozoite 9 days pi
 sporozoite 10 days pi
 hypozoite 6-7 days pi
 hypozoite 9 days pi

select all | clear all

Example showing one gene that would meet search criteria
(Dots represent this gene's expression values for selected samples)

Up or down regulated

Average Expression Value Comparison

Average Expression Value Reference

Average Expression Value Comparison

Average Expression Value Reference

For each gene, the search calculates:

$$\text{fold change}_{\text{up}} = \frac{\text{average expression value in comparison}}{\text{average expression value in reference}}$$

$$\text{fold change}_{\text{down}} = \frac{\text{average expression value in reference}}{\text{average expression value in comparison}}$$

and returns genes when $\text{fold change}_{\text{up}} \geq 2$ or $\text{fold change}_{\text{down}} \geq 2$.

You are searching for genes that are up or down regulated between at least two reference samples and at least two comparison samples.

Run Step

Remember that the previous search was a list of *P. falciparum* genes and this RNA-Seq was from *P. cynomolgy*. What you would like to do is convert the *P. cynomolgy* genes into *P. falciparum* genes. To do this follow these steps:

- hover your mouse over the RNA-seq step then click on the edit option on that step.
- In the popup window, click on the **orthologs** link.

Unnamed Search Strategy *

Step 1: Signal Pep (44,582 Genes)
Step 2: Transmb Dom (75,167 Genes)
Step 3: SNPs (1,604 Genes)
Step 4: PcyM Liver HvsS RNA-Seq (2,236 Genes)

View | Analyze | Revise | Make nested strategy | Insert step before | Orthologs | Delete

Details for step PcyM Liver HvsS RNA-Seq (fc)
2,236 Genes

Experiment: Liver stage hypozoite vs schizont transcriptomes (primary culture) unstranded
Direction: up or down regulated
Reference Samples: sporozoite 6-7 days pi, sporozoite 9 days pi, sporozoite 10 days pi
Operation Applied to Reference Samples: average
Comparison Samples: hypozoite 6-7 days pi, hypozoite 9 days pi
Operation Applied to Comparison Samples: average
fold difference >= 2
Floor = 10 reads
Protein Coding Only: protein coding

Give this search a weight

- c. In the next window select which organism(s) you would like to transform to. For this exercise select *P. falciparum* 3D7 and click on run step.
- d. Did you get results now?

Organism

Note: You must select at least 1 values for this parameter.
1 selected, out of 45

[add these](#) | [clear these](#) | [select only these](#)
[select all](#) | [clear all](#)

3d7 (arrow)

Plasmodium falciparum (arrow)

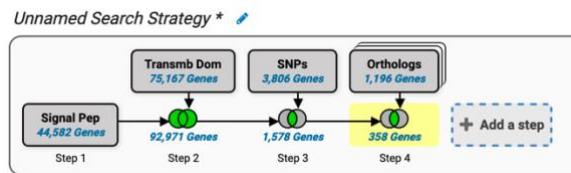
[add these](#) | [clear these](#) | [select only these](#)
[select all](#) | [clear all](#)

Syntenic Orthologs Only?

no (arrow)

Run Step

13. Next identify how many of these genes do not have orthologs in mammals. To do this



add a step for genes based on orthology phylogenetic profile. Again you can filter the searches by typing the word “phylogenetic”.

On the next page select *P. falciparum* 3D7 the configure the phylogenetic profile by

Add a step to your search strategy ?

Combine with other Genes

① Choose how to combine with other Genes

4 INTERSECT 5 4 UNION 5 4 MINUS 5 5 MINUS 4

② Choose which Genes to combine. From...

A new search An existing strategy My basket

phy (arrow)

Orthology and synteny
Orthology Phylogenetic Profile

Transform into related records

Use Genomic Colocation to combine with other features

finding Mammalia under Chordata which are under Metazoa. Click twice on the circle next to Mammalia – it should become a red x (See image below).

The screenshot shows the 'Add a step to your search strategy' interface. In the search bar, '3d7' is typed, with a red arrow pointing to it. Below the search bar, there's a list of selected organisms: Plasmodium falciparum and Plasmodium falciparum 3D7. Underneath this, there's a section titled 'Select orthology profile' with a detailed tree view of taxonomic groups. A red arrow points to the 'Mammalia (MAMM)' node under Chordata. At the bottom of the interface, there are several buttons for combining genes, transforming records, and using genomic colocation.

14. Determine if a mutation in any of these genes affects fitness. Click on add step and find the search for phenotype evidence.

This screenshot shows the 'Add a step to your search strategy' interface. On the right, there's a search bar with 'phen' typed in, and a red arrow points to it. Below the search bar, there are three tabs: 'Choose how to combine with other Genes', 'Choose which Genes to combine. From...', and 'Choose which Genes to combine. To...'. The 'From...' tab is selected. On the left, there are three boxes showing workflows for 'Combine with other Genes', 'Transform into related records', and 'Use Genomic Colocation to combine with other features'. Each workflow has a 'Step 5' and 'Step 6' button.

15. Select the P. falciparum piggyBac insertion mutagenesis (John Adams) experiment.

This screenshot shows the 'Search for Genes by Phenotype Evidence' interface. At the top, it says 'The results will be intersected with the results of Step 5.' Below this, there's a legend with four categories: 'Association to Genomic Segments' (red), 'Curated Phenotype' (blue), 'Similarity' (green), and 'Similarity of Association' (yellow). A red circle highlights the 'Curated Phenotype' button. The main table lists various datasets, and the last row, 'piggyBac insertion mutagenesis (John Adams)', also has its 'Curated Phenotype' button highlighted with a red circle.

16. On the next page select the Mutant Fitness Score (MFS) option and choose any score range – generally the more negative the bigger the effect is on fitness. For this example a score range of -4.078 to -3.07 was chosen.



Explore your final results. Do they make sense/plausible? Note that you can revise any of the steps in the strategy to explore the data further. You can also save your strategy and share it with others or make it public. Here is a link to this search strategy:

<https://plasmodb.org/plasmo/app/workspace/strategies/import/fd387e8d3acda856>

Public Strategies

Users can share their strategies publicly so that others may use them. The public strategies link is located under the **About** menu followed by **Community** and **Public Strategy**. A table of available public strategies will appear and there is a filter box located at the top of the public strategies so that you can search for the author or subject of the strategy among other items. The public strategies for PlasmoDB, as an example, are located at:

<https://plasmodb.org/plasmo/app/workspace/strategies/public>

Regular Expressions & Genomic Colocation

Note: this exercise uses different VEuPathDB resources as an example database, but the same functionality is available on all VEuPathDB resources.

Learning objectives:

- Run a regular expression search on amino acid sequences
- Run a regular expression search on nucleotide sequences
- Use the genomic colocation search

Protein or nucleotide sequences can be identified using the regular expression searches in VEuPathDB. This search is very useful to identify patterns of sequences.

Searches can be accessed from categorized menus in the left search for panel (A) or from the searches menu in the header (B).

Accessing the protein motif pattern search:

- Click on the *Genes* category then click on the sequence analysis category

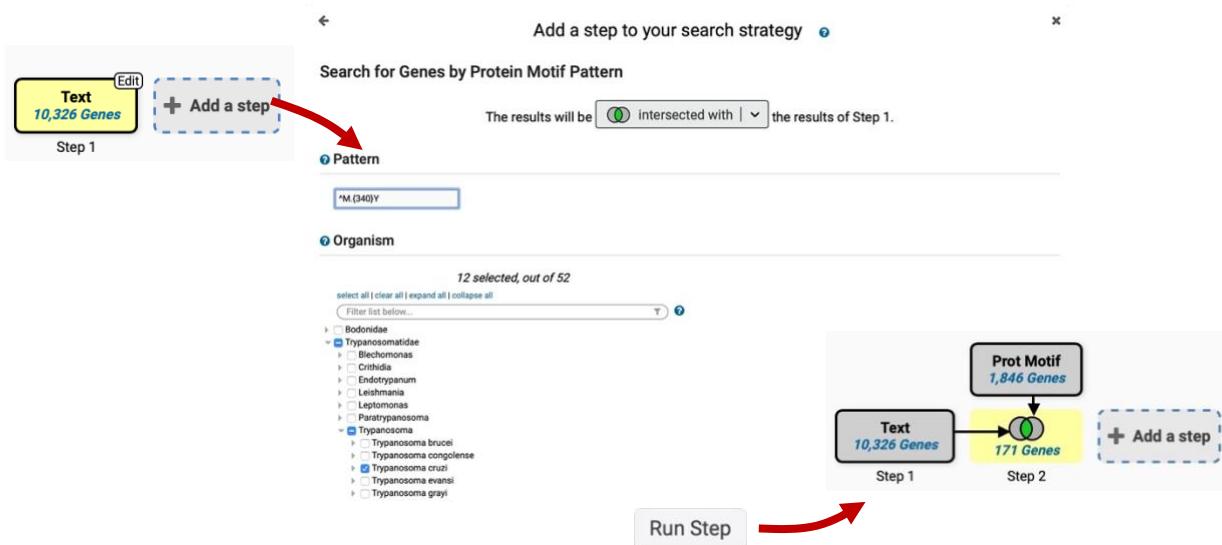
Accessing the DNA motif pattern search:

- Click on the *Genomic segments* category

Note: the appendix at the end of this document includes additional regular expression help.

Using regular expressions to find motifs in TriTrypDB: finding active trans-sialidases in *T. cruzi* (TriTrypDB).

- T. cruzi* has an expanded family of trans-sialidases. In fact, if you run a text search for any gene with the word “trans-sialidase” in its *product description*, you return over 10000 genes among the strains in the database!!! Try this and see what you get.
- Not all of the genes returned in (a) are predicted to be active. It is known that active trans-sialidases have a signature tyrosine (Y) at position 342 in their amino acid sequence. Add a motif search step to the text search in ‘a’ to identify only the active trans-sialidases.
- Write a regular expression that defines a protein sequence that starts with a methionine, and is followed by 340 of any amino acids, followed by a tyrosine ‘Y’.



Find Cryptosporidium genes with the YXXΦ receptor signal motif (CryptoDB)

The YXXΦ (Y=tyrosine, X=any amino acid, Φ=bulky hydrophobic [phenylalanine, tyrosine, threonine]) motif is conserved in many eukaryotic membrane proteins that are recognized by adaptor proteins for sorting in the endosomal/lysosomal pathway. This motif is typically located in the c-terminal

end of the protein. *****Note:** do not look for the ϕ symbol on your keyboard – this will not work. Rather you should use the amino acid symbols.

- d. Use the “protein motif pattern” search to find all *Cryptosporidium* proteins that contain this motif anywhere in the terminal 10 amino acids of proteins. (hint: for your regular expression, remember that you want the first amino acid to be a tyrosine, followed by any two amino acids, followed by any bulky hydrophobic amino acid (phenylalanine, tyrosine, threonine).
- e. How many of these proteins also contain at least one transmembrane domain.

Identify Genes based on Protein Motif Pattern

Pattern

Organism

11 selected, out of 14

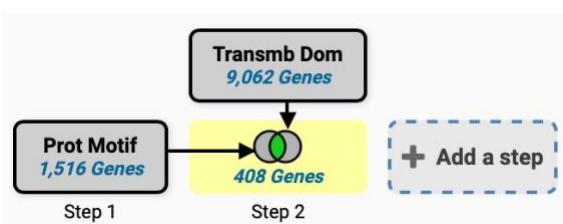
select all | clear all | expand all | collapse all

Filter list below...

- Apicomplexa
 - Coccidia
 - Euoccidiidae
 - Cryptosporidiidae
 - Cryptosporidium
 - Gregarinina
 - Chromerida

select all | clear all | expand all | collapse all

Get Answer



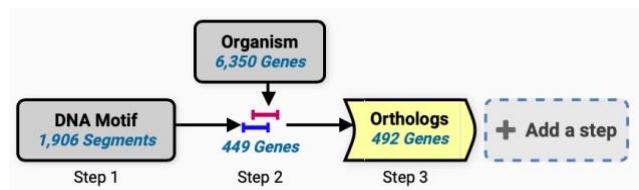
- f. What would happen if you revise the first step (the motif pattern step) to include genes with the sorting motif in the C-terminal 20 amino acids? (hint: edit the first step and modify your regular expression).

[Note: if you need help with the regular expression the answers are in appendix B.](#)

Find fungal genes downstream of a regulatory DNA motif (FungiDB).

Transcriptional start sites are often located within a certain distance upstream of the genes or gene clusters that they regulate. In fungi, DNA motifs are also important for regulation of processes linked to host cell invasion or production of secondary metabolites. Readily available genomic data facilitate the discovery of regulatory motifs via examination of orthologous sequences.

The goal of this exercise is to identify all genes harboring upstream CACGTG motif, known for its role in transcriptional regulation. We will start our search in an extensively studied model organism *Saccharomyces cerevisiae* and expand our search to *Fusarium graminearum*.



Here is a summary of the search strategy:

Find the CACGTG DNA motif in the *Saccharomyces cerevisiae* genome.

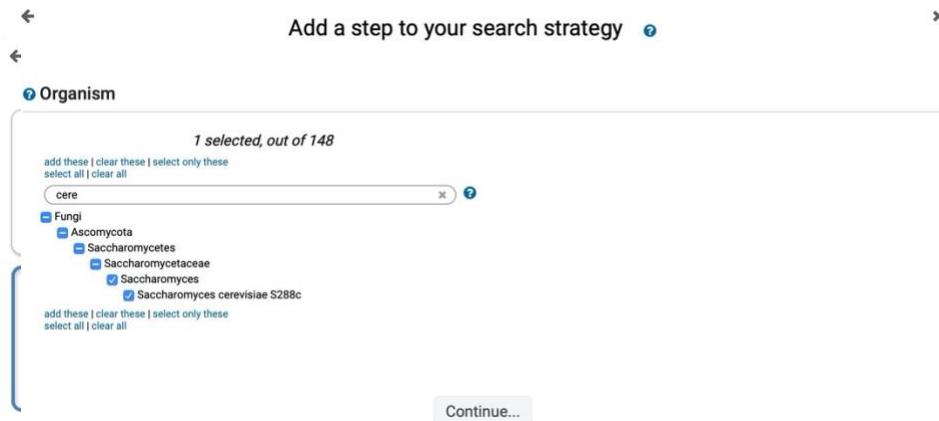
1. Select the “Search for genomic segments (DNA motif)” menu from the Search menu and look for CACGTG in *S. cerevisiae*.
2. Your search returns over 1900 DNA segments containing GACGTG motif. Next, let’s look for putative regulatory targets of this motif by searching for genes that are located 600bp downstream of this sequence.

The screenshot shows the "Identify Genomic Segments based on DNA Motif Pattern" search interface. On the left, a sidebar titled "Search for..." lists various search categories: Genes, Organisms, Popset Isolate Sequences, Genomic Sequences, Genomic Segments (selected), DNA Motif Pattern (selected), Genomic Location, SNPs, ESTs, Metabolic Pathways, and Compounds. A blue arrow points from the "Genomic Segments" and "DNA Motif Pattern" items in the sidebar to the corresponding fields in the main search form. The main search form has sections for "Organism" and "Pattern". In the "Organism" section, a dropdown menu is open, showing a tree structure of taxonomic ranks. The path selected is "cer" under "Fungi" > "Ascomycota" > "Eurotiomycetes" > "Onygenales" > "Onygenaceae" > "Byssomyces" > "Byssomyces ceratinophila" > "Byssomyces ceratinophila isolate UAMH 5669". The "Pattern" section contains a text input field with "CACGTG" typed into it. A yellow box highlights the "Filter to find organism" button above the dropdown and the "Type sequence pattern" input field below it. A "Get Answer" button is located at the bottom right of the search form.

Identify genes with the CACGTG motif located 600bp upstream of an open reading frame.

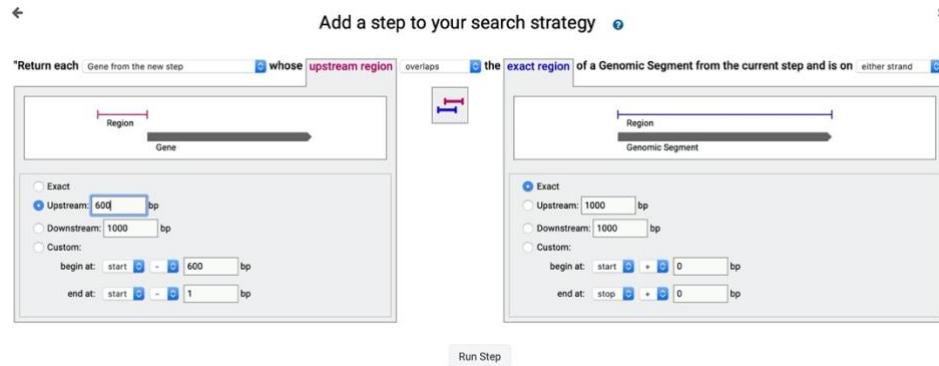
EuPathDB offers a colocation function to identify genomic features within a specified distance of each other. Run a search for all genes in *Saccharomyces cerevisiae* and use the colocation tool to identify genes that contain the CACGTG motif in their upstream regions. Follow these steps:

3. Click “Add Step”. Choose the option on the left called “Use Genomic Colocation to combine with other features” then select the



organism gene search which can be found under the *Taxonomy* category. On the next page select *Saccharomyces cerevisiae* from the taxonomy browser and click on continue.

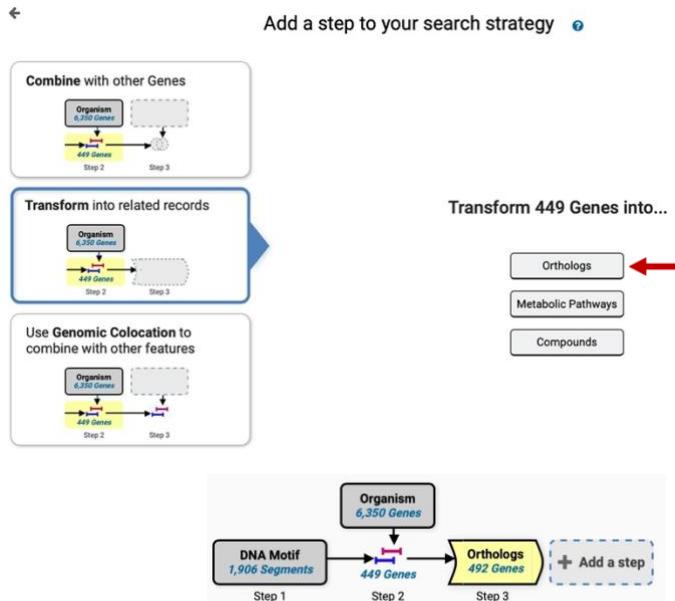
4. Configure the parameters on the next page to return each gene from step 2 whose upstream region (600bp) overlaps the exact region of a Genomic Segment in Step1 (CACGTG) and is on either strand.



g. Identify orthologs *S. cerevisiae* genes in *Fusarium graminearum*.

All VEuPathDB sites offer tools to transform results between record types. The “Transform by Orthology” tool uses orthology clusters assigned by the OrthoMCL algorithm to enable transformation of a list of genes from one or more species to another (or more) species.

- click on add step then select the **Transform** into related records option from the left side of the popup. Next click on the *Orthologs*



option.

- Select *F. graminearum* in the next popup window and click on *Run Step*.
-

Appendix A

Regular expression help

The following codes can be used to represent classes of amino acids.

AA property	Amino acids	Code
Acidic	DE	0
Alcohol	ST	1
Aliphatic	ILV	2
Aromatic	FHWY	3
Basic	KRH	4
Charged	DEHKR	5
Hydrophobic	AVILMFYW	6
Hydrophilic	KRHDENQ	7
Polar	CDEHKKNQRST	8
Small	ACDGNPSTV	9
Tiny	AGS	B
Turnlike	ACDEGHKNQRST	Z
	ACDEFGHIKLM	
Any	NPQRSTVWY	X

The following is a simple explanation of regular expressions.

Perl regular expressions are terms used for pattern matching in text strings, e.g. '**aadgt**', '**aa+dgt**', '**a/d/c**', '**[mac]a**'.

Because nucleotide and amino acid sequences are text strings, regular expressions are very useful for finding motifs within sequences.

Motifs often include repetitive or ambiguous assignments at some locations. The rules and special characters used in regular expressions help define the full set of strings that match the motif pattern.

The following is a description of some of these characters and examples of how they are used.

Although regular expressions seem complicated at first, they are very useful and easy to understand after going through some examples.

Special Characters

- . Match any character.
- + Matches "one or more of the preceding characters".
- * Matches "any number of occurrences of the preceding character", including 0.
- ? Matches "zero or one occurrences of the preceding character".
- [] Matches any character contained in the brackets.
- [^] Match any character *except* those in the brackets.
- {n} Matches when the preceding character, or character range, occurs exactly n times.
- {n,} Matches when the preceding character occurs at least n times.
- {n,m} Matches when the preceding character occurs at least n times, but no more than m times.

Here are some examples of searches.

ad+f (1 or more occurrences of 'd') would match any of the following:

adf
addf
addd
f
add
dddf

...

ad*f (0 or more occurrences of 'd') would match:

a
f
a
d
f
a
d
d
f
adddf

...

ad?f (0 or 1 occurrence of 'd') would match:

a
f
a
d
f

a[yst]c would match:

a
t
c
a
s
c
a
y
c

Specify the number of occurrences of a residue.

P{1,5} would match P from 1 to 5 times.

.{1,30} would match any amino acid 1 to 30 times so you could find a motif within 30 amino acids of something like the beginning.

Pattern Anchors

- ^ Match only at the beginning of the string.
- \$ Match only at the end of the string.

Here are examples of expressions using pattern anchors.

^mdef (e.g. a protein sequence **starting with** 'mdef') would match:

- mdef
- mdefab
- mdefared

fadfk **but**

not match:

- edefa
- emdefa
- eeeemdef

kdel\$ (searches for proteins **ending with** 'kdel', a standard ER retention signal) would match:

- eeeekdel
- kdel

but not match :

- edefkdell
 - akdeleef
- Appendix B

Answers to exercise 2:

A: YXXΦ in the terminal 10 amino acids à ReEX = Y..[FTY].{0,6}\$

B: YXXΦ in the terminal 20 amino acids à ReEX = Y..[FTY].{0,16}\$

Variant calling in VEuPathDB galaxy (Part 1)

Learning objectives:

1. Retrieve DNA sequence data from the sequence repository EBI and upload data to VEuPathDB Galaxy using Globus Data Transfer;
2. Name a new project/history;
3. Deploy a Variant calling workflow in the VEuPathDB Galaxy.

Galaxy is an open, web-based platform for data-intensive biomedical research. Galaxy allows you to perform, reproduce, and share complete analyses without the use of command-line scripting. VEuPathDB developed its Galaxy instance in collaboration with Globus Genomics (VEuPathDB Galaxy). To learn how to use Galaxy, follow this link to access tutorials prepared by the Galaxy Training Network:

https://wiki.galaxyproject.org/Learn#Galaxy_101

There are different ways to get data into Galaxy. In this exercise we will use Globus Data Transfer to get data from the EBI server using a unique project ID.

Retrieve DNA sequence data from the sequence repository and upload data to VEuPathDB Galaxy using Globus Data Transfer option.

- a. Click on the “Globus Data Transfer” menu on the left to expand the Data Transfer section.
- b. Click on the “Get Data via Globus from the EBI server” link.

The screenshot shows the VEuPathDB Galaxy Site interface. On the left, there is a sidebar with a 'globus Genomics' logo and a search bar. Below the search bar, under 'VEUPATHDB APPLICATIONS', are links for 'VEuPathDB Export Tools', 'VEuPathDB OrthoMCL Tools', and 'VEuPathDB RNA-Seq Tools'. Under 'DATA TRANSFER', there is a red box around the 'Globus Data Transfer' link, which has a red arrow pointing to it. Other options listed under 'DATA TRANSFER' include 'Get Data via Globus High speed file upload', 'Get Flowcell sample FastQ per lane via Globus Transfer FASTQ from Globus to Galaxy', 'Get Data via Globus from the EBI server using your unique file identifier', 'Get Data with BioProject ID from the EBI server using SRA ID', 'Get Data via Globus from the EBI server (collections) using your unique file identifier', and 'Get BDBag from MINID to collection transfer data given a MINID to a collection'. The main content area features a 'Welcome to the VEuPathDB Galaxy Site' header and a sub-section titled 'With VEuPathDB Galaxy you can:' with a numbered list of six items. At the bottom, there are sections for 'Get started with VEuPathDB pre-configured workflows:' and 'OrthoMCL'.

Welcome to the VEuPathDB Galaxy Site
A free, interactive, web-based platform for large-scale data analysis

With VEuPathDB Galaxy you can:

1. Start analyzing your data now with pre-configured workflows. All VEuPathDB genomes are pre-loaded.
2. Perform large-scale data analysis with no prior programming or bioinformatics experience.
3. Create custom workflows using an interactive workflow editor. [Learn how](#)
4. Export your results to VEuPathDB, so that you can explore your data with our tools, such as JBrowse and search strategies. See [this tutorial](#).
5. View your results on Galaxy or download results to your computer.
6. Keep data private, or share data with colleagues or the community.

To learn more about Galaxy, visit the [public Galaxy resources](#).

Get started with VEuPathDB pre-configured workflows:

OrthoMCL

This workflow uses BLASTP and the OrthoMCL algorithm to assign your set of proteins to OrthoMCL groups. Version OG6r1 is the latest set of groups (as of April 2020), but you can also select the previous set (OG5). [Explore this tutorial to learn more.](#)

- Workflow to map your proteins to OrthoMCL groups

- c. Enter ENA sample ID and define the dataset type to be transferred into the VEuPathDB Galaxy workspace.

The ENA ID should start with the letters ‘SAM’. For this exercise, we will use SAMN01815907, which is a paired-ended dataset. Take care to specify

whether a dataset is a single or paired-ended as incorrect selection will cause the upload to fail.

- d. Once the form is properly filled, click on the “Execute” button to start the data transfer process.

- e. When the job has been successfully deployed and added to the queue, the screen will refresh, and the added job will appear in the history on the right.

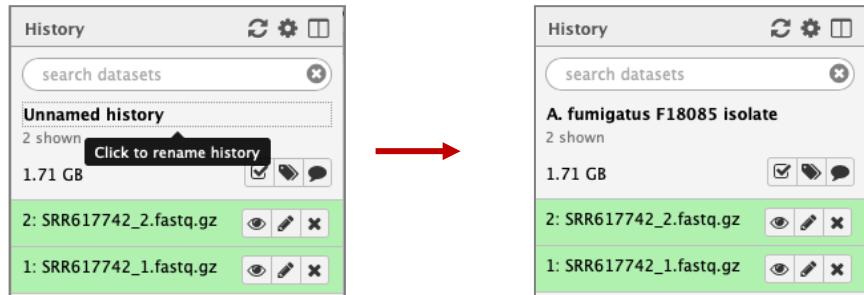
Note: new jobs are highlighted in grey, in progress – yellow, and those completed are in green.

Notice that there are two files appearing in the history on the right. This is because the uploaded data is paired-ended.

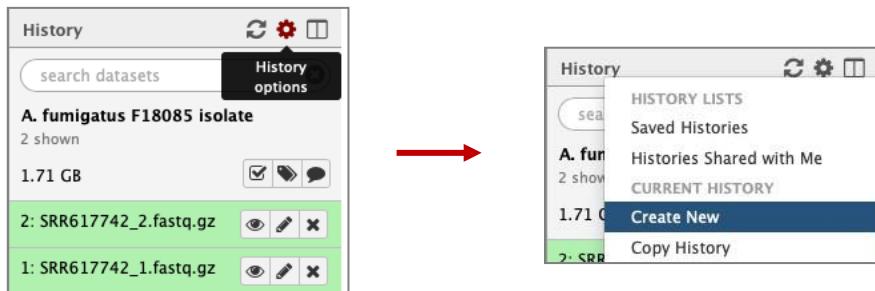
Rename your history.

By default, all new jobs will be added to the current history on the right. Unless renamed, the history will show up in your history as “Unnamed history”. Let’s rename the history to help us track this project in the future.

- f. Click on the “Unnamed history” and type “A. fumigatus F18085 isolate”, and then press “enter” to rename this history.



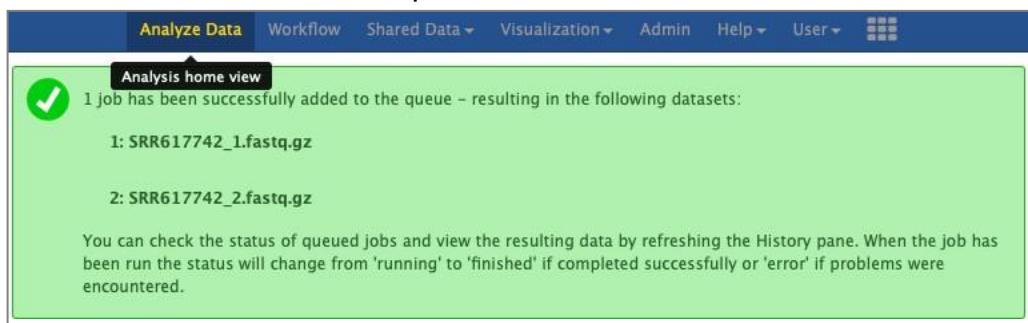
Note: if you would like to start a new project/history, click on the wheel button at the top of the history section and select “Create a new history”.



Deploy a Variant calling workflow.

VEuPathDB Galaxy main landing page has several workflows for variant calling.

- g. To navigate to the main page, click on the “Analyze Data”, which is located in the main menu at the top.



- h. Scroll down to the Variant calling section and choose the workflow for paired-end reads.

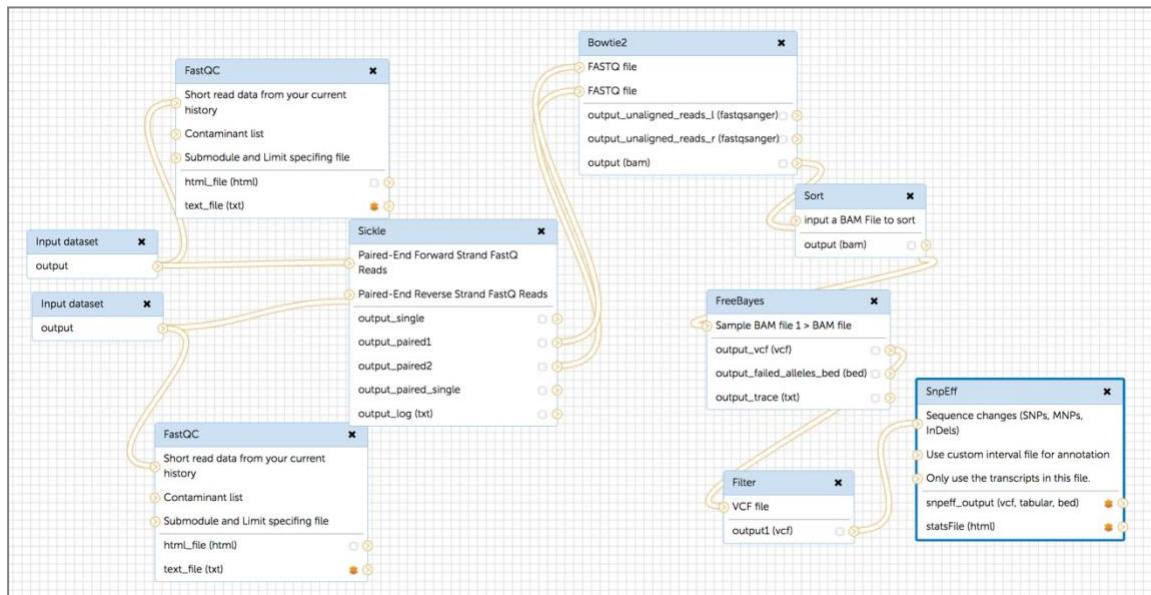
Variant calling

Use the following workflows to analyze your FASTQ files. The workflows use Sickle for preparation of reads, Bowtie2 for mapping reads to a VEuPathDB reference genome, Freebayes for variant detection, SnpEff to evaluate the effect of variants, and SnpSift for filtering types of variants. Choose the appropriate workflow based on your input data. A VCF file is generated that can be analyzed in Galaxy or downloaded to your computer. NOTE: Export of VCF files to VEuPathDB will be available soon.

- [Workflow for single-end reads](#)
- [Workflow for paired-end reads](#)

The pre-configured variant calling workflows include the following steps:

- Determine quality of the reads and generate reports (FastQC);
- Trim reads based on their quality scores (Sickle);
- Align reads to a reference genome using Bowtie2 and generate coverage plots ;
- Sort alignments with respect to their chromosomal positions (Sort);
- Detect variants (FreeBayes);
- Filter SNP candidates (Filter);
- Analyze and annotate variants, and calculate the effects of SNPs via SnpEff.



- i. Click on the workflow for paired-end reads and set workflow parameters.
 - Make sure that the input steps for paired-end data are set to the xxxx_1.fastq.gz and xxxx_2.fastq.gz file (by default the same file will be selected in both fileds).

The screenshot shows the 'Workflow: imported: EuPathDB_Workshop_VariantCalling_PairedEnd' interface. In the 'History Options' section, there is a 'Send results to a new history' button with 'Yes' and 'No' options. Below it, two input dataset sections are shown: '1: Input dataset - 1' containing '1: SRR617742_1.fastq.gz' and '2: Input dataset - 8' containing '2: SRR617742_2.fastq.gz'. To the right, a 'History' panel displays a dataset named 'A. fumigatus F18085 isolate' with a size of 1.71 GB. The 'Run workflow' button is located at the top right of the main workflow area.

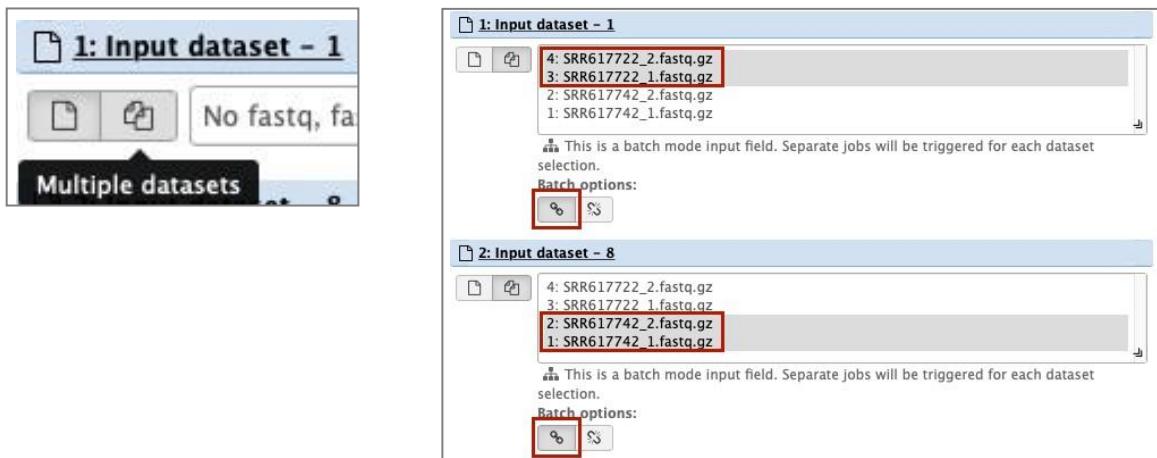
- Select the correct reference genome.
 - Select *Aspergillus fumigatus* Af293 as a reference genome (steps: Bowtie2, FreeBayes, SnpEff).

The screenshot shows the 'FreeBayes – Bayesian genetic variant detector' configuration interface. It includes a 'Choose the source for the reference list' section with 'Locally cached' and 'Sample BAM file' options. Under 'Sample BAM file', '1: Sample BAM file' is selected, showing 'BAM file' and 'Output dataset 'output' from step 7'. In the 'Using reference genome' section, 'FungiDB-29_AfumigatusAf293_Genome' is selected. The 'Run workflow' button is located at the top right of the main workflow area.

- Choose to deploy the analysis within the same history and click on the Run workflow button.

The screenshot shows the final configuration of the 'Workflow: imported: EuPathDB_Workshop_VariantCalling_PairedEnd'. The 'History Options' section has a 'Send results to a new history' button with 'Yes' and 'No' options, where 'No' is highlighted with a red box. The 'Run workflow' button is located at the top right of the main workflow area.

Note: You can use the same workflow to analyze multiple samples in batches. The upload steps remain the same, however, when setting up the workflow, click on multiple dataset button within the input dataset section.



Related sites of interest to our communities

- [Previous EuPathDB Workshops](#)
- [Companion](#)
- [OrthoMCL](#)
- [GeneDB](#)
- [ModBase at UCSF](#)
- [Tetrahymena Genome](#)
- [The Arabidopsis Information Resource](#)
- [NAR Database Summary Paper Categories](#)

VEuPathDB Publications and Citations

This [PubMed filter](#) provides a current list of the most recent publications about VEuPathDB resources

Recent publications include:

Omar Harb, Jessica C. Kissinger and David S. Roos on behalf of the EuPathDB group

ToxoDB: the functional genomic resource for *Toxoplasma*

Toxoplasma gondii 3rd edition, Edited by Louis Weiss and Kami Kim (2020)

<https://doi.org/10.1016/B978-0-12-815041-2.00023-2>

Susanne Warrenfeltz and Jessica Kissinger on behalf of the EuPathDB Consortium

[**Accessing *Cryptosporidium* omic & isolate data via CryptoDB.org**](#)

Methods in Molecular Biology, Editor, Jan Mead (2020)

https://doi.org/10.1007/978-1-4939-9748-0_22

Omar S. Harb and David S. Roos on behalf of the EuPathDB Consortium

ToxoDB: Functional Genomics Resource for *Toxoplasma* and Related Organisms

Methods in Molecular Biology, Editor, Christopher J. Tonkin (2020)

https://doi.org/10.1007/978-1-4939-9857-9_2

Susanne Warrenfeltz, Evelina Y Basenko, Kathryn Crouch, Omar S. Harb, Jessica C. Kissinger, Achchuthan Shanmugasundram and Fatima Silva-Franco

EuPathDB: the eukaryotic pathogen genomics database resource

Methods in Molecular Biology, Editor, Martin Kollmar (2018)

https://doi.org/10.1007/978-1-4939-7737-6_5

Evelina Y. Basenko, Jane A. Pulman, Achchuthan Shanmugasundram, Omar S. Harb, Kathryn Crouch, David Starns, Susanne Warrenfeltz, Cristina Aurrecoechea, Christian J. Stoeckert, Jr., Jessica C. Kissinger, David S. Roos and Christiane Hertz-Fowler

FungiDB: An Integrated Bioinformatic Resource for Fungi and Oomycetes. (2018)

J. Fungi 2018, 4(1), 39

<https://doi.org/10.3390/jof4010039>

Cristina Aurrecoechea; Ana Barreto; Evelina Y. Basenko; John Brestelli; Brian P. Brunk; Shon Cade; Kathryn Crouch; Ryan Doherty; Dave Falke; Steve Fischer; Bindu Gajria; Omar S. Harb; Mark Heiges; Christiane Hertz-Fowler; Sufen Hu; John Iodice; Jessica C. Kissinger; Cris Lawrence; Wei Li; Deborah F. Pinney; Jane A. Pulman; David S. Roos; Achchuthan Shanmugasundram; Fatima Silva-Franco; Sascha Steinbiss; Christian J. Stoeckert Jr; Drew Spruill; Haiming Wang; Susanne Warrenfeltz; Jie Zheng

EuPathDB: the eukaryotic pathogen genomics database resource

[Nucleic Acids Research 2017 doi: 10.1093/nar/gkw1105](https://doi.org/10.1093/nar/gkw1105)

Warren, A. S., Aurrecoechea, C., Brunk, B., Desai, P., Emrich, S., Giraldo-Calderón, G. I., Harb, O., Hix, D., Lawson, D., Machi, D., Mao, C., McClelland, M., Nordberg, E., Shukla, M., Vosshall, L. B., Wattam, A. R., Will, R., Yoo, H. S., & Sobral, B.

RNA-Rocket: an RNA-Seq analysis resource for infectious disease research

Bioinformatics, 1 May 2015; 31(9), 1496–1498

<https://doi.org/10.1093/bioinformatics/btv002>

Giraldo-Calderón GI, Emrich SJ, MacCallum RM, Maslen G, Dialynas E, Topalis P, Ho N, Gesing S; VectorBase Consortium, Madey G, Collins FH, Lawson D.

VectorBase: an updated bioinformatics resource for invertebrate vectors and other organisms related with human diseases

Nucleic Acids Res. 2015 Jan;43(Database issue):D707-13

<https://doi.org/10.1093/nar/gku1117>

Publications that use our resource

Our project resources are cited in >20,000 publications. You can view the [publications that cite us](#) in Google Scholar.

Release Notes

Release notes are generated for each project site within VEuPathDB with each new release. They are located in the “New” tab on the right-hand expandable panel of each project site. An example for VectorBase is shown below.

Example release notes - VectorBase 48 Released

(27 Aug 2020)

We are pleased to announce the release of VectorBase 48.

Beginning with VectorBase 48, the URL <https://VectorBase.org> directs you to the newest version of our interface. The site contains previous VectorBase.org data plus some new tools and functions in an updated and streamlined interface that emphasizes easy access to help information. Navigate to <https://legacy.VectorBase.org> for the previous version of the interface. The legacy interface will remain active for at least one release.

Release 48 Webinar

Join us on September 3, 2020 at 10am EDT for our Release 48 Highlights webinar where we will demonstrate and discuss new data and features in VEuPathDB 48. [Register here](#)

New features in VectorBase 48

- Gene pages: The sequence section now contains links to copy the sequence to your clipboard. Sequences are copied in FASTA format with the gene ID in the def line.
- Gene pages: It is now possible to scroll and zoom within the JBrowse views that are present in the 'Gene models' and 'Protein features and properties' sections.
- Contact Us form: A new option on the Contact Us form makes it possible to paste a screenshot into the form when sending us questions or suggestions.
- Search result page: The Download and Add to Basket options from a search result tab now appear as buttons for better visibility.
- Error messages: The text in our error messages was reviewed and edited to improve clarity. And we've made it easier to report errors by adding links in error messages to the Contact Us form.

Omics data sets

- Fifteen additional microarray datasets are available in the merged VectorBase site. These are marked as 'New' in the [Identify Genes based on Microarray Evidence](#) search page.
- The results of nine comparative genomic study analyses (VCF files) are available in [JBrowse](#). The data are categorized under Genetic Variation in the Select Tracks tool.

Population Biology

This release represents our second largest abundance data release ever – over 219,000 new non-zero abundance records. It is also our largest increase in pathogen status assays, doubling our records to over 174,000 assays.

45 projects have been added. You can use any of the Project IDs from the list below, in the Search box of the MapVEu (formerly PopBio map) tool: <https://vectorbase.org/poppbio-map/web/>

- Indoor Human landing catches from Western Cameroon (VBP0000626: [map](#), 52 collections, 166 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Uganda by PRISM ICEMR group, 2018 (VBP0000627: [map](#), 1113 collections, 7791 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Marion County, Indiana, USA. 2019. (VBP0000628: [map](#), 1758 collections, 7539 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Manatee County, Florida, USA. 2017 (VBP0000629: [map](#), 1845 collections, 10390 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Manatee County, Florida, USA. 2018 (VBP0000630: [map](#), 1217 collections, 8462 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Manatee County, Florida, USA. 2019 (VBP0000631: [map](#), 1859 collections, 11223 samples, 0 phenotypes, 0 genotypes)
- Maryland Dept. of Agriculture mosquito surveillance, Maryland, USA, 2016 (VBP0000632: [map](#), 558 collections, 4373 samples, 0 phenotypes, 0 genotypes)
- Maryland Dept. of Agriculture mosquito surveillance, Maryland, USA, 2019 (VBP0000634: [map](#), 757 collections, 7016 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Lucas County, Ohio, USA. 2018 (VBP0000635: [map](#), 1438 collections, 6626 samples, 0 phenotypes, 0 genotypes)
- Supplemental Maryland Dept. of Agriculture mosquito surveillance, Maryland, USA, 2016 (VBP0000636: [map](#), 45 collections, 259 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Lucas County, Ohio, USA. 2019 (VBP0000637: [map](#), 1554 collections, 7260 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Hernando County, Florida, USA. 2018 (VBP0000638: [map](#), 361 collections, 2167 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Hernando County, Florida, USA. 2019 (VBP0000639: [map](#), 513 collections, 2656 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2016 (VBP0000664: [map](#), 1750 collections, 4745 samples, 0 phenotypes, 0 genotypes)
- Dynamic of resistance alleles of two major insecticide targets in *Anopheles gambiae* (s.l.) populations from Benin, West Africa (VBP0000640: [map](#), 28 collections, 28 samples, 54 phenotypes, 0 genotypes)
- Contrasting resistance patterns to type I and II pyrethroids in two major arbovirus vectors *Aedes aegypti* and *Aedes albopictus* in the Republic of the Congo, Central Africa (VBP0000641: [map](#), 7 collections, 17 samples, 45 phenotypes, 12 genotypes)
- Maryland Dept. of Agriculture mosquito surveillance, Maryland, USA, 2017 (VBP0000633: [map](#), 614 collections, 5500 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Norfolk County Mosquito Control, in Massachusetts, USA, 2019 (VBP0000644: [map](#), 934 collections, 4297 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in Desplaines Valley Mosquito Abatement, Illinois, USA, 2019 (VBP0000645: [map](#), 3006 collections, 25974 samples, 0 phenotypes, 0 genotypes)

- Mosquito surveillance in Salt Lake County, Utah, USA. 2018 (VBP0000646: [map](#), 1585 collections, 6134 samples, 0 phenotypes, 0 genotypes)
- Mosquito infection assays for Plasmodium falciparum in Uganda by PRISM ICEMR group, 2015-2017 (VBP0000647: [map](#), 7621 collections, 92524 samples, 92524 phenotypes, 0 genotypes)
- Mosquito infection assays for Plasmodium falciparum in Uganda by PRISM ICEMR group, 2018 (VBP0000648: [map](#), 379 collections, 1497 samples, 1497 phenotypes, 0 genotypes)
- Mosquito surveillance in Uganda by PRISM ICEMR group, 2017 - 2019 (VBP0000649: [map](#), 2706 collections, 6791 samples, 0 phenotypes, 0 genotypes)
- Susceptibility of *An. gambiae* s.l. from Côte d'Ivoire to Insecticides used on Insecticide-Treated Nets: Evaluating the Additional Entomological Impact of Piperonyl Butoxide and Chlорfenapyr (VBP0000650: [map](#), 15 collections, 165 samples, 165 phenotypes, 0 genotypes)
- Mosquito abundance in Central Mozambique (VBP0000651: [map](#), 28 collections, 55 samples, 0 phenotypes, 0 genotypes)
- Spatial and temporal distribution of *Anopheles arabiensis* larvae (Sudan) (VBP0000652: [map](#), 3293 collections, 3296 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2017 (VBP0000653: [map](#), 1623 collections, 5221 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2019 (VBP0000654: [map](#), 1376 collections, 4778 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2007 (VBP0000655: [map](#), 602 collections, 1408 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2008 (VBP0000656: [map](#), 692 collections, 1627 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2009 (VBP0000657: [map](#), 831 collections, 2389 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2011 (VBP0000659: [map](#), 1064 collections, 3050 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2012 (VBP0000660: [map](#), 1482 collections, 3972 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2013 (VBP0000661: [map](#), 1819 collections, 5103 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2014 (VBP0000662: [map](#), 1488 collections, 5151 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2015 (VBP0000663: [map](#), 1625 collections, 4856 samples, 0 phenotypes, 0 genotypes)
- Montana State University mosquito data for 2015 (VBP0000666: [map](#), 223 collections, 1458 samples, 0 phenotypes, 0 genotypes)
- Montana State University mosquito data for 2018 (VBP0000667: [map](#), 239 collections, 1640 samples, 0 phenotypes, 0 genotypes)
- Montana State University mosquito data for 2016 (VBP0000668: [map](#), 283 collections, 1790 samples, 0 phenotypes, 0 genotypes)
- Montana State University mosquito data for 2019 (VBP0000669: [map](#), 80 collections, 639 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance in South Walton County, Florida, USA. 2019 (VBP0000670: [map](#), 856 collections, 5264 samples, 0 phenotypes, 0 genotypes)
- Mosquito surveillance, Ada County Weed, Pest, and Mosquito, in Idaho, USA, 2010 (VBP0000658: [map](#), 798 collections, 1976 samples, 0 phenotypes, 0 genotypes)
- Larval habitats seasonality and species distribution (VBP0000671: [map](#), 4278 collections, 8866 samples, 0 phenotypes, 0 genotypes)

- Comparative toxicity of larvicides and growth inhibitors on Aedes aegypti from select areas in Jamaica (VBP0000642: [map](#), 6 collections, 238 samples, 238 phenotypes, 0 genotypes)
- Variation in Malaria Transmission Dynamics in Three Different Sites in Western Kenya (VBP0000665: [map](#), 133 collections, 1701 samples, 0 phenotypes, 0 genotypes)

Note: Our bimonthly database releases incorporate new data and correct errors in old data when necessary. Changes in annotation and new experimental data may slightly alter your search results by increasing or decreasing the number of hits. When search parameters change with a new release, we invalidate (\emptyset) the search and ask you to rerun it. When IDs are updated or removed, we map the old IDs to the new ones, remove the old IDs from your Basket, and leave your Favorites page alone.

Analyses methods:

VEuPathDB draws data from many sources. To facilitate comparisons across data sets, we analyze all data with standardized, data type-specific analyses. All data of one type are analyzed with the same workflow. Although our results may show some differences from an author's publication, our re-analysis of the data makes it feasible to compare data sets from very different sources and to update the data analysis with contemporary methods. For transparency, the methods we use to analyze data are presented below. They are also located online at: *Link*: <https://beta.veupathdb.org/veupathdb.beta/app/static-content/methods.html>

Genome analyses

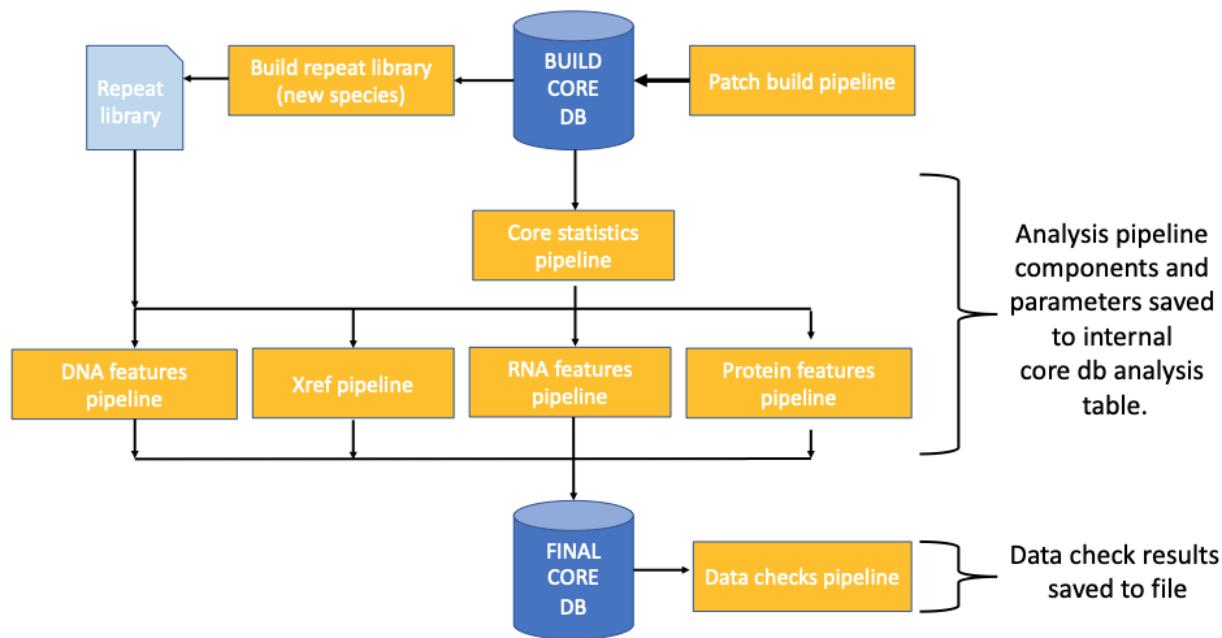
VEuPathDB employs the [Ensembl genome analysis](#) pipelines for analyzing genomic sequence to enhance annotations. While most of the genomic sequence (FASTA) are integrated into VEuPathDB from an INSDC repository, genome annotation (GFF3) may come from either the INSDC repository or a community submission.

Core database pipelines (next figure)- Primary genomic sequence and structural annotation data are loaded into a core database and run through 6 pipelines: core statistics, DNA feature annotation, [external cross reference](#) annotation, [RNA gene](#) annotation, [repeat feature](#) annotation, and [protein feature](#) annotation. The main pipelines applied to the core database and their components are listed in table 1.

Table 1 - Core database analysis pipelines and hive components.

Pipeline	Hive pipeline component	Git repository link
Core db build	Bio::EnsEMBL::EGPipeline::LoadGFF3 ::LoadGFF3	Under code review, details coming soon
Protein features	InterProScan	https://github.com/Ensembl/ensembl-production

DNA features	Bio::EnsEMBL::EGPipeline::DNAFeatures::*	Under code review, details coming soon
RNA features	Bio::EnsEMBL::EGPipeline::RNATFeatures::*	Under code review, details coming soon
Xrefs	Bio::EnsEMBL::EGPipeline::Xref::*	Under code review, details coming soon
Core Statistics	shortnoncodingdensity	https://github.com/Ensembl/ensembl-production
Core Statistics	snpdensity	https://github.com/Ensembl/ensembl-production
Core statistics	codingdensity	https://github.com/Ensembl/ensembl-production
Core statistics	percentgc	https://github.com/Ensembl/ensembl-production
Core statistics	percentagerepeat	https://github.com/Ensembl/ensembl-production
Core statistics	pseudogenedensity	https://github.com/Ensembl/ensembl-production
Core statistics	longnoncodingdensity	https://github.com/Ensembl/ensembl-production



Core database analysis pipelines and hive components

Example ehive pipelines, modules, programs and parameter data from coredb analysis table

Table 2 - Example ehive pipelines, modules, programs and parameter data from coredb analysis table.

Pipeline	Program	Program version	Parameters	Ehive module (Bio::EnsEMBL::)	Database	Database version
DNA features	dustmasker	NULL	NULL	Analysis::Runnable::DustMasker		
DNA features	trf	4	25 7 80 10 40 500 -d -h	Analysis::Runnable::TRF		
DNA features	RepeatMasker	4.0.5	-nolow-gccalc-species "Aedes aegypti" - engine crossmatch -q	Analysis::Runnable::RepeatMasker		
DNA features	RepeatMasker	4.0.5	-nolow -gccalc -lib "location1" -engine crossmatch	Analysis::Runnable::R		

			-q	epeatMasker		
DNA features	RepeatMasker	4.0.5	-nolow -gccalc -lib "location2" -engine crossmatch -q	Analysis::Runnable::RepeatMasker		
	NULL	NULL	NULL	EGPipeline::LoadGFF3: :LoadGFF3		
RNA features	Infernal	1.1		Analysis::Runnable::C MScan		
RNA features	tRNAscan-SE	1.23		Analysis::Runnable::t RNAscan		
RNA features	Infernal	1.1		Analysis::Runnable::C MScan		
RNA features	rfam_12.2_gene	NULL	NULL	EGPipeline::RNAFeatures ::CreateCmscanGenes		
RNA features	mirbase_gene	NULL	NULL	EGPipeline::RNAFeatures ::CreateMirbaseGenes		
RNA features	trnascan_gene	NULL	NULL	EGPipeline::RNAFeatures ::CreateTrnascanGenes		
Xref	xrefchecksum	NULL	NULL	EGPipeline::Xref::LoadUniParc		
Xref	xrefuniparc	NULL	NULL	EGPipeline::Xref::LoadUniProt		
Xref	gouniprot	NULL	NULL	EGPipeline::Xref::LoadUniProtGO		
Xref	xrefuniprot	NULL	NULL	EGPipeline::Xref::LoadUniProtXrefs		
DNA features	blastp	NULL	-word_size 3 -num_alignments 100000 -num_descriptions 100000 -lcase_masking -seg yes -num_threads 3	Analysis::Runnable::Bla stEG		

DNA features	blastp	NULL	-word_size 3 -num_alignments 100000 -num_descriptions 100000 -lcase_masking -seg yes -num_threads 3	Analysis::Runnable::BlastEG		
Protein Features	InterProScan	5.37-76.0	NULL		Prosite patterns	2019_01
Protein Features	InterProScan	5.37-76.0	NULL		SFLD	4
Protein Features	InterProScan	5.37-76.0	NULL		CDD	3.17
Protein Features	InterProScan	5.37-76.0	NULL		Gene3D	4.2.0
Protein Features	InterProScan	5.37-76.0	NULL		HAMAP	2019_01
Protein Features	InterProScan	5.37-76.0	NULL		PANTHER	14.1
Protein Features	InterProScan	5.37-76.0	NULL		ncoils	2.2.1
Protein Features	InterProScan	5.37-76.0	NULL		Prosite profiles	2019_01
Protein Features	InterProScan	5.37-76.0	NULL		Pfam	32
Protein Features	InterProScan	5.37-76.0	NULL		PRINTS	42
Protein Features	InterProScan	5.37-76.0	NULL		Smart	7.1
Protein Features	InterProScan	5.37-76.0	NULL		SuperFamily	1.75
Protein Features	InterProScan	5.37-76.0	NULL		TIGRfam	15
Protein Features	InterProScan	5.37-76.0	NULL		InterPro-2GO	NULL
Protein Features	InterProScan	5.37-76.0	NULL		PIRSF	3.02
Protein Features	InterProScan	5.37-76.0	NULL		SignalP	4.1

Protein Features	InterProScan	5.37-76.0	NULL		TMHMM	2.0c
Protein Features	InterProScan	5.37-76.0	NULL		Seg	NULL
Protein Features	InterProScan	5.37-76.0	NULL		MobiDBLite	2
Protein Features	InterProScan	5.37-76.0	NULL		InterPro-2Pathway	NULL

*Location 1

"/homes/jallen/scratch/vb/recent_assemblies/aedes_aegypti/RepeatModeler/aedes_aegypti.rm.lib"

*Location 2 "/nfs/panda/ensemblgenomes/vectorbase/data/tefam/aegypti_tefam.lib"

Supplements to the EBI Pipelines

VEuPathDB supplements the EBI pipeline with workflows that produce data for EST alignments, Open reading frames, and synteny (Table).

EST alignments: BLAT is applied to EST sequences that have been blocked using RepeatMasker.

Open reading frame generation: Open reading frames are generated from genomic DNA or EST sequences. The analysis produces a gff file containing the ORFs (50 or more amino acid translations of the input nucleic acid) for the 6 reading frames. The translations in all 6 reading frames do not necessarily begin with MET but always end at a stop codon. ORF names are in the form template-frame-start-end, e.g. AAEL01000396-5-5847-4366.

Synteny: VEuPathDB uses an in-house script called runMercator to run pair-wise alignments that employs Mercator and MAVID for comparative genome analysis. Mercator generates orthology maps using genomes and exon coordinates to create exon translations for protein BLAT alignments. The orthology maps are used as a guide for MAVID which also uses a phylogenetic newick tree to generate gene alignments.

Details for the supplements to the EBI pipelines

Table 3: Details for supplements to the EBI Pipeline for genomes

Data	Program	Parameters or VEuPathDB GitHub repository	Version
EST alignments	BLAT	-nohead -maxIntron=1000 -t=dna -q=dna -dots=10	35
Open Reading Frames	orfFinder	--minPepLength 50 https://github.com/VEuPathDB/.../orfFinder	

Synteny	Mercator and MAVID	<p>-p <PATH TO MERCATOR DIRECTORY> -t <tree string> -m <MAVID_EXE> -c <CNDSRC_DIR> -d draftGenome1... -d draftGemoneN -n nonDraftGeome1... -n nonDraftGenomeN -r referenceGenome https://github.com/VEuPathDB/.../runMercator</p>	Mercator mapmaker 0.4 (2016-01-21) Mavid Version 2.0.4
---------	--------------------	---	---

In-house genome analyses in Lieu of the EBI Pipeline

On rare occasions the EBI pipeline cannot be applied to a genome. For example, genomes that are not housed at an INSDC repository cannot be analyzed by the EBI pipeline. VEuPathDB uses the following in-house analyses in lieu of the EBI pipeline.

BLAT against NRDB: For every genome, VEuPathDB runs BLAT alignments of the annotated proteins against the GenBank Non-Redundant Protein Sequence Database (NRDB) to identify possible relationships and alignments outside the scope of VEuPathDB-supported organisms.

Compute open reading frames: Open reading frames are generated from genomic DNA or EST sequences. The analysis produces a gff file containing the ORFs (50 or more amino acid translations of the input nucleic acid) for the 6 reading frames. The translations in all 6 reading frames do not necessarily begin with MET but always end at a stop codon. ORF names are in the form template-frame-start-end, e.g. AAEL01000396-5-5847-4366.

DNA repeats: The Tandem Repeats Finder program locates and displays tandem repeats in genomic sequences.

EST alignments: BLAT is applied to EST sequences that have been blocked using RepeatMasker.

Protein domain annotations: InterProScan scans protein sequences against the protein signatures of the InterPro member databases and generates a file containing the domain matched, description of the InterPro entry, GO descriptions and E-values.

Signal peptide prediction: Signal P is used to identify signal peptides and their likely cleavage sites. A signal peptide is a short peptide present at the N-terminus of most newly synthesized proteins that are destined towards the secretory pathway.

Syntenic sequences: VEuPathDB uses an in-house script called runMercator to run pair-wise alignments that employs Mercator and MAVID for comparative genome analysis. Mercator generates orthology maps using genomes and exon coordinates to create exon translations for

protein BLAT alignments. The orthology maps are used as a guide for MAVID which also uses a phylogenetic newick tree to generate gene alignments.

Transmembrane domain prediction: TMHMM is used to predict transmembrane domain presence and topology from protein sequences.

tRNA gene prediction: tRNAScan identifies transfer RNA genes in transcript or genome sequences.

Details for the VEuPathDB in-house pipelines

Table 4: Details for genome analyses in lieu of EBI Pipeline

Data	Program	Parameters or VEuPathDB GitHub repository	Version
BLAT against NRDB	BLAT	-nohead -maxIntron=1000 -t=dna -q=prot -dots=10 -minScore=25 -minIdentity=20	35
Computed Open Reading Frames	orfFinder	--minPepLength 50 https://github.com/VEuPathDB/.../orfFinder	
DNA Repeat regions	Tandem Repeats Finder	2 7 7 80 20 50 500	4.04
EST alignments	BLAT	-nohead -maxIntron=1000 -t=dna -q=dna -dots=10	35
Signal peptide predictions	SignalP	-t euk -f short -m nn+hmm -q -trunc 70	3.0
Synteny	Mercator and MAVID	-p <PATH TO MERCATOR DIRECTORY> -t <tree string> -m <MAVID_EXE> -c <CNDSRC_DIR> -d draftGenome1... -d draftGemoneN -n nonDraftGeome1... -n nonDraftGenomeN -r referenceGenome https://github.com/VEuPathDB/.../runMercator	Mercator mapmaker 0.4 (2016-01-21) Mavid Version 2.0.4

Transmembrane domain prediction	TMHMM	Nice -short	2.0c

Proteomics

VEuPathDB integrates the results of proteomics experiments as peptides aligned to a reference genome or as abundance data assigned to a gene. We do not reanalyze the raw mass spec data but instead use an in-house plugin that loads found peptides or abundance data from tab delimited input files of a specific format.

[Details for the VEuPathDB in-house proteomics pipeline](#)

RNA-Sequence

VEuPathDB integrates RNA-Seq data from many different experiments and analyzes all data with the same EBI RNA-Seq analysis pipeline. The RNA sequence data that we integrate is processed at EBI.

The following is a general outline of the analysis process.

- Trim poor quality data (Trimmomatic)
- HiSAT2 alignment to a reference genome
- HT-Seq-count to tally aligned reads per gene
- Convert to transcripts per kilobase million (TPM)
- DESeq2 to determine differential expression

[EBI RNA-Seq pipeline details](#)

ChIP-Sequence

VEuPathDB integrates ChIP-Seq data from many different experiments and sources. Details coming soon.

Copy Number Variation

VEuPathDB analyzes whole genome resequencing data to estimate each gene's copy number in resequenced strains. Details coming soon.

Genetic Variation and SNP calling

VEuPathDB analyzes whole genome resequencing data to call single nucleotide polymorphisms of isolates. Details coming soon.

Protein Array data

VEuPathDB integrates protein array data from serum antibody microarray experiments. Analysis details coming soon.

Metabolic Pathways

VEuPathDB integrates metabolic pathways from KEGG and MetCyc. Metabolic pathways are associated with genes via Enzyme Commission annotations. Details coming soon.

Dataset descriptions:

This document provides a high-level overview of the software infrastructure utilized by the VEuPathDB BRC to load, integrate and provide data to users. Please check [a list of all the data sets](#) loaded in our VEuPathDB sites utilizing this infrastructure.

Link: <https://beta.veupathdb.org/veupathdb.beta/app/search/dataset>AllDatasets/result>

All loaded datasets are described appropriately with links to publications and repositories as appropriate.

- Data Set (Name)
- Organism(s) (source or reference)
- Category (type of data set, e.g. RNA-Seq)
- Description
- Release # /Date
- Summary
- Release Policy
- Publications
- Contact
- Contact Institution

Technical infrastructure and software documentation:

Link: <https://beta.veupathdb.org/veupathdb.beta/app/static-content/infrastructure.html>

The above document describes the overall infrastructure of our data loading system and websites with links to appropriate GitHub repositories for users who want to explore using VEuPathDB infrastructure to build their own genomics database and website.

Browser Compatibility Statement

We recognize that our users access VEuPathDB using various Internet Browsers and Operating Systems. Our goal is to ensure that you have the best possible experience on VEuPathDB, but it is impossible to develop applications that work identically, efficiently and effectively on all web browsers.

Based on our site usage statistics we support the following browsers used by greater than 95% of our visitors:

- Firefox
- Safari
- Chrome

Feel free to [contact us](#) about any browsing issues you might come across.

Data Loading and Database Schema

We use the [Genomics Unified Schema \(GUS\) database schema](#) and data loading infrastructure and its framework available at [GusAppFramework](#). This includes not only a comprehensive database schema for integrating and representing genomic and functional (or post) genomic data but also tools for loading said data into that system. We have made some extensions to the schema and tools for VEuPathDB specific purposes primarily to generate de-normalized views of the data for query optimization purposes.

Our data are all stored in Oracle12c databases. Our software infrastructure also supports PostgreSQL but we have some Oracle specific SQL constructs in our model that would need to be changed in order to run successfully in PostgreSQL.

We load all data using an in house engineered workflow system called [ReFlow](#). Briefly, ReFlow is engineered to be an efficient graph-based workflow system. In it each step (node in the graph) has the ability to be undone and subsequently rerun with updated data. This was a significant requirement as it enables us to undo entire genomes when the annotation or underlying sequence changes. This results in automated removal of all data dependent on that genome. When the step is re-run with the new annotation, all dependent data are recomputed and reloaded automatically, thus greatly improving our ability to keep these complex databases up-to-date.

The ReFlow workflow system utilizes another piece of software developed at the University of Pennsylvania to schedule, manage and monitor running tasks called [DistribJob](#). DistribJob distributes tasks generated from a large input dataset such as a set of sequences to compute nodes in a cluster for analysis and retrieves and collates the results in an efficient manner. We

automate the running of large compute tasks on compute clusters located at the University of Pennsylvania and the University of Georgia.

Code Availability

To facilitate greater transparency and tool reuse, our codebase has been migrated to the GIT repository (<https://github.com/VEuPathDB>).

Web Presentation System and User Interfaces

Our websites are based on code that we developed and have released to the community called the Strategies-WDK (Strategies Web Development Kit) which enables the graphical strategies search system. You can download the software and see documentation for this toolkit at [Strategies-WDK](#). This toolkit enables us to represent our data as an XML model which is then turned into the web interfaces that are presented to users using these tools.

Software Code Repository

To facilitate greater transparency and tool reuse, our codebase has been migrated to the GIT repository (<https://github.com/VEuPathDB>).

System Hardware and Third-Party Software

VEuPathDB maintains redundant database and content web servers at the University of Pennsylvania and the University of Georgia to minimize interruptions for our users during maintenance periods. Additionally VEuPathDB compute and data loading servers are located at the University of Pennsylvania.

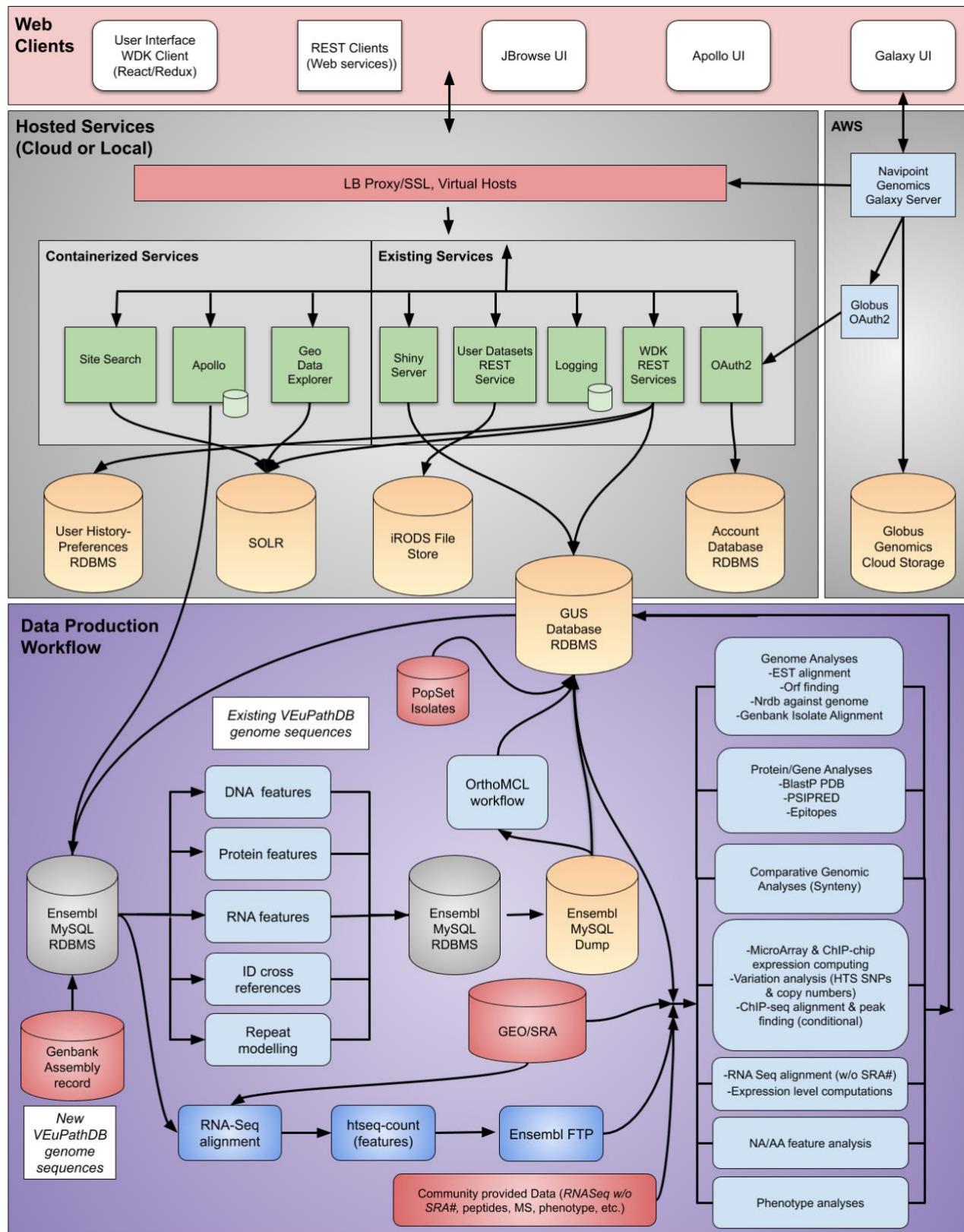
Server configurations are coordinated and deployed through Puppet automation software (<http://puppetlabs.com/>). Custom infrastructure software is versioned and deployed through standard RPM/YUM mechanisms. When appropriate, software builds are automated using Jenkins Continuous Integration Server (<http://jenkins-ci.org/>)

System infrastructure statistics (CPU load, I/O, etc) are gathered with collectd (<http://collectd.org/>) and in-house applications and feed to Graphite (<http://graphite.wikidot.com/>) for human review. Nagios (<http://www.nagios.org/>) provides notifications of system degradations.

Both Universities also maintain large compute clusters that are heavily utilized by VEuPathDB in order to analyze and load incoming data in a timely fashion. The linked document below describes our actual hardware and includes a list of third-party software required in order to analyze, load and present data via our websites.

Overview of the VEuPathDB Data Production Workflow and Architecture

The complete pathway from data acquisition to web presentation and utilization by users is detailed in the next figure. The data production workflow (bottom, purple rectangle) begins with data acquired from numerous sources (indicated in red) including direct deposition by the community. These data are processed through two main pipelines (indicated in blue), one targeting genome features and RNA-Seq mapping using Ensembl pipelines and the second focused on comparative genomics and analysis of other types of omics-scale data including proteomics and high-throughput phenotypic screens among others. Quality-controlled analysed data are loaded into the GUS database (indicated in gold spanning the purple data production and grey services rectangles). The GUS database links the data to the servers and services (green) and resources that allow the user community to securely access, interrogate, visualize, download and further analyse data, including their own data via the provided web clients (indicated in the pink rectangle).



VEuPathDB Website Privacy Policy

We do not use or share any of your personal information for any purpose unrelated to the functionality of the websites; however, we do collect some information to help us understand how our sites are being used and to improve community support.

UPDATED: April 12, 2020

Introduction

VEuPathDB (also referred to as ‘we’ throughout this document) is committed to protecting its users’ data and privacy. The purpose of this page is to provide you with information about how the data we collect from users of VEuPathDB websites is used or shared. We may update this Privacy Notice from time to time. We encourage you to visit this page frequently and take note of the date updated field above.

We do not use or share any of your personal information for any purpose unrelated to the functionality of the websites; however, we do collect some information to help us understand how our sites are being used in order to improve community support and to enhance the VEuPathDB community’s experience when visiting our sites.

Information Automatically Collected

When you browse VEuPathDB sites, certain information about your visit will be collected. We automatically collect and store the following type of information about your visit:

- The IP address of the client making the request. Often the IP address is that of your personal computer or smart phone; however, it might be that of a firewall or proxy your internet provider manages.
- The operating system and information about the browser used when visiting the site.
- The date and time of each visit.
- Pages visited.
- The address of a referring page. If you click a link on a website that directs you to a VEuPathDB page, the address of that originating web page will be collected. This “referrer” information is transmitted as part of the browser and server communications; it is not based on any marketing or partnering agreements with the referring site.

This automatically collected information does not identify you personally unless you include personally identifying information in a support form request; see the “Contact Us” policy below for details. We use this information to measure the number of visitors

to our site. The aggregate data may be included in prospectuses and reports to funding agencies.

Information You Directly Provide

The Basket, Favorites, Public Strategies, Gene/Sequence Comment and GBrowse Track features of the VEuPathDB websites require that you register for an account. A valid email address is required so we can send you your temporary account password. An anonymous email service can be used if you do not want to provide personally identifying information.

Your email address will be used to send you infrequent alerts if you subscribe to receive them. We do not sell or distribute email addresses to third parties.

We also ask for your name and institution during account registration. If you add a comment to a Gene or a Sequence, your name and institution will be displayed with the comment. If you make one of your strategies public, your name and institution will be displayed with it. We do not routinely verify the validity of names and institutions associated with comments or public strategies; however, we will delete accounts or comments if we believe them to be fraudulent based on inappropriate activity or posted content. We will not sell or distribute your name or institution to third parties.

When you log in, the client IP address is recorded. This IP address can be correlated with the address automatically collected as noted above. If your user profile personally identifies you, then it may be possible to associate you with your detailed activity on VEuPathDB web sites.

“Contact Us” Form

The header on each web page includes a “Contact Us” link to a form where users can submit questions, error reports, feature requests, and dataset proposals. Submissions through this form are emailed to VEuPathDB staff and recorded in a project management application accessible only by VEuPathDB staff.

The form includes a field for an email address. If the email address identifies you personally, say if you use your institutional email, then your correspondence with us will likewise be linked to you. A valid email is not strictly required, although we cannot reply to you without one.

When you submit the form, your IP address and browser version will be recorded for internal use. In the case of reported bugs or other site errors, this information may be used by technical staff to help locate your session in the server logs to aid in troubleshooting the issue. This does have the side effect of making it possible to associate an IP address with an email address which may, in turn, personally identify you. However, VEuPathDB does not publicly release this information.

How VEuPathDB Uses Cookies

VEuPathDB uses cookies to associate multiple requests by your web browser into a stateful session. Cookies are essential to track the state of query strategies, gene baskets and authentication.

Some cookies persist only for a single session. The information is recorded temporarily and is erased when the user quits the session or closes the browser. Others may be persistently stored on the hard drive of your computer until you manually delete them from a browser folder or until they expire, which can be months after they were last used.

Cookies can be disabled in your browser (refer to your browser's documentation for instructions); however, the majority of the website functionality will be unavailable if cookies are disabled.

Google Analytics

Google Analytics provides aggregate measurements of website traffic including counts of page hits and unique users along with statistics on countries of origin.

The raw measurements and statistics are only available to approved VEuPathDB staff. Aggregated data may be included in prospectuses and reports to funding agencies.

Third-Party Websites and Applications

Third-party websites and applications are not exclusively operated or controlled by VEuPathDB. By using these third-party websites, individuals may be providing nongovernmental third-parties with access to personally identifying information.

Twitter

VEuPathDB maintains a presence on Twitter in the form of a [VEuPathDB branded page](#). This page allows for a direct connection with end users to promote information related VEuPathDB services and to disseminate educational information on research publications, news and events related to the biology of eukaryotic pathogens. Postings may also include information about planned service maintenance and outages.

Twitter collects profile information such as name and email address about users who register to use this third-party website. Depending on the user's privacy settings, this information may be displayed on the user's profile page or in the user's tweets which may be retweeted on VEuPathDB's page. The VEuPathDB Twitter account may post the authors and institutions of publicly published scientific papers and news articles. VEuPathDB does not actively collect or maintain personally identifying information through its use of Twitter. VEuPathDB will redact or refrain from retweeting a posting that contains obviously identifiable personal information. A Twitter account is not

required to read VEuPathDB postings on Twitter. VEuPathDB does not collect or use personal information outside of Twitter's site.

Twitter is hosted and maintained by a third party which may use browser tracking and related technologies to collect information about visitors to twitter.com and its affiliates. Refer to Twitter's privacy statement, <https://twitter.com/en/privacy>, for more information.

Facebook

VEuPathDB maintains a presence on Facebook in the form of a [VEuPathDB branded page](#). This page allows for a direct connection with end users to promote information related VEuPathDB services and to disseminate educational information on research publications, news and events related to the biology of eukaryotic pathogens. Postings may also include information about planned service maintenance and outages.

Like Twitter, Facebook collects profile information, including name and email address, from its users. Depending on the user's privacy settings, this information may be displayed on the user's profile page along with any activity such as comments or "likes" on the VEuPathDB Facebook page or in posts that VEuPathDB shares on Facebook. VEuPathDB does not collect or use any personally identifying information outside of our Facebook page. To understand how Facebook collects and uses personal information, refer to their data policy page, <https://www.facebook.com/policy.php>.

YouTube

VEuPathDB maintains a presence on YouTube in the form of a [VEuPathDB branded page](#). This page provides tutorials on the use of our websites.

YouTube also requires some information when users create an account, including an email address, and users may choose to provide a name and other identifying information in their public profile. Depending on their individual privacy settings, some personally identifiable information may be available to other users, including VEuPathDB. However, VEuPathDB does not collect or use any of that information outside of its YouTube interactions. You can view videos without signing in to an account, but you must be a registered user in order to comment. As YouTube is a Google service, the VEuPathDB youtube channel is subject to [Google's privacy policy](#).

Globus Genomics

The VEuPathDB Galaxy Data Analysis Service is a workspace for large-scale data analyses. Developed in partnership with [Globus Genomics](#), workspaces offer a private analysis platform with published workflows and pre-loaded annotated genomes. The workspace is accessed through the "Analyze My Experiment" tab on the home page of any VEuPathDB resource and can be used to upload your own data, compose and run custom workflows, retrieve results and share workflows and data analyses with colleagues.

The VEuPathDB Galaxy Data Analysis Service is hosted by Globus Genomics, an affiliate of Globus. The first time you visit VEuPathDB Galaxy you will be asked to sign up with Globus in order to set up your private Galaxy workspace. Linking your Globus account with your VEuPathDB account is necessary so that input data and analysis results can be transferred between the two systems. We encrypt data transfers and storage, but ultimately, we cannot guarantee the security of data transmissions among VEuPathDB, Globus and affiliates, Amazon Cloud Services, and the user. It is your responsibility to back up your data and obtain any required permissions from your study and/or institution prior to uploading data for analyses on the VEuPathDB Galaxy platform. Do not use, transmit, upload or share any human identifiable information in the files you analyze. VEuPathDB, Globus and affiliates, the University of Georgia, the University of Pennsylvania, the University of Liverpool, and Amazon Cloud Services do not take any responsibility and are not liable for the loss and/or release of any data you analyze via the VEuPathDB Galaxy platform. We encourage you to review the [Globus' privacy policy](#).

Your Rights based on the General Data Protection Regulation (GDPR)

To read more about GDPR please check the [GDPR website](#).

1. The right of transparency and modalities. The privacy policy should be clear and easy to follow in explaining what data we collect and how we use it.
2. The right to be informed about when data is gathered. This is described in the privacy policy, during the registration process (if you choose to register), site banner and an email sent out to all registered users on May 25, 2018.
3. The right of access. You can ask for what specific data we have about you and how we use it.
4. The right to rectification. We will correct any errors in your personal data that you point out to us.
5. The right to be forgotten. We are happy to delete your account and info when you make such a request.
6. The right to restrict processing. You have the right to request that we restrict the use of your data.
7. The right for notification obligation regarding rectification/erasure/restriction.
8. The right to data portability.
9. The right to object to the processing of your personal data at any time.
10. The right in relation to automated decision making and profiling. Basically, you have the right not to be subject to decisions based solely on automated processing which significantly affect you.

To make any of the above stated requests or if you have any questions please email us at help@VEuPathDB.org

VEuPathDB Accessibility Conformance

We strive to make the VEuPathDB website accessible. Wherever possible, our sites have alternative text for images, the pages can be magnified and read by screen readers. The results of some searches are less accessible as they are dynamic and user-specific, thus alternative text describing images cannot be generated. We update our Voluntary Product Accessibility Template, VPAT, 508 accessibility report annually. The more recent version is located at:
https://veupathdb.org/documents/VEuPathDB_Section_508_BRC4.pdf

VEuPathDB personnel

VEuPathDB Management

Beatrice Amos, Annotation Manager
Cristina Aurrecoechea, User Interface and Portal Manager
Bob Belnap, Systems and Databases Manager
John Brestelli, Data Development Manager
Brian Brunk, VEuPathDB Senior Manager
Mark Caddick, Wellcome Trust PI; Co-I NIAID BRC Contract
George Christophides, Co-I, NIAID BRC Contract
Kathryn Crouch, Co-I, Wellcome Trust
Jeremy DeBarry, Project Coordinator
Steve Fischer, Software and Infrastructure Manager
Paul Flieck, Co-I, NIAID BRC Contract
Omar Harb, Director of Scientific Outreach & Education
Jessica C Kissinger, Joint-PI, NIAID BRC Contract; WT Co-PI
Dan Lawson, Project Coordinator
Wei Li, Data Loading Manager
Mary Ann McDowell, Joint-PI, NIAID BRC Contract
David S Roos, Joint-PI, NIAID BRC Contract; WT Co-PI
Chris J Stoeckert, Co-I, NIAID BRC Contract

To contact any one of us please use the [contact us form](#).

Current VEuPathDB Team members

Beatrice Amos⁴, Rachel Ankirskiy¹, Cristina Aurrecoechea¹, Matthieu Barba⁹, Ana Barreto³, Evelina Basenko⁴, Wojtek Bazant², Dan Beiting², Bob Belnap¹, Ulrike Böhme⁵, John Brestelli³, Brian Brunk², Mark Caddick⁴, Danielle Callan², Mikkel Christensen⁹,

George Christophides⁸, Kathryn Crouch⁶, Katie Cybulski⁷, Elaine Daugan⁴, Jeremy DeBarry¹, Ryan Doherty³, Yikun Duan², Dave Falke¹, Steve Fischer³, Paul Flicek⁹, Bindu Gajria², Gloria I. Giraldo-Calderón⁷, Omar S. Harb², Elizabeth Harper², Danica Helb², Mark Hickman², Connor Howington⁷, Sufen Hu², Jay Humphrey¹, John Iodice³, John Judkins², Sarah Kelly⁸, Jessica C. Kissinger¹, Dae Kun Kwon⁷, Kris Lamoureux¹, Daniel Lawson⁸, Wei Li², Brianna Lindsay², Jamie Long², Bob MacCallum⁸, Gareth Maslen⁹, Mary Ann McDowell⁷, Greg Milewski², Jarek Nabrzyski⁷, David S. Roos², Samuel Rund⁷, Steph Wever Schulman², Achchuthan Shanmugasundram⁴, Vasili Sitnik⁹, Drew Spruill¹, David Starns⁴, Christian J. Stoeckert Jr.³, Sheena Shah Tomko², Haiming Wang¹, Susanne Warrenfeltz¹, Robert Wieck⁷, Mariann Winkelmann², Lin Xu², Jie Zheng³.

Previous VEuPathDB Team members, 2004-2020

Antelmo Aguilar⁷, James Allen⁹, Alexis Allot⁹, Nora Besansky⁷, Austin Billings², Sanjay Boddu⁹, Steve Bogol⁷, Ewan Birney⁹, Andrew Brockman⁸, Robert Bruggner⁷, Ja'Shon Cade³, David Campbell⁷, Cristian Cocos², Frank Collins (VectorBase Principal Investigator, 2004-2018)⁷, Kathy Couch¹, Greg Davies⁷, Ale Diaz Miranda², Emmanuel Dialynas, Jennifer Dommer³, Vicky Dritsou, Scott Emrich¹⁰, Xin Gao², William Gelbart¹², Sandra Gesing⁷, Alan Gingle¹, Greg Grant³, Matt Guidry¹, Martin Hammond⁹, Mark Heiges¹, Christiane Hertz-Fowler (Principal Investigator, WT 2008-2019)⁴, Nicholas Ho⁸, Daniel Hughes⁹, Frank Innamorato³, San James¹⁴, Amie Jaye⁸, Fotis Kafatos⁸, Paul Kersey⁹, Ioannis Kimitzoglou⁸, Nathan Konopinski⁷, Carolyn Knoll², Eileen T. Kraemer¹, Nick Langridge⁹, Cris Lawrence², Neil Lobo⁷, Christos (Kitsos) Louis¹¹, Ross Madden⁶, Greg Madey⁷, Elisabetta Manduchi³, Karine Megy⁹, John A. Miller⁶, Elvira Mitraka¹¹, Vishal Nayak³, Cary Pennington¹, Deborah F. Pinney³, Brian Pitts¹, Jane A. Pulman⁴, Caleb Reinking⁷, Seth Redmon, Chris Ross¹, Andrew Sheehan⁷, Fatima Silva⁴, Ganesh Srinivasamoorthy¹, Scott Szakonyi⁷, Pantelis Toplais¹¹, Ryan Thibodeau¹, Charles Treatman², Betsy Wenthe¹, Matt Vander Werf⁷, Maggie Werner-Washburne¹³, Patricia L. Whetzel³, Derek Wilson⁹, Andrew Yates⁹

¹University of Georgia, Athens, GA 30602, USA

²University of Pennsylvania, Philadelphia, PA 19104, USA

³University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA

⁴University of Liverpool, UK

⁵Wellcome Sanger Institute, Hinxton, CB10 1RQ, UK

⁶Wellcome Centre for Integrative Parasitology, University of Glasgow, UK

⁷University of Notre Dame, Notre Dame, IN 46556, USA

⁸Imperial College London, South Kensington, London SW7 2BU, UK

⁹European Bioinformatics Institute, Hinxton, CB10 1SD, UK

¹⁰University of Tennessee, Knoxville, TN 37996, USA

¹¹Institute of Molecular Biology and Biotechnology-FORTH, Heraklion, Crete, Greece

¹²Harvard University, Cambridge, MA 02138, USA

¹³University of New Mexico, Albuquerque, NM 87131, USA

¹⁴Makerere University and Infectious Diseases Research Collaboration (IDRC),
Kampala, Uganda

VEuPathDB acknowledgements

- All the community members who have contributed data (often pre-publication), entered user comments or sent us their suggestions.
- Scientists who provided images to be used as a template for the logos in our websites and for images used in the header section of our sites:

AmoebaDB:

William Petri

Serge Ankri

Craig Roberts

Fiona Henriquez

Hugo Aguilar-Díaz

CryptoDB:

Boris Striepen

FungiDB:[ZyGoLife](#)

Jason Stajich

Zachary Lewis

GiardiaDB:

Fran Gillin

Tineke Lauwaet

Barbara Davids

Scott Dawson

MicrosporidiaDB:

Gira Bhabha

Pattana Jaroenlak (Michael)

Damian Ekiert

Michael Cammer

PiroplasmaDB:

Ellen Yeh

Lowell Kappemeyer

Audrey Lau

Dirk Dobbelaere

Manoj Duraisingh

Brendan Elsworth

Caroline Keroack

Isabelle Coppens

PlasmoDB:

Lawrence Bannister

Lewis Tilney

Pedro Moura

ToxoDB:

David Roos

TrichDB:

Antonio Pereira-Neves

Marlene Benchimol

TriTrypDB:

Rick Tarleton

Richard Wheeler

Leandro Lemgruber Soares

Margaret Mullin

Camila Silva Gonçalves

Wanderley de Souza

Maria Cristina Machado Motta

VEuPathDB Community Representatives

VEuPathDB encourages community members to provide feedback about our resources. We get feedback from many community members including those listed below who have been active in our open community meetings. We encourage you to get involved.

Feel free to [contact us](#) any time.

Amoeba: Open community call and Carol Gilchrist, Upi Singh

Cryptosporidium: Gregory Buck, Guy Robinson, Karin Troell, Sumiti Vinayak, Jonathon Wastling, Giovanni Widmer, Lihua Xiao, Guan Zhu

Fungi: Bridget Barker, Elaine Bignell, Katherine Borkovich, Michael Bromley, Christina Cuomo, Tamara Doering, Jay Dunlap, Michael Freitag, Louise Glass, Kim Hammond-Kosack, Guilhem Janbon, Seogchan Kang, Theo Kirkland, Corby Kistler, Jennifer Lodge, Robin May, Jessie Uehling, Sinem Beyhan, Douglas Lake, Natalie Mitchell, Maureen Donlin, Vera Meyer, Marc Orbach, Nadia Ponts, Antonis Rokas, Jason Stajich, Matt Sachs, George R. Thompson, Martin Urban, Nathan P. Wiederhold

Giardia: Scott Dawson, Fran Gillin, Adrian Hehl, Aaron Jex, Hilary Morrison, John Samuelson, Cornelia Spycher, Staffan Svard

Microsporidia: James Becnel, Nicolas Corradi, Elizabeth Didier, Patrick Keeling, Emily Troemel, Louis Weiss

Piroplasma: Open community call and Choukri Mamoun

Plasmodium: John Adams, Chris Janse, Rays H.Y. Jiang, Shahid Khan, Stuart Ralph, Akhil Vaidya, Andy Waters

Toxoplasma: John Boothroyd, Jon P. Boyle, Vern B. Carruthers, Marc-Jan Gubbels, Kami Kim, Markus Meissner, Jeroen Saeij, Lilach Sheiner, Ross Waller, Michael White

Trichomonas: Jane Carlton, Patricia Johnson, Steven Sullivan, Jan Tachazy

Trypanosoma/Kinetoplastids: Fernan Aguero, Vivian Bellofatto, Richard Burchmore, George Cross, Angela Cruz, Antonio Estevez, Mark Field, Catarina Gadelha, Eva Gluenz, Keith Gull, John Kelly, Annette MacLeod, Jeremy Mottram, Torsten Ochsenreiter, Marc Ouellette, Barbara Papadopoulou, Laurie Read, Sergio Schenkman, Rick Tarleton, Brent Weatherly, Bill Wickstead, Michael (Mick) Urbaniak

Vectors: Gregory Dasch, Jeff Grabowski, María de Lourdes Muñoz, Monika Gulia-Nuss, Sukanya Narasimhan, Kristin Michel, Michael Povelones, Igor Sharakhov, Ronald van Rij, Rob Waterhouse

Previous Scientific Working Group

VEuPathDB wishes to acknowledge previous scientific working group members. They provided regular feedback oversight and guidance.

Lyric Bartholomay	Michael Gottlieb	Malcolm McConville	John Taylor
Matt Berriman	Keith Gull	Nicola Mulder	Jake Tu
Bill Black	Matthew Hahn	Uli Munderloh	Brett Tyler
John Boothroyd	Adrian Hehl	Daniel Neafsey	Kenneth Vernick
Greg Buck	Steve Higgs	Kenneth Olson	Sarah Volkman
Geraldine Butler	Catherine Hill	Bill Petri	Jonathon Wastling

Angela Cruz	Marcelo Jacobs-Lorena	Barry Pittendrigh	Scott Weaver
George Dimopoulos	Anthony (Tony) James	Jeffrey Powell	Louis Weiss
Martin Donnelly	Pedro Lagerblad de Oliveira	Hillary Ranson	Dyann Wirth
Patrick Duffy	Greg Lanzaro	Alexander Raikhel	Jennifer Wortman
Pascale Gaudet	Daniel Masiga	Lincoln Stein	Guilun Yan

Website usage statistics

We collect project-specific website usage statistics using the program awstats. The links to all statistics (they are continually updated) are listed below. A sample website usage statistics report, as prepared by awstats is shown for an example of the types of data that are available. Users can visit the site and select custom reporting periods.

Website usage links:

- <https://amoebadb.org/awstats/awstats.pl>
- <https://cryptodb.org/awstats/awstats.pl>
- <https://fungidb.org/awstats/awstats.pl>
- <https://giardiadb.org/awstats/awstats.pl>
- <https://hostdb.org/awstats/awstats.pl>
- <https://microsporidiadb.org/awstats/awstats.pl>
- <https://piroplasmadb.org/awstats/awstats.pl>
- <https://plasmodb.org/awstats/awstats.pl>
- <https://toxodbd.org/awstats/awstats.pl>
- <https://trichdb.org/awstats/awstats.pl>
- <https://tritrypdb.org/awstats/awstats.pl>
- <https://orthomcl.org/proxystats/awstats.pl?config=orthomcl.org> (this link will be updated when OrthoMCLDB moves into the new UI)
- <https://veupathdb.org/awstats/awstats.pl>
- <https://vectorbase.org/awstats/awstats.pl>

Sample awstats report from FungiDB.org

Statistics for: fungidb.org

Last Update: 17 Sep 2020 - 01:00
Reported period: Sep 2020

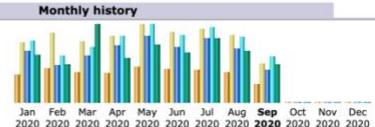


Summary

Reported period	Month Sep 2020
First visit	01 Sep 2020 - 00:01
Last visit	17 Sep 2020 - 00:59
Unique visitors	2,423
Number of visits	5,080
Pages	743,449
Hits	1,034,297
Bandwidth	23.66 GB

* Not viewed traffic includes traffic generated by robots, worms, or replies with special HTTP status codes.

Monthly history

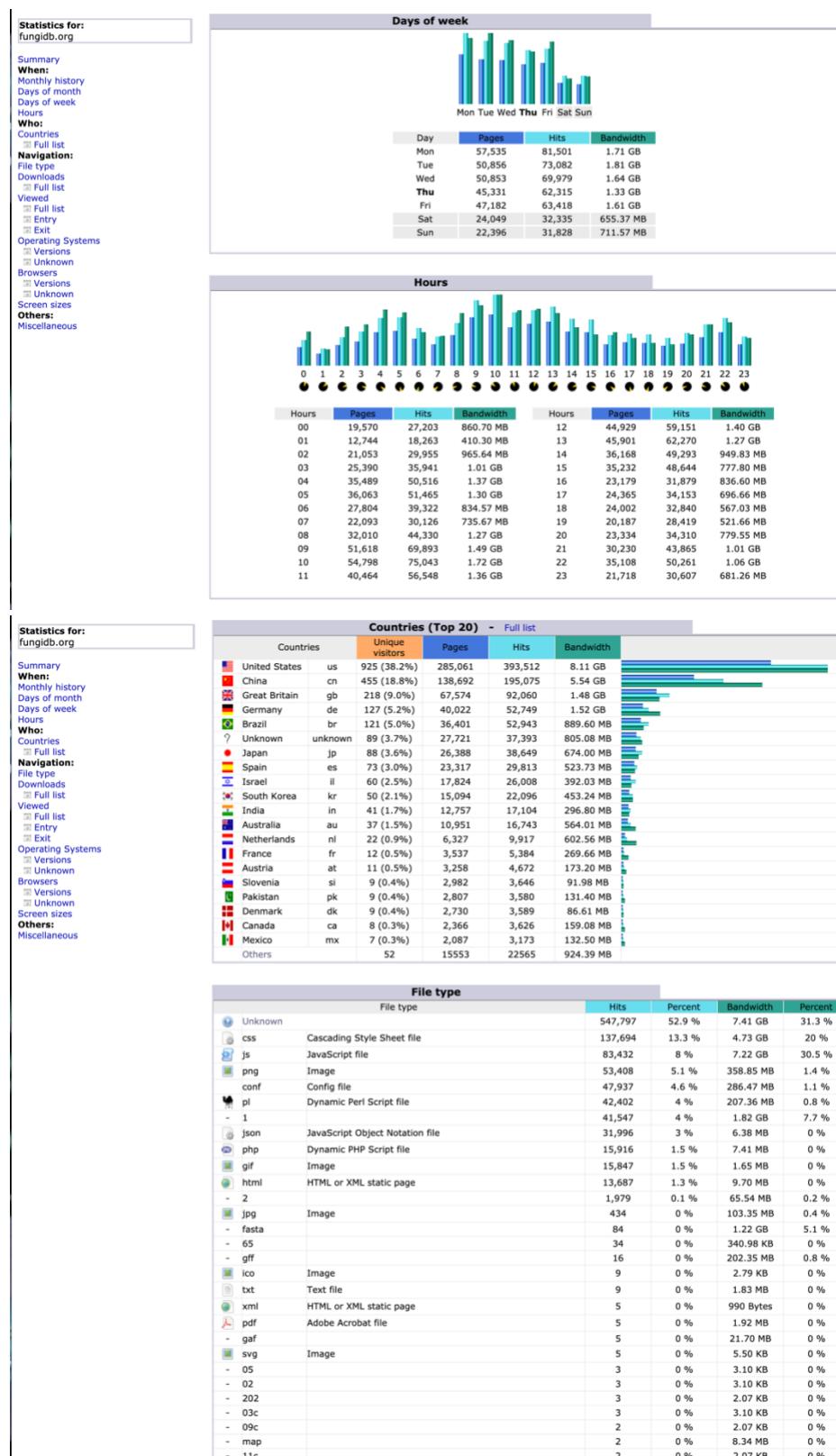


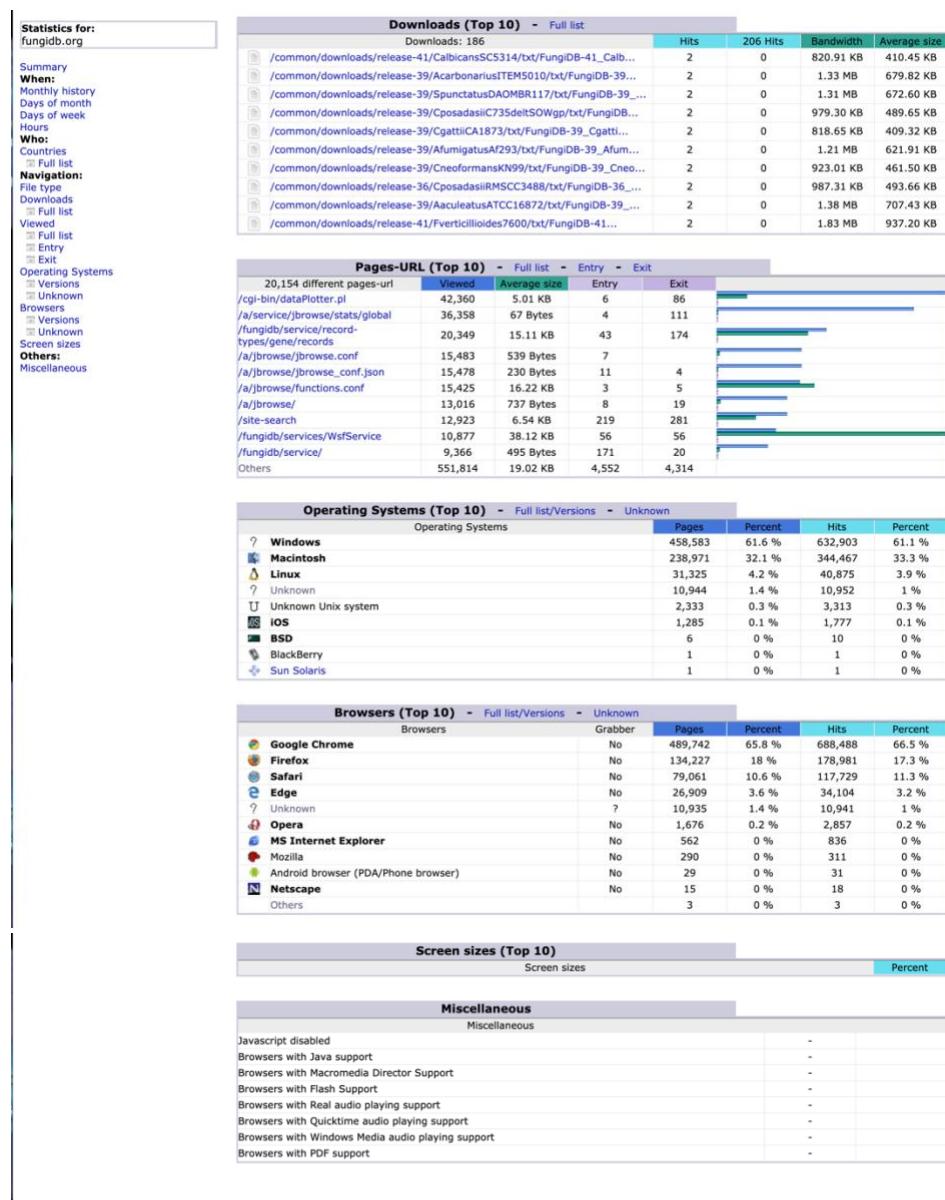
Month	Unique visitors	Number of visits	Pages	Hits	Bandwidth
Jan 2020	3,644	7,822	1,143,872	1,387,925	29.51 GB
Feb 2020	4,429	9,010	847,864	1,044,521	23.63 GB
Mar 2020	3,965	7,925	1,052,667	1,237,807	47.52 GB
Apr 2020	3,807	8,657	1,260,543	1,476,644	27.14 GB
May 2020	4,626	10,120	1,474,243	1,739,026	34.97 GB
Jun 2020	4,296	9,220	1,245,047	1,489,054	30.23 GB
Jul 2020	4,189	9,495	1,417,204	1,679,054	38.86 GB
Aug 2020	3,826	8,711	1,213,322	1,467,604	31.25 GB
Sep 2020	2,423	5,080	743,449	1,034,297	23.66 GB
Oct 2020	0	0	0	0	0
Nov 2020	0	0	0	0	0
Dec 2020	0	0	0	0	0
Total	35,205	76,040	10,398,211	12,554,932	286.77 GB

Days of month



Day	Number of visits	Pages	Hits	Bandwidth
01 Sep 2020	304	56,276	79,716	1.68 GB
02 Sep 2020	333	47,162	65,008	1.68 GB
03 Sep 2020	370	63,500	89,315	2.08 GB
04 Sep 2020	356	43,424	57,688	2.02 GB
05 Sep 2020	181	17,145	22,268	540.29 MB
06 Sep 2020	222	24,581	34,378	798.19 MB
07 Sep 2020	334	43,870	59,084	1.62 GB
08 Sep 2020	336	47,544	65,825	1.99 GB
09 Sep 2020	437	52,328	72,512	1.96 GB
10 Sep 2020	420	71,755	96,542	1.87 GB
11 Sep 2020	344	50,941	69,148	1.20 GB
12 Sep 2020	204	30,954	42,403	770.46 MB
13 Sep 2020	194	20,212	29,278	624.94 MB
14 Sep 2020	379	71,201	103,919	1.80 GB
15 Sep 2020	344	48,749	73,707	1.77 GB
16 Sep 2020	312	53,069	72,418	1.26 GB
17 Sep 2020	10	732	1,088	47.24 MB
18 Sep 2020	0	0	0	0
19 Sep 2020	0	0	0	0
20 Sep 2020	0	0	0	0
21 Sep 2020	0	0	0	0
22 Sep 2020	0	0	0	0
23 Sep 2020	0	0	0	0
24 Sep 2020	0	0	0	0
25 Sep 2020	0	0	0	0
26 Sep 2020	0	0	0	0
27 Sep 2020	0	0	0	0
28 Sep 2020	0	0	0	0
29 Sep 2020	0	0	0	0
30 Sep 2020	0	0	0	0
Average	298	43,732	60,841	1.39 GB
Total	5,080	743,449	1,034,297	23.66 GB





VEuPathDB Glossary

Please also check the [NCBI glossary](#)

- **3-Frame translation (forward)**
Translation of a nucleotide sequence in all three possible reading frames in one direction, usually "on the top [strand](#)" of DNA.
 - **3-Frame translation (reverse)**
Translation of a nucleotide sequence in all three possible reading frames in the reverse direction, usually "on the bottom [strand](#)" of DNA.
 - **AA sequence**

Amino acid sequence.

- **Affymetrix genotyped SNP probes**

Probes on Affymetrix [SNP](#) (single nucleotide polymorphism) arrays, which are used for [SNP genotyping](#). See [Affymetrix microarray technology](#) and www.affymetrix.com

- **Affymetrix microarray technology**

[Microarray](#) manufacturing technology developed by Affymetrix. Combines semiconductor fabrication techniques, solid phase chemistry, combinatorial chemistry, molecular biology, and robotics to generate a photolithographic manufacturing process in which oligonucleotides are synthesized directly on a chip.

See www.affymetrix.com

- **Affymetrix probes**

Probe on an Affymetrix [microarray](#) designed to determine whether or not the complementary sequence of RNA or DNA is present in a sample. Generally 25 nucleotides in length (25-mers), their short length provides higher specificity than longer probes. See [Affymetrix microarray technology](#) and www.affymetrix.com

- **Amitochondriate**

Eukaryotic organism that lacks a [mitochondrion](#). Examples include Giardia and other parasites such as Trachipleistophora and Entamoeba. However, most of these organisms contain what appear to be mitochondrial remnants as well as mitochondrial [genes](#) in their nuclear genomes.

- **Annotation**

Identified feature within a sequence, such as a known or predicted [gene](#), domain, motif, post-translational modification, etc.

- **Annotation density**

Level to which a nucleotide or protein sequence has been annotated.

See [Annotation](#).

- **ApiCyc**

Database/utility on EuPathDB used for searching and visualizing metabolic pathway information for organisms in EuPathDB; derivative database generated by analyzing various genomes (for example from Plasmodium, Cryptosporidium, and Toxoplasma) with SRI International's pathway tools.

- **ApiDots alignments**

Consensus sequences found in the ApiDots database and generated by clustering and assembling Apicomplexan mRNA and [EST](#) sequences. These consensus sequences were subjected to database searches against protein and protein domain sequences.

- **Apicoplast**

Nonphotosynthetic plastid found in almost all protozoan parasites belonging to the phylum Apicomplexa that have been examined. The apicoplast is surrounded by four membranes, giving rise to the theory that its presence in the Apicomplexa is the result of a secondary endosymbiosis (acquired by the engulfment of an ancestral alga and retention of the algal plastid). Similar to other endosymbiotic organelles (mitochondria, chloroplasts), the apicoplast contains its own genome as well as proteins that are encoded in the nucleus and post-translationally imported. The apicoplast is a vital organelle to the parasite's long-term survival.

- **Attribute**

Inherent characteristic or feature; in a database, a data item related to a database object. For example, attributes of [genes](#) can include features such as introns and untranslated regions (UTRs).

- **BLAST**

Basic local alignment search tool, a [sequence similarity](#) search tool used to quickly find local alignments between a [query](#) sequence and sequences in nucleotide or protein databases. Different versions of this search tool are available to match the types of query sequence and database used. See [blastn](#), [blastp](#), [blastx](#), [tblastn](#), and [tblastx](#).

- **Boolean**

System of logical thought developed by George Boole (1815-1864). In Boolean searching, an "and" operator between two words or values (for example, "apple AND orange") generates a search for items in a database containing both of the words or values. Similarly, an "or" operator between two words or values (for example, "apple OR orange") generates a search for items containing either word.

- **CDS**

Coding sequence. Region of nucleotides that corresponds to the sequence of amino acids in a predicted protein and that includes start and stop codons. Unexpressed sequences (for example, the 5'-UTR, the 3'-UTR, and introns) are not included within a CDS. The CDS usually does not correspond to the actual mRNA sequence.

- **Centromere**

Region of the [chromosome](#) or chromosomal structure essential for division and retention of the chromosome within the cell; point of a chromosome where the spindle fibers attach to pull the chromosome apart during cell division.

- **Chromosome**

Macromolecule of DNA constituting the physical organization of DNA in a cell.

- **Coil**

Three-dimensional spiral structure in protein macromolecules.

- **Contig**

Contiguous genomic sequence assembled from overlapping primary sequences representing overlapping regions of a particular [chromosome](#).

- **CryptoCyc**

Database/utility built by analyzing the Cryptosporidium genome with SRI International's pathway tools; used for searching and visualizing Cryptosporidium metabolic pathway information.

- **Curated annotation**

[Annotation](#) made under the supervision of a curator as opposed to a purely computational prediction. Curated predictions often contain combinations of different types of evidence to support the annotation.

- **DNA/GC content**

Content of guanine (G) and cytosine (C) in a fragment of DNA or a genome. Because GC pairs are more thermostable compared to the AT pairs, it was commonly believed that GC content played a vital part in adaptation to high temperatures, a hypothesis that has been refuted. In the genome browser, the DNA/GC content track displays a GC content graph of the reference sequence at low magnifications and the DNA sequence itself at higher magnifications.

- **Dalton**

Unit of mass abbreviated Da and used to express atomic and molecular masses.

- **Deprecated gene**

[Genes](#) with little or no evidence (similarities / expression) that overlapped (or were subsumed by) larger genes for which there was evidence such as protein similarities, expression evidence from [EST alignments](#), SAGE or [proteomics](#) data were marked as deprecated. These genes will likely be removed in a subsequent release (and in GenBank) unless additional evidence is provided indicating they should be moved into the real gene category.

- **EC numbers**

Enzyme Commission numbers. EC numbers constitute the numerical classification scheme for enzymes based on the chemical reactions they catalyze. EC numbers do not refer to the enzymes, but to the reactions they catalyze.

- **EST**

Expressed sequence tag. Short (typically 100-500 base pairs) partial [cDNA](#) produced by single-shot sequencing of a cloned mRNA (cDNA) and often used to identify [gene](#) transcripts.

- **EST alignments**

Alignments of expressed sequence tags (ESTs) with a corresponding genomic region. For example, in ToxoDB you can visualize [EST](#) alignments by clicking on the "View this sequence in the genome browser" link and turning on the EST Alignments track. Useful for identifying intron boundaries.

- **EST clusters**

Groups of homologous, overlapping [EST](#) sequences created to reduce redundancy of the EST database.

- **Expression level**

Level at which an mRNA or protein is present in a sample. Value can be absolute or relative to other mRNA or protein species in the sample.

- **Expression profile**

Pattern of expression of one or more [genes](#) or proteins over time or over a set of experimental conditions (for example, during development or treatment, or as a result of a genetic mutation such as a knockout).

- **Expression profile correlation**

Method for correlation of [gene expression profiles](#) with gene ontology (GO) [annotations](#) developed for the purpose of identifying groups of genes, pathways, and processes reacting in concert to experimental perturbations.

- **Expression timing**

Timing of [gene](#) expression during a developmental, metabolic, regulatory, or other biological process or response.

- **GBrowse**

Interactive genome browser developed by the Generic Model Organism Database (GMOD) project (www.gmod.org) that can be customized to show selected chromosomal features as well as display user-provided [annotations](#).

- **GBrowse track**

In the [GBrowse](#) viewer, a line of data that corresponds to a particular type of genomic information or feature and that is distinguished by a particular shape or color.

- **GLEAN gene**

Predicted [gene](#) sequence generated by GLEAN, an algorithm that integrates different sources of gene structure evidence (for example, gene model predictions, [EST](#) and protein sequence alignments to the genome, and SAGE or peptide tags) to produce a consensus gene prediction in the absence of known genes.

- **GO**

[Gene](#) Ontology project. Collaborative project that has developed three structured, controlled vocabularies (ontologies) to describe gene products in terms of their associated biological processes, cellular components, and molecular functions in a species-independent manner. The use of a consistent vocabulary allows genes from different species to be compared based on their GO [annotations](#).

- **GO component**

[Gene](#) Ontology term used to describe a cellular component, or the location where a gene product may act, rather than physical features of proteins or RNAs. For example, membrane (GO:0016020), extrinsic to membrane (GO:0019898), and integral to membrane (GO:0016021).

- **GO function**

[Gene](#) Ontology term used to describe the molecular function of a gene product, the jobs that it performs, or the "abilities" that it has (for example, transporting compounds, binding to things, holding things together, and changing one thing into another). This is different from the biological processes the gene product is involved in, which involve more than one activity.

- **GO process**

[Gene](#) Ontology term used to describe a biological process, a recognized series of events, or molecular functions associated with a gene product. A biological process is not equivalent to a pathway, though some [GO terms](#) do describe pathways.

- **GO term**

[Gene](#) Ontology term. The building blocks of the Gene Ontology, each term is assigned to one of the three ontologies: molecular function, cellular component, or biological process. Each [GO](#) term consists of a unique alphanumerical identifier, a common name, synonyms (if applicable), and a definition. When a term has multiple meanings depending on species, the GO uses a "sensu" tag to differentiate among them. For example, the enzyme fumarase has the GO term GO:0004333, fumarate hydratase activity (fumarase activity), catalysis of the reaction: (S)-malate = fumarate + H₂O.

- **GPI anchor**

C-terminal post-translational modification of many eukaryotic proteins. The two fatty acids within the glycophosphatidylinositol (GPI) moiety anchor the protein to the outer leaflet of the plasma membrane. GPI-anchored proteins are believed to be involved in signal transduction and immune responses, as well as the pathobiology of many parasites.

- **Gametocyte**

Eukaryotic germ cell that divides by mitosis to generate other gametocytes or by meiosis to generate gametes. Male gametocytes are called spermatocytes, and female gametocytes are called oocytes. Term often used to describe gametocytes of Plasmodium.

- **GenBank protein record**

Protein sequence file in the GenBank database generally derived by [translation](#) of a related nucleotide record.

- **GenPept protein**

Protein record from the GenPept database at the [NCBI](#) GenBank, which contains inferred [translations](#) of [protein-coding](#) sequences.

- **Gene**

Fundamental physical and functional unit of heredity. Ordered sequence of nucleotides located in a particular position on a particular [chromosome](#) that encodes a specific functional product, such as a protein or RNA molecule. A gene may have a number of parts, including the [promoter](#) region, untranslated regions (5' and 3' [UTRs](#)), introns, and exons.

- **GeneDB**

Project developed by the Sanger Institute Pathogen Sequencing Unit (PSU) and aimed at developing and maintaining curated database resources for all projects handled by the PSU. The database is accessible at [www.genedb.org](#)

- **Genetic markers**

Known DNA sequences that can be identified by a simple assay. Generally genetic variations caused by mutation or alterations in loci that can be observed, examples include restriction length polymorphisms (RFLPs), short tandem repeats (STRs), variable number tandem repeats (VNTRs), short DNA sequences surrounding single base-pair changes (single nucleotide polymorphisms, or [SNPs](#)), or longer [microsatellite](#) sequences.

- **Genomic context**

Location of a [gene](#) in the genome, which can influence the expression of the gene and functional interactions of the gene expression products. In this database, genes are depicted on individual gene pages with their surrounding genomic region and [annotations](#).

- **Genotyped SNPs**

Single nucleotide polymorphisms (SNPs) identified during [genotyping](#) of individual organism strains. See [SNP genotyping](#).

- **Genotyping**

Process of determining the genotype of an individual with a biological assay using PCR, DNA sequencing, or hybridization to DNA [microarrays](#) or beads. Provides a measurement of the genetic variation between members of a species.

- **HMM**

Hidden Markov model. Statistical model in which the system being modeled is assumed to be a Markov process with unknown parameters, and the challenge is to determine the hidden parameters from the observable parameters. The extracted model parameters can then be used to perform further analysis, for example for pattern recognition applications.

- **Homolog**

Related by evolutionary descent, either between species (ortholog) or within a species (paralog).

- **Hydropathy**

Hydropathicity. Degree to which a peptide or protein is likely to be soluble in water. Protein hydropathy plots can be useful in predicting [transmembrane domains](#), potential antigenic sites, and regions that are likely to be exposed on the protein's surface.

- **Intergenic region**

Stretch of DNA located between [genes](#) that may act to regulate gene expression.

- **Intersect**

Similar to the [Boolean](#) operator "AND", an action used in the [query](#) history to find items that are common to two query result sets. For example, to search for items that appear in both result set X and result Y, type, "X INTERSECT Y".

- **Join**

Similar to the [Boolean](#) operator "UNION", an action used in the [query](#) history to combine query sets. For example, to combine result sets X and Y, type, "X JOIN Y".

- **JBrowse**

Interactive genome browser ([jbrowse.org](#)) developed by the Generic Model Organism Database (GMOD) project ([www.gmod.org](#)) that can be customized to show selected chromosomal features as well as display user-provided [annotations](#).

- **KEGG map**

Metabolic or regulatory interaction pathway generated by the Kyoto Encyclopedia of [Genes](#) and Genomes (KEGG) or by the use of their tools ([www.genome.jp/kegg](#)).

- **Locus**

Position on a [chromosome](#) of a [gene](#), feature (such as a [telomere](#)), or other chromosomal marker; also, the DNA at that position. Use of this term is sometimes restricted to mean expressed DNA regions.

- **Low complexity**

Pertaining to sequence regions that have an unusually repetitive nature (for example, a protein sequence of low complexity might look like PPTDPPPKKDGGPPL, and a low-complexity nucleotide sequence might be AAATAAAAAAAAAATAAAAAAATTA). Low-complexity regions can create problems in [sequence similarity](#) searching by causing artifactual hits. For this reason, filters are often used to remove low-complexity sequences. Low-complexity regions also contribute to antigenic variation in apicomplexan parasites.

- **Mass spec data**

Mass spectrometry data. Mass spectrometry is an analytical technique used to measure the mass-to-charge ratios of small molecules in several applications, including identification of proteins or peptides. In our databases, the "Identify [Genes](#) by Mass Spec Evidence" [query](#) is used to identify genes that have evidence of protein expression based on mass spec data.

- **Metabolic pathways**

Series of chemical reactions occurring within a cell and often catalyzed by enzymes. In a pathway, a molecule is often changed or modified into another product, which can be stored by the cell, used as a metabolic product, or used to initiate another pathway.

- **Microarray**

Microscopic array of biological molecules (for example, DNA or protein) used to determine the presence and/or amount (referred to as quantitation) of other biomolecules (other proteins, transcripts, etc.) in biological samples.

- **Microsatellite**

Polymorphic [locus](#) in nuclear and organellar DNA that consists of repeating units of 1-4 base pairs in length. Mostly neutral and codominant, microsatellites are used as molecular markers and to study [gene](#) dosage (looking for duplications or deletions of a particular genetic region). Also known as simple sequence repeats (SSRs).

- **Microsatellite map**

Map of [microsatellite](#) locations and linkages on a genome.

- **Mitochondrion**

Organelle responsible for respiration in a eukaryotic cell. Proteins required for mitochondrial function are encoded both in the nucleus and within the smaller mitochondrial genome.

- **Motif search**

Tool used to identify and locate sequence patterns (motifs) in protein and nucleic acid sequences. In our databases, this flexible search can be based on the general characteristics of the pattern and not solely on specific sequences (for example, Cys-[9-11 amino acids]-Cys or Leu-Leu-[basic residue]-Val). This allows the user to [query](#) using previously undescribed motifs.

- **NCBI**

National Center for Biotechnology Information. Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease.

- **Nonredundant protein DB (NRDB)**

Peptide sequence database containing all nonidentical protein sequences from the same species extracted from GenPept, [NCBI](#) RefSeq, Swiss-Prot, and PRF databases. Used for [BLAST](#) protein database searches because its smaller size results in shorter search times and more meaningful statistics.

- **Nucmer**

NUCleotide MUMmer. Part of the MUMmer alignment package, an alignment tool used for the rapid alignment of very large DNA and amino acid sequences.

- **ORF**

Amino acid sequences computed by translating the six frames of raw genomic sequence using the standard genetic code. We save the translated sequences having at least 50 amino acids. The sequences are not annotated nor human reviewed. ORF's do not necessarily begin with methionine residues, but they all terminate with a stop codons.

- **Oligo**

Oligonucleotide. Short sequence of DNA or RNA, typically of 20-70 nucleotides.

- **Oligonucleotide microarray**

Collection of microscopic oligonucleotide spots arrayed on a solid surface by covalent attachment to chemically suitable matrices. Used for expression profiling,

the monitoring of the [expression levels](#) of hundreds or thousands of [genes](#) simultaneously. Probes for oligonucleotide [microarrays](#) are designed to match parts of known or predicted mRNAs.

- **Open Reading Frame**

See [ORF](#)

- **OrthoMCL**

Genome-scale algorithm for grouping orthologous protein sequences. It provides [genes](#) shared by two or more species/genomes and also genes representing species-specific gene expansion families. Therefore, it serves as a utility for automated eukaryotic genome [annotation](#) and phylogenetic profiling.

Available at [orthomcl.org](#)

- **Ortholog**

Same [gene](#) in different organisms; having evolved from the same ancestral [locus](#).

- **Ortholog group**

Orthologous [genes](#) shared by group of organisms; the group of genes can also contain [paralogs](#).

- **Orthology-based phylogenetic profile**

Tool used to find [genes](#) that are present or not present in a desired group of organisms on the tree of life (currently computed for 81 complete genomes). The user has control over the profile and over whether or not genes must be found in any particular group of organisms (for example, in Apicomplexa but not in mammals). Taxa can also be marked as indifferent (for example, it does not matter if the gene is also found in plants). [Ortholog](#) and [paralog](#) relationships are determined using the [OrthoMCL](#) algorithm.

- **PATS**

Neural network analysis tool that identifies amino acid sequences within a [query](#) sequence that are potentially targeted to the [apicoplast](#) matrix of *Plasmodium falciparum*.

- **PDB 3D structure**

Three-dimensional macromolecular structure in the Protein Data Bank (PDB) ([www.pdb.org](#)) obtained by one of three methods: X-ray crystallography (over 80%), solution nuclear magnetic resonance (NMR) (about 16%), or theoretical modeling (2%). A few structures were determined by other methods.

- **PROSITE motif**

Protein sequence pattern or profile derived from multiple alignments of homologous sequences and stored in the PROSITE database ([prosite.expasy.org](#)), an annotated collection of motif descriptors dedicated to the identification of protein families and domains.

- **Paralog**

Related by [gene](#) duplication within a genome; originated by duplication and then diverged from the parent copy by mutation and selection or drift.

- **Pearson correlation**

Pearson Product Moment Correlation, the most common measure of the correlation between two variables. Reflects the degree of linear relationship between two variables and ranges from +1 to -1, with a correlation of +1 indicating a perfect positive linear relationship between variables.

- **Peptide mass fingerprinting**

Analytical technique for protein identification wherein an unknown protein of interest is cleaved into peptides by a protease such as trypsin, and the peptides resulting from this cleavage are analyzed using a mass spectrometric method such as MALDI-TOF or ESI-TOF. The masses derived for the peptides are then compared to a database containing known protein sequences or even to the genome. Computer programs theoretically cut the protein sequences in the database into peptides with the same protease (for example trypsin), and calculate the absolute masses of the peptides from each protein. They then compare the masses of the peptides of the unknown protein to the theoretical peptide masses of each protein encoded in the genome. The results are statistically analyzed to find the best match.

- **Pfam domain**

Conserved protein region in the Pfam database ([Pfam.org](#)), a collection of multiple sequence alignments and hidden Markov models covering many common protein families. The alignments may represent evolutionarily conserved structures that may shed light on protein function. Profile hidden Markov models (profile [HMMs](#)) built from the Pfam alignments can be useful for associating a new protein to a known protein family, even if the homology is weak. Unlike standard pairwise alignment methods (for example, [BLAST](#) and FASTA), Pfam HMMs deal sensibly with multidomain proteins.

- **Phylogeny**

Historical relationships among lineages of organisms or their parts, including their [genes](#).

- **PlasmoAP**

Algorithm/tool that predicts the likelihood that a protein sequence is targeted to the [apicoplast](#). It provides the position of [signal peptide](#) cleavage sites in amino acid sequences if targeting is predicted.

- **PlasmoCyc**

Database/utility built by analyzing the genomes of the Plasmodium species in EuPathDB with SRI International's pathway tools; used for searching and visualizing Plasmodium metabolic pathway information.

- **ProDom**

Database of protein domain families generated from the global comparison of all available protein sequences ([prodom.prabi.fr](#)).

- **Promoter**

Regulatory region of DNA located upstream (towards the 5' region) of a [gene](#) and providing a control point for regulated gene transcription.

- **Protein-coding**

Capable of encoding a protein sequence; generally refers to a sequence of DNA.

- **Proteomics**

The large-scale study of proteins, particularly of the full set of proteins encoded by a genome.

- **Pseudogene**

Defunct relatives of known [genes](#) that have lost their [protein-coding](#) ability or are otherwise no longer expressed in the cell. Although they may have some gene-like features (such as [promoters](#), CpG islands, and splice sites), they are nonetheless

considered nonfunctional due to their lack of protein-coding ability resulting from various genetic disablements (stop codons, frameshifts, or a lack of transcription) or their inability to function as an RNA (such as with [rRNA](#) pseudogenes).

- **PubCrawler**

Free service that scans daily updates to the [NCBI](#) Medline (PubMed) and GenBank databases and alerts users to any relevant updates. Available at pubcrawler.gen.tcd.ie

- **Query**

Sequence or term used in a database search. For example, the sequence submitted for a [BLAST](#) search is the query sequence.

- **RNA predictions**

Predictions of [genes](#) that encode nonprotein-encoding RNA's such as [tRNA](#), snoRNA, [rRNA](#), etc.

- **RefSeq mRNA**

Nonredundant mRNA sequence in the RefSeq database. RefSeq mRNA sequences with an NM_XXXXXX accession are curated sequences and are, therefore, considered more reliable than those with XM_XXXXXX accessions (predicted mRNA sequences).

- **RefSeq noncoding RNA**

Nonredundant noncoding RNA (ncRNA) sequence in the RefSeq database. RefSeq [ncRNA](#) sequences with an NR_XXXXXX accession are curated sequences and are, therefore, considered more reliable than those with XR_XXXXXX accessions (predicted ncRNA sequences).

- **RefSeq protein**

Nonredundant protein sequence in the RefSeq database. RefSeq protein sequences with NP_XXXXXX accessions are curated sequences and are, therefore, considered more reliable than those with XP_XXXXXX accessions (predicted protein sequences).

- **RefSeq**

[NCBI](#) reference sequences. A curated nonredundant collection of sequences representing genomes, transcripts, and proteins as annotated by NCBI (available at www.ncbi.nlm.nih.gov/refseq). The [annotation](#) in these records is often different from the original GenBank submission, which may not be updated every time new information is obtained.

- **Repeat regions**

Sequences present in many identical or highly similar copies in the genome.

- **SAGE tags**

Serial analysis of [gene](#) expression (SAGE) tags. Short (14-nucleotide) sequences found within mRNA, the relative abundance of which indicates the level of expression of the mRNA containing that tag.

- **SNP**

Single nucleotide polymorphism. Small genetic changes or variations that can occur within a DNA sequence, for example when a single nucleotide, such as an A replacing one of the other three nucleotide letters C, G, or T. Most SNPs are found outside of coding sequences, but SNPs found within a coding sequence are more likely to alter the biological function of a protein. SNPs may be synonymous

(generating a conservative change not altering the amino acid sequence) or they can be nonsynonymous and change the amino acid that is encoded.

- **SNP density**

Amount or number of single nucleotide polymorphisms (SNPs) in a region of the genome.

- **SNP genotyping**

Identifying and mapping single nucleotide polymorphisms (SNPs) in an effort to determine the genotype members of a species. [SNPs](#) usually consist of two alleles (where the rare allele frequency is less than 1%), are evolutionarily conserved, and are the most common type of genetic variation. See [Genotyping](#).

- **Scaffolds**

In genomic mapping, a series of [contigs](#) that are in the right order and orientation, but not necessarily connected in one continuous stretch of sequence.

- **Sequence similarity**

Degree of similarity between two or more protein or nucleotide sequences.

- **Signal peptide**

Short (3-60 amino acids) peptide sequence that directs the co-translational import of a protein to certain organelles or for secretion.

- **SignalP**

Program that predicts the presence and location of [signal peptide](#) cleavage sites in amino acid sequences from Gram-positive and -negative prokaryotes and eukaryotes. The method incorporates a prediction of cleavage sites and a signal peptide/nonsignal peptide prediction based on a combination of several artificial neural networks and hidden Markov models (HMMs).

- **Strand**

One half of the DNA helix; string or stretch of covalently linked nucleotides.

- **Subtract**

Similar to the [Boolean](#) operator "NOT", an action used in the [query](#) history to remove items from one result set that occur in another result set. For example, to remove items that exist in set Y from those in set X, type, "X SUBTRACT Y".

- **Syntenic**

Loci located on the same [chromosome](#) but not necessarily linked. For example, [genes](#) that are part of a syntenic group share a common chromosomal location. Also used to refer to conservation of gene order across species.

- **TIGR**

The Institute for Genomic Research, a nonprofit center dedicated to deciphering and analyzing genomes.

- **TM domains**

See [transmembrane domain](#).

- **TWINSCAN gene models**

[Gene](#) models generated using TWINSCAN, a tool that integrates traditional probability models (such as those underlying GENSCAN and FGENESH) with information from alignments between two genomes. TWINSCAN is based on the idea that functional sequences show different patterns of evolutionary conservation than sequences under little selective pressure, such as the central regions of introns.

TWINSCAN is designed for the analysis of high-throughput genomic sequences containing an unknown number of genes.

- **Telomere**

Nucleoprotein complexes that constitute the physical ends of linear eukaryotic [chromosomes](#) and that have important functions, primarily in the protection, replication, and stabilization of the chromosome ends. Telomeres often contain lengthy stretches of tandemly repeated simple DNA sequences composed of a G-rich [strand](#) and a C-rich strand (called terminal repeats). These terminal repeats are highly conserved. Sequences adjacent to the telomeric repeats are often highly polymorphic and rich in repetitive elements (termed subtelomeric repeats); in some cases, [genes](#) have been found in the proterminal regions of chromosomes.

- **TigrScan gene**

[Gene](#) model generated using TigrScan, a gene-finding tool based on the generalized hidden Markov model (HMM) framework, similar to GENSCAN and Genie. It is highly reconfigurable and includes software for retraining.

- **ToxoCyc**

Database/utility built by analyzing the *Toxoplasma gondii* genome with SRI International's pathway tools; used for searching and visualizing *Toxoplasma* metabolic pathway information.

- **Translation**

Synthesis of protein from an mRNA template.

- **Transmembrane domain**

Three-dimensional protein structure that is thermodynamically stable in a membrane. This may be a single alpha helix, a stable complex of several transmembrane alpha helices, a transmembrane beta barrel, a beta-helix of gramicidin A, or any other structure. Transmembrane domains average 20 amino acid residues in length, though they may be much smaller or much longer.

- **Transmembrane protein**

Protein that spans an entire biological membrane.

- **UTR**

Untranslated region. Section of messenger RNA (mRNA) that either precedes (5' UTR) or follows (3' UTR) the coding region and is not itself translated. The UTR contains several regulatory regions, including the polyadenylation (polyA) site in the 3' UTR, sequences involved in the initiation of [translation](#) (in the 5' UTR), and binding regions for proteins and other regulatory molecules in both the 3' and 5' UTR.

- **UniGene**

Project and database at [NCBI](#) aimed at defining [gene](#)-oriented clusters of expressed sequence tags (ESTs). Sets of [ESTs](#) are clustered based on strong sequence homology in an attempt to define a specific, nonredundant cluster for each transcript in a tissue or genome. Each UniGene cluster contains sequences that represent a unique gene in addition to information about the tissue types in which the gene has been expressed and map location.

- **Wildcard character**

Character used to substitute for any other character(s) in a string.

- **Xenolog**

[Gene](#) found in an unrelated species and that is related by gene transfer rather than common vertical descent.

- **blastn**

Version of the basic local alignment search tool (BLAST) used to compare a nucleotide [query](#) sequence against a nucleotide sequence database.

- **blastp**

Version of the basic local alignment search tool (BLAST) used to compare a protein [query](#) sequence against a protein sequence database.

- **blastx**

Version of the basic local alignment search tool (BLAST) used to compare a nucleotide [query](#) sequence translated in all reading frames against a protein sequence database.

- **cDNA**

Complementary DNA. DNA molecule synthesized by the enzyme reverse transcriptase using an mRNA as template.

- **cDNA microarray**

Collection of microscopic [cDNA](#) spots commonly representing single [genes](#) and arrayed on a solid surface (commonly glass slides) by covalent attachment to chemically suitable matrices. Used for expression profiling, the monitoring of the [expression levels](#) of hundreds or thousands of genes simultaneously.

- **ePCR**

Electronic PCR (polymerase chain reaction). Computational procedure used to check for uniqueness in spacing and number of primer binding sites within DNA sequences. Searches for subsequences that closely match a set of PCR primers and that have the correct order, orientation, and spacing to make a PCR product. Used to check the expected length of a PCR product, which can provide information regarding unexpected repetitive sequences.

- **ncRNA**

Noncoding RNA. Any RNA that is not translated into a protein. Includes transfer RNA (tRNA), ribosomal RNA (rRNA), small RNAs such as snoRNAs, microRNAs, siRNAs and piRNAs, as well as long ncRNAs.

- **rRNA**

Ribosomal RNA. Component of the ribosomes, which function in protein synthesis.

- **snRNA**

Small nuclear RNA. Class of small RNA molecules found within the nucleus, transcribed by RNA polymerase II or III, and involved in a variety of important processes such as RNA splicing (removal of introns from hnRNA), regulation of transcription factors (7SK RNA) or RNA polymerase II (B2 RNA), and maintaining [telomeres](#). They are always associated with specific proteins, and the complexes are referred to as small nuclear ribonucleoproteins (snRNP) or snurps. These elements are rich in uridine. A large group of snRNAs known as small nucleolar RNAs (snoRNAs) are small RNA molecules that play an essential role in RNA biogenesis and chemical modification of ribosomal RNAs (rRNAs) and other RNA [genes](#) (tRNA and snRNAs). They are located in the nucleus and the Cajal bodies of eukaryotic cells (the major sites of RNA synthesis).

- **tRNA**

Transfer RNA. Small RNA chain (73-93 nucleotides) that transfers a specific amino acid to a growing polypeptide chain at the ribosomal site of protein synthesis during [translation](#). A three-base region, the anticodon, pairs to the corresponding three-base codon region on the template mRNA. Each type of tRNA molecule can be attached to only one type of amino acid, but because the genetic code contains multiple codons that specify the same amino acid, tRNA molecules bearing different anticodons may also carry the same amino acid.

- **tblastn**

Version of the basic local alignment search tool (BLAST) used to compare a protein [query](#) sequence against a translated nucleotide sequence database.

- **tblastx**

Version of the basic local alignment search tool (BLAST) used to compare the six-frame [translations](#) of a nucleotide [query](#) sequence against the six-frame translations of a nucleotide sequence database.