# VEuPathDB

**Eukaryotic Pathogen, Vector & Host Informatics Resources**

# VEuPathDB Annual Workshop

Pre-workshop module

June 2021

**Note**: the exercises in this pre-workshop module cover some of the basic functionality in VEuPathDB.

AmoebaDB   CryptoDB   FungiDB   GiardiaDB   HostDB   MicrosporidiaDB   OrthoMCL DB

PiroplasmaDB   PlasmoDB   ToxoDB   TrichDB   TriTrypDB   VectorBase

# Site Search

*Note: this exercise uses PlasmoDB.org 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

The site search is located in the header of any VEuPathDB site and is available from every page. The site search queries the databases for your term or ID and returns a list of pages and documents that contain your query term.



1. Enter the word *kinase* in the site search window (arrow in the image below). Then click enter on your keyboard or click on the search icon (square in the image below).

2. The site search returns a categorized list of pages and documents that contain your term. Site search results are summarized by category, with a details panel on the right. Changing the panel on the left will populate the details panel with that list. What is the total number of results with the word kinase? Are all the results genes? Explore the filter panel on the left side of the webpage.



Results summarized in categories

Details panel with information about each item returned

3. Filter the results so that you only view gene results (hint: click on the word *genes* in the *Filter results* section; arrow in image above). 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 Filter Fields section expands to reveal additional filtering options. Select the *Product descriptions* field and Choose *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).

4. How many of the above genes are found in *Plasmodium falciparum* 3D7? How did you find this number? Hint: explore the *Filter organisms* section of the results filter. There is a search option to aid navigation through the organism tree (left) or the tree can be expanded to find the organism of interest (right). Select the correct organism and apply the filter.

5. Export the results to a search strategy. Use the blue *Export as a search strategy* button at the top right-hand side of the results.



6. Return to the site search results page. You can achieve this in two ways: 1. Your previous results and filter settings were preserved and can be accessed by clicking on the 'back to results' arrow in the site search window. 2. Click on your browser's back arrow. Notice that

7. Clear all filters. You can achieve this in two ways: 1. You can click on each of the clear filter options in the filter results panel (boxes below). 2. You can click on the *clear filters option* in the site search window, which serves to Clear All filters.

**1**

Filter results

| | |
|---|---|
| Genome | **Clear filter** |
| Genes | |

Filter Gene fields     **Clear filter**

select all | clear all

| | |
|---|---|
| ☐ Alternate product descriptions | 3 |
| ☐ EC descriptions and numbers | 217 |
| ☐ GO terms | 185 |
| ☐ Orthologs | 158 |
| ☐ PDB chains | 123 |
| ☑ Product descriptions | 138 |
| ☐ PubMed | 123 |
| ☐ Rodent malaria phenotype | 42 |
| ☐ User comments | 51 |

Filter organisms     **Clear filter**

select all | clear all | expand all | collapse all

*Type a taxonomic name* 🔍 ❓

| | |
|---|---|
| ▾ ⊟ Plasmodiidae | 5,829 |
|    ☐ Hepatocystis sp. ex Piliocolobus tephrosceles 2019 | 132 |
|    ▸ ⊟ Plasmodium | 5,697 |

**2**     kinase     [Clear filters] 🔍

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

☑ Hide zero counts

Filter results

| | |
|---|---|
| Genome | |
| Genes | 11,811 |
| Population biology | |
| Popset isolate sequences | 352 |
| Metabolism | |
| Metabolic pathways | 309 |
| Compounds | 80 |
| Data access | |
| Data sets | 1 |
| Searches | 3 |

Filter fields

*Select a result filter above*

Filter organisms

select all | clear all | expand all | collapse all

*Type a taxonomic name* 🔍 ❓

| | |
|---|---|
| ▾ ☐ Plasmodiidae | 11,812 |
|    ☐ Hepatocystis sp. ex Piliocolobus tephrosceles 2019 | 254 |
|    ▸ ☐ Plasmodium | 11,558 |

☐ Hide zero counts

Filter results

| | |
|---|---|
| Genome | |
| Genes | 11,811 |
| Genomic sequences | 0 |
| Organism | |
| Organisms | 0 |
| Transcriptomics | |
| ESTs | 0 |
| Population biology | |
| Popset isolate sequences | 352 |
| Field samples | 0 |
| Metabolism | |
| Metabolic pathways | 309 |
| Compounds | 80 |
| Data access | |
| Data sets | 1 |
| Searches | 3 |
| Instructional | |
| Tutorials | 0 |
| Workshop exercises | 0 |
| About | |
| News | 0 |
| General info pages | 0 |

Filter fields

*Select a result filter above*

Filter organisms

select all | clear all | expand all | collapse all

*Type a taxonomic name* 🔍 ❓

| | |
|---|---|
| ▸ ☐ Haemoproteidae | 0 |
| ▸ ☐ Plasmodiidae | 11,812 |

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.

10. Try searching for a specific gene ID.  Enter the gene ID below in the site search window:    *PF3D7_0310100*



When the query ID has an exact match in the database, the site search returns a card at the top of the details panel for easy access to the gene page.  The site search also returns other pages that contain the query ID.  Click on the Gene ID to go the gene page.

# Exploring the Gene Page

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

**Learning objectives**

<u>Gene pages</u>:
- Become familiar with the information in gene pages
- Navigate to and from the gene pages
- Use the contents section of the gene page
- Interact with gene page subsections

1. **Navigation to the Gene pages**
   For this exercise visit the gene page for TGME49_222020 (phosphoglycerate kinase PGKII).  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.



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?
   - Is this gene protein coding?
   - What do the shortcuts do?



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?
- How many exons does the gene have?
- Does this gene have an available community annotation?
- How long is the transcript? You can determine transcript length by expanding the Transcripts section.



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 cl ick 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**
   a. Expand the "Orthologs and Paralogs in ToxoDB" section.
   b. Select a few genes from the table using the checkbox.
   c. Scroll to the bottom of the table and click on the Run Clustal Omega button.

6. **Exploring the genetic variation section**

| | | | | | TgCatPRC2 | | |
|---|---|---|---|---|---|---|---|
| ☐ | TGVAND_222020 | phosphoglycerate kinase PGKII | Toxoplasma gondii VAND | no | yes | no |
| ☑ | TGVAND_318230 | phosphoglycerate kinase PGKI | Toxoplasma gondii VAND | no | no | no |
| ☑ | TGVEG_222020 | phosphoglycerate kinase PGKII | Toxoplasma gondii VEG | no | yes | no |
| ☐ | TGVEG_318230 | phosphoglycerate kinase PGKI | Toxoplasma gondii VEG | no | no | no |
| ☐ | TGP89_222020 | phosphoglycerate kinase PGKII | Toxoplasma gondii p89 | no | yes | no |

▾ SNPs



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:
   a. 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 (https://TriTrypdb.org) 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

## 1.    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.



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.

## 2. Getting around JBrowse.

a. Use any of the above described JBrowse linking strategies to get to the genome browser.

b. Once in JBrowse examine the following features:

   i. 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?

   ii. 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.

   iii. The **genome view**: this is where the data tracks are displayed.

c. 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.

## 3. **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.



a. Make sure the *Trypanosoma brucei brucei* TREU927 genome is selected from the genome menu.

b. 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.



c. 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?



d. 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).



e. 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? (*Tip1:* 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).

Track   View   Help

100,000    200,000    300,000    400,000    500,000    600,000    Click and drag    900,0

Tb927_02_v5.1    Tb927_02_

1,035,000    1,040,000    1,045,000    1,050,000
Transcripts (UTRs in Gray when available)    1,045,540    in to    1,048,100
region

nRNA    Tb927.2.5800:mRNA    Tb
Tdp1), putative    sedoheptulose-1,7-bisphosphatase    hyp

Tb927.2.5760:mRNA    Tb927.2.5810:mRNA
Flagellar Member 8    Holliday-junction resolvase-like of SPT6/SH2 domain co

## 4.    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).

Select Tracks    Help

My Tracks    Back to browser    Clear All Filters    Contains text  splice    7 matching tracks
Currently Active
Recently Used

| Name | Category | Subcategory | Dataset | Track Type | RNASeq Alignment | RNASeq Strand |
|---|---|---|---|---|---|---|
| Bloodstream and Procyclic for spliced leader transcriptomes (927, 427)(2014) - Splice Sites | Gene Models | Splice Sites | ... | Segments | ... | ... |
| Bloodstream and procyclic form spliced leader transcriptomes (427, Antat) (2010) - Splice Sites | Gene Models | Splice Sites | ... | Segments | ... | ... |
| Curated Poly A Sites from bloodstream and procyclic forms | Gene Models | Poly A Sites | ... | Segments | ... | ... |
| Procyclic form spliced leader transcriptome - Poly A Sites | Gene Models | Poly A Sites | ... | Segments | ... | ... |
| Procyclic form spliced leader transcriptome - Splice Sites | Gene Models | Splice Sites | ... | Segments | ... | ... |
| Spliced Leader and Poly A Sites from bloodstream and procylic forms - Splice Sites | Gene Models | Splice Sites | ... | Segments | ... | ... |
| Unified Spliced Leader Addition Sites | Gene Models | Splice Sites | ... | Segments | ... | ... |

Category
7  Gene Models
Subcategory
2  Poly A Sites
5  Splice Sites
Dataset
7  (no data)
Track Type
7  Segments
RNASeq Alignment
7  (no data)
RNASeq Strand
7  (no data)

500,000   550,000   600,000   650,000   700,000   750,000   800,000   850,000   900,000   950,000   1,000,000   1,050,000   1,100,

Tb927_02_v5.1    Tb927_02_v5.1:1044924..1048082 (3.16 Kb)    Go

1,046,500    1,047,000    1,047,500

Tb927.2.5800:mRNA
sedoheptulose-1,7-bisphosphatase

Tb927.2.5800(1)    Tb927.2.5810(3
Tb927.2.5800(11)    Tb927.2.5810(1
Tb927.2.5800(5)
Tb927.2.5800(1)    Tb9
Tb927.2.5800(12)
Tb927.2.5800(1)

Tb927.2.5800(11) details    ×

Location:    1046195
Gene ID:    Tb927.2.5800
UTR Length: 162
Count:    128.25999999999996
Note:    The overall count is the sum of the count per million for each sample.

| Sample | Count per million |
|---|---|
| T.brucei 427 cBF | 7.32 |
| T.brucei 5-SL-end-enriched cDNA | 62.91 |
| T.brucei 927 PCF | 6.72 |
| T.brucei 927 slBF | 8.15 |
| T.brucei Alba 1+ (29-13 RNAi) | 13.9 |
| T.brucei Alba 1- (29-13 RNAi) | 4.85 |
| T.brucei Alba 3_4+ (29-13 RNAi) | 2.02 |
| T.brucei bloodstream (Lister 427) | 4.66 |
| T.brucei bloodstream long slender (Antat1.1) | 5.6 |
| T.brucei bloodstream short stumpy (Antat1.1) | 4.74 |
| T.brucei curated splice site cDNA | 7.39 |

OK

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?

5. **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 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*? *Crithidia*?

- 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?

**6. 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 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.

## 7. Retrieving data from and uploading your own tracks to JBrowse

a. 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".

23

v. In the next popup window you can keep everything as the default and either save or view the sequence.



b. 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.

# Strategies Tutorial

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

**Learning objectives:**
- Build a multistep strategy
- Use the Text, GO Term, RNA-Seq, and SNP searches
- Combine search results using Boolean operators and the colocation tool
- Transform genes of one organism into their orthologs in another organism
- Infer expression timing from a well-studied organism onto another organism that lacks data.

In this tutorial you will find genes expressed in gametocytes that are likely proteases and have variation in their upstream regions, possibly affecting promoter regions and other gene features. The strategy you build will combine three different searches that query *P. falciparum* data, then transform the *P. falciparum* genes returned by those searches into their *P. vivax* orthologs and look for SNPs in the upstream regions of the *P. vivax* genes. The ortholog transform enables you make inferences about genes in *P. vivax*, an organism with limited functional data, based on existing data in the closely related and well-studied *P. falciparum*. The *P. vivax* genes returned by the search are likely to share two biological properties, proteolytic activity and expression in gametocytes, and show variation in their upstream regions, possibly affecting promoter regions or other gene features.

## Strategies Overview:

The strategy system offers over 100 structured searches that can be combined to produce multi-step strategies. Each search queries a specific data set and **returns a list of IDs** that share the biological characteristic defined by the data.

Searches are accessible from the 'Search For…' menu on the home page and from the 'Searches' dropdown menu in the header of every page. Searches listed under Genes will return a list of gene IDs, while searches listed under 'SNPs' or 'Metabolic Pathways' will return record IDs representing SNPs, or metabolic pathways.

The 5 searches you will use in this tutorial are:

1.  <u>Identify Genes by Text (product name, notes, etc.)</u> –  The search compares your term against the text in the fields that you specify, returning the IDs of gene records that have a match.
2.  <u>Identify Genes by GO Term</u> – Returns genes that have your specified Gene Ontology (GO) Term(s) or ID(s) assigned to them.
3.  <u>Identify Genes based on RNA Seq Evidence</u> – PlasmoDB integrates raw RNA sequencing data from many different experiments and analyzes all data according to the same workflow to produce expression values.  This search returns genes based on their transcript expression as measure by RNA sequencing.
4.  <u>Transform by Orthology</u> – PlasmoDB integrates ortholog profiles from OrthoMCL. The OrthoMCL algorithm clusters proteins into ortholog groups based on BLAST similarity across at 150 genomes that span the tree of life. The transform we perform here will convert a list of genes in one organism to their orthologs in a different organism.  In this case, we will transform a list of *P. falciparum* genes into their *P. vivax* orthologs.
5.  <u>Identify SNPs based on Differences within a Group of Isolates</u> – PlasmoDB integrates whole genome resequencing of isolates and analyzes each isolate for single nucleotide polymorphisms compared to a reference genome.  This search returns SNPs that are shared between all the *P. vivax* isolates that are integrated in PlasmoDB.

## Before we get started… a few words about combining search results:

Each search returns a list of IDs.  When two searches are combined, the two result sets (list of IDs) are merged.  The table shows the 5 options for combining search results.

| Operator | : | Combined Result will contain: |
|---|---|---|
| ○ ◑ 1 INTERSECT 2 | : | IDs in common between the two lists |
| ◉ ◑ 1 UNION 2 | : | IDs from list 1 and list 2 |
| ○ ◑ 1 MINUS 2 | : | IDs unique to 1 |
| ○ ◑ 2 MINUS 1 | : | IDs unique to 2 |
| ○ ⊢—⊣ 1 **Relative to** 2 | : | IDs whose features are near each other (collocated) in the genome |

If the searches return the same type of genomic feature they can be combined using any of the 5 operators (i.e. search 1 returns genes, search 2 returns genes as in screenshot group A below).



However, searches that return different genomic features will yield no results when combined with intersect, union or minus operators.  This is illustrated in screenshot groupings C and D below.  Because genes and SNPs are different genomic features, there are no IDs in the list of genes (Step 1) that are present in the list of SNPs (Step 2). To combine a search that returns genes with a search that returns SNPs, you must use the collocation option (1 relative to 2).  We know the genomic location of each gene and each SNP and the colocation option is designed to return features based on their relative genomic location, i.e. SNPs that are near or within genes.



28

## Build the Strategy:

**Find *P. vivax* genes that are possible proteases, likely expressed during the gametocyte stages and contain SNPs in their upstream regions.** This search strategy employs 4 searches, an ortholog transform and the colocation tool to integrate SNP information.  Steps 1 and 2 return *P. falciparum* proteases using two different lines of evidence – a text search in step 1 and a Gene Ontology (GO) term search in step 2. These searches are combined with a union to obtain a more comprehensive list of possible proteases. Step 3 returns genes with evidence for expression during the gametocyte stages based on RNA sequencing data collected in *P. falciparum*.  Steps 2 and 3 are combined using the intersect operator to produce a list of genes that have BOTH biological properties: these genes are likely proteases with evidence for expression during gametocyte stages.  In the next step, the *P. falciparum* genes returned in the step 3 result are transformed into their *P. vivax* orthologs.  This results in a set of 125 *P. vivax* genes with suspected protease activity and expression in gametocytes based on annotation and experimental evidence from *P. falciparum*, an organism for which more complete annotation and functional genomics data is available. In Step 5 we look for single nucleotide polymorphisms (SNPs) among isolates of *P. vivax* and collocate these SNPs to the upstream regions of the *P. vivax* genes. The final result is a set of 32 *P. vivax* genes that are likely proteases expressed in the gametocyte stage and that have SNPs in their upstream regions. Your strategy should look like this when you are done:



### Step by Step Instructions

1.  **Run a text search using protease as the text term.**

    Identify Genes by Text (product name, notes, etc.):  Using the Text Search, find genes whose records contain the term 'protease'.  To reach the text search, click on the link in the home page 'Search For…' menu. The page opens showing a list of parameters that are needed to query the data.  Every search is loaded with default parameters so that you can click Get Answer and run the search. Change the Text term to 'protease' and click Get Answer to initiate the search. The search results are displayed in the My Strategies section which consists of a strategy panel followed by a filter table and a result table.

    **Navigation:**   >PlasmoDB    >Search for Genes   >Text   > Text (product name, notes, etc.)

## Identify Genes based on Text (product name, notes, etc.)

🔄 Reset values

❓ **Organism**

*46 selected, out of 46*

select all | clear all | expand all | collapse all

[Filter list below...] ▼ ❓

☑ Hepatocystis sp. ex Piliocolobus tephrosceles 2019
▸ ☑ Plasmodium ⬅ **Choose all organisms**

select all | clear all | expand all | collapse all

❓ **Text term (use * as wildcard)**

[Protease] ⬅ **Enter protease**

❓ **Fields**

☑ Alternate product descriptions
☑ EC descriptions and numbers ⬅ **Leave all fields checked. We will use the default setting here.**
☑ Epitopes from IEDB
☑ External links
☑ Gene ID
☑ Gene name or symbol
☑ Gene type
☑ Genomic sequence ID
☑ GO terms
☑ InterPro domains
☑ Metabolic pathways
☑ Names, IDs, and aliases
☑ Notes from annotators
☑ Organism
☑ Ortholog group
☑ Orthologs
☑ PDB chains
☑ Product descriptions
☑ PubMed
☑ Rodent malaria phenotype
☑ Transcripts
☑ User comments

select all | clear all

**Click Get Answer to initiate the search**

[Get Answer]

**Parameters:**

| Organism | : | Default - all |
|---|---|---|
| **Text term (use * as wildcard)** | : | protease |
| **Fields** | : | Default - all |

**Results and strategy:** You created a one-step strategy by running the text search. The strategy returns 4320 genes that are annotated with the word 'protease'. This annotation could appear in any field that you searched. You can analyze this result by exploring the hits. Look at the data in the columns of the result table. You can add more data with the Add Columns button. Clicking a gene ID in the first column will take you to that gene's record page. Please explore your results to see if they make sense. For example, gene product names might contain the word 'protease'. Functional data assigned to the genes (GO terms and EC numbers) may indicate protease activity.

**Add a step choosing to run a search for genes annotated with the biological process gene ontology term – GO:0006508: proteolysis.** Gene Ontology annotations offer a second line of evidence for finding proteases. The ontologies are a controlled vocabulary for describing the molecular function, biological process and subcellular location of a gene product. GO annotations in PlasmoDB were either provided by the sequencing and annotation centers or inferred based on a gene's similarity to protein domains from the InterPro databases. The GO Term search returns a gene if it is annotated with the GO term that you are looking for. Let's use that search to look for genes annotated with GO:0006508: proteolysis. We will union the text search results with our GO term results when we combine the results of the two searches.

**Navigation:** Add Step   >Combine with other Genes   >1 union 2   > A new search   >GO Term

**Text**
**4,320 Genes**

**+ Add a step**

Step 1

Add a step to your search strategy

**Combine** with other Genes

Text
4,320 Genes

Step 2

**Transform** into related records

Text
4,320 Genes

Step 2

Use **Genomic Colocation** to combine with other features

Text
4,320 Genes

Step 2

① Choose *how* to combine with other Ge...

○ 1 INTERSECT 2   ● 1 UNION 2   ○ 1 MINUS 2   ○ 2 MINUS 1

② Choose *which* Genes to combine. From...

● A new search   ○ An existing strategy   ○ My basket

GO

Function prediction
🔍 GO Term
Text
🔍 Text (product name, notes, etc.)

Search for and choose the GO Term search.

**Add Step**

Add Step 2 : GO Term

❷ Organism

*0 selected, out of 46*

select all | clear all | expand all | collapse all

Filter list below...

☐ Hepatocystis sp. ex Piliocolobus tephrosceles 2019 [Reference]
▸ ☐ Plasmodium

select all | clear all | expand all | collapse all

❷ Evidence

☑ Curated
☑ Computed

❷ Limit to GO Slim terms

○ Yes
● No

❷ GO Term or GO ID

Begin typing to see suggestions...

Begin typing to see suggestions to choose from (CTRL or CMD click to select multiple)

❷ GO Term or GO ID wildcard search

N/A

**Run Step**

❷ Give this search a name (optional)

❷ Give this search a weight (optiona...

Which organism is chosen by default for this search?  Click 'select all' to run the search on all

Begin typing Proteolysis and then choose the correct GO term from the list

Click Run Step to initiate the search

32

**Parameters:**

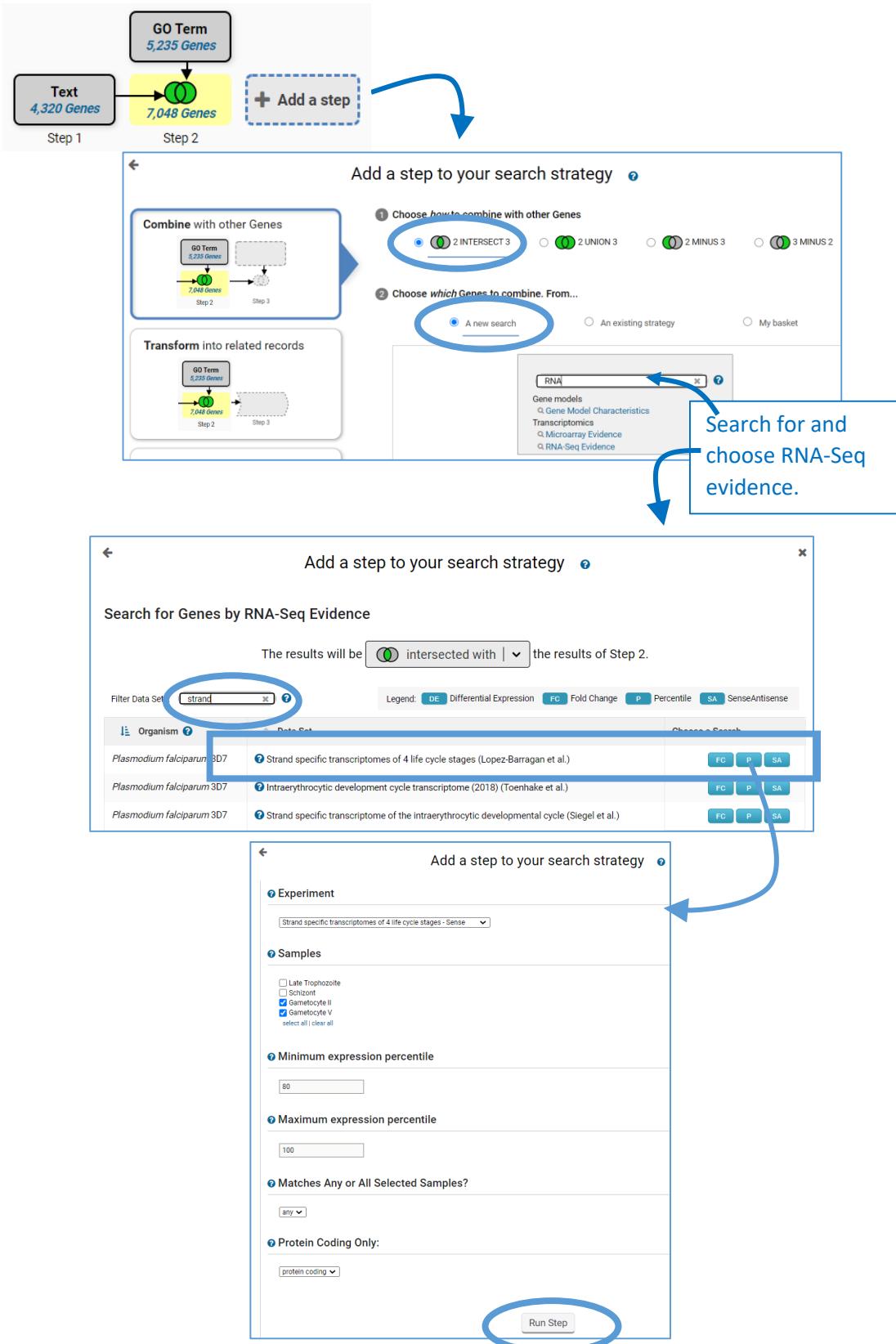| | | |
|---|---|---|
| **Organism** | : | Choose All |
| **Evidence** | : | Default |
| **Limit to GO Slim Terms?** | | Default |
| **GO Term or GO ID** | : | GO:0006508 : proteolysis |
| **Free Text (use '*' for wildcard)** | : | N/A |

**Combine:**



**Strategy Result:** The GO term search returned 5,235 genes annotated with the proteolysis GO term. The union of the text and GO search returns 7,048 genes that are suspected to have proteolytic activity.



2. **Add a step choosing to run a search for genes based on Transcript Expression using RNA-Seq Evidence.** Since PlasmoDB has integrated several RNA sequencing data sets you must first choose what data set (experiment) to search before you are taken to the search form to choose parameters. Use the Filter Data set tool to choose the Percentile search (P) for 'Strand specific Transcriptomes of 4 life cycle stages (Lopez-Barragan et al)'. This data set contains the RNA sequencing analysis of two gametocyte samples. Running the percentile search using the default parameters will return the genes whose expression levels are in the top 20% for those samples.
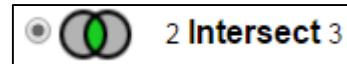
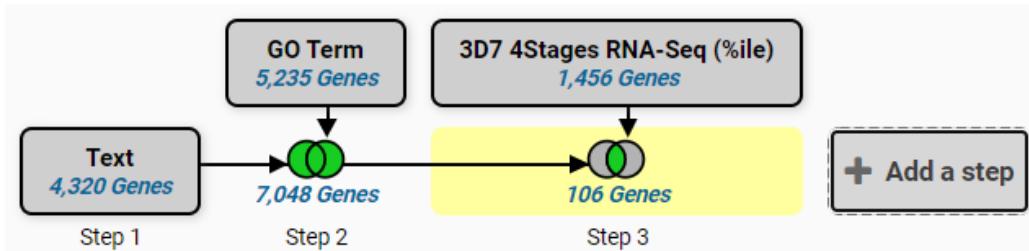**Navigation**: Add Step  >Combine with other Genes  >2 intersect 3  >A new search  >RNA Seq Evidence

Search for and choose RNA-Seq evidence.

**Parameters:**

| Experiment | : | Strand specific transcriptomes of 4 life cycle stages sense strand |
|---|---|---|
| **Samples** | : | Gametocyte II, Gametocyte V |
| **Minimum expression percentile** | : | default |
| **Maximum expression percentile** | : | default |
| **Matches Any or All Selected Samples?** | : | default |
| **Protein Coding Only:** | : | default |

**Combine:** Intersecting this search with the previous result will produce a list of genes that are common to both result lists.
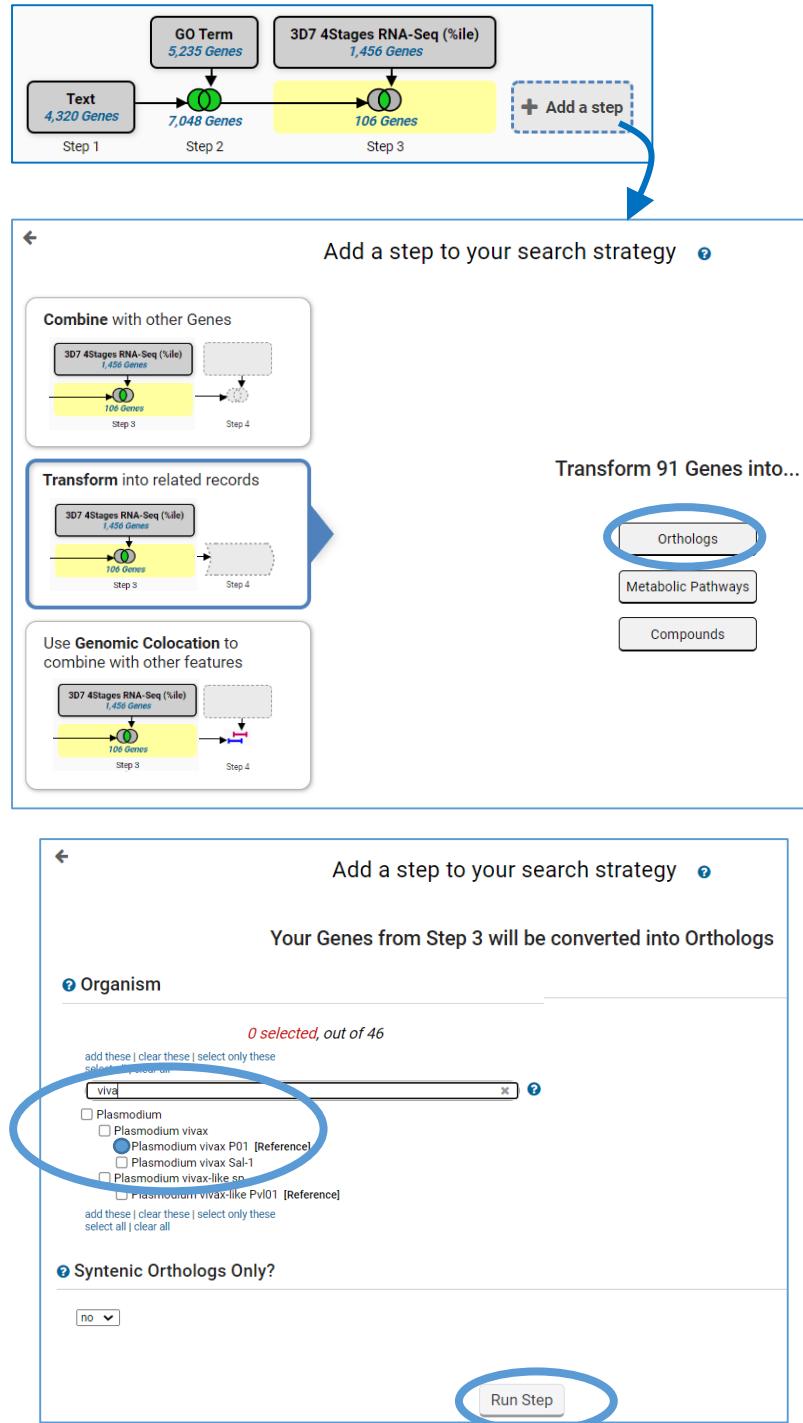


**Strategy result:** We have a three-step strategy that returns 106 *P. falciparum* genes that are suspected proteases with evidence for expression in gametocytes based on RNA Sequencing data. Explore your gene list!!



3. **Add a step to the strategy that transforms the 106 *P. falciparum* genes into *P. vivax* genes.**
   *P. falciparum* is a well-studied organism with active curatorial efforts and large amounts of functional data. For example, PlasmoDB has 18 RNA sequencing and 11 microarray data sets integrated for *P. falciparum*, but only 4 RNA-Seq and 2 microarray for *P. vivax*. A researcher interested in *P. vivax* can take advantage of the *P. falciparum* data by creating a strategy based on *P. falciparum* data to retrieve genes with the biological properties they are interested in, and then transforming the results to their *P. vivax* orthologs.
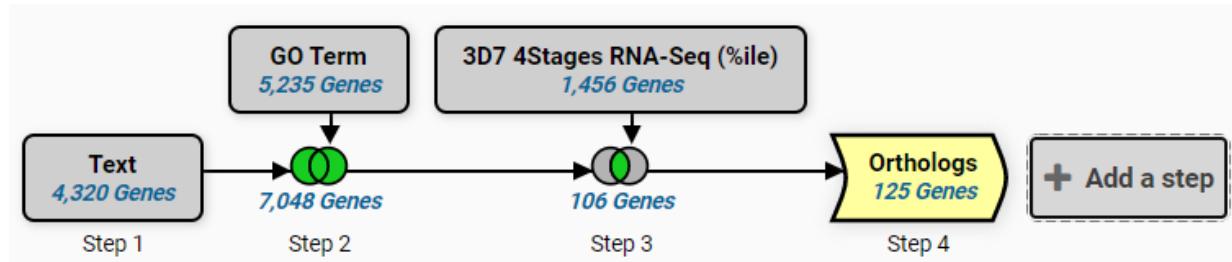
   **Navigation:** >Add Step   >Transform into related records   >Orthologs

**Parameters:** Choose only *P. vivax* P01 in the Organism parameter of the Add Step Popup.

**Combine:** The ortholog transform function does not combine lists, but instead transforms the results into orthologs from a different species.

**Strategy Result:** We have a four-step strategy that returns 125 *P. vivax* genes that are suspected proteases expressed in gametocytes based on *P. falciparum* RNA Sequencing data.



4.  **Add a step to the strategy that returns *P vivax* SNPs and collocate those SNPs to the upstream 1000bp of the *P. vivax* genes in step 4.** We can look for variation (SNPs) associated with the genes from Step 4.  PlasmoDB integrates whole genome resequencing data from many isolates, and PlasmoDB contains 195 data sets from whole-genome sequencing of *P. vivax* isolates. PlasmoDB analyzes the whole genome sequencing reads by aligning them to the reference genome and then examines the genome one base at a time to find bases in the isolate that do not match the reference sequence. The SNPs are loaded in the database along with other information such as how many sequencing reads supported the SNP call and the genomic location of the SNP.  The search we will use analyzes whole genome resequencing data from all *P. vivax* isolates to find SNPs shared between all isolates. You will notice that initiating the search does not immediately bring up the result, but instead leads you to the colocation tool.

    **Navigation**: >Add Step   >Use Genomic Colocation   >A new search   >Differences Within a Group of Isolates

**Parameters:**

| | | |
|---|---|---|
| **Organism** | : | *P. vivax* P01 |
| **Isolates** | : | Default = All Isolates (195) |
| **Read frequency threshold** | : | Default - 80% |
| **Minor allele frequency >=** | : | Default - 0 |
| **Percent isolates with a base call >=** | : | Default - 70 |

**Colocation:** Because this search returns SNPs and not genes, the only option for combining the two result lists is by relative genomic location.  Arrange the statement in the Colocation popup to read: **Return each Gene from step 4 whose upstream 1000bp region overlaps the exact region of a SNP in Step 5 and is on either strand**.  Remember to indicate that you want to locate the SNPs in the upstream region of the gene.

**Strategy: Congratulations!** You have completed the strategy and have a list of 32 *P. vivax* genes that are possible proteases, are likely expressed in gametocytes and have upstream SNPs.

This link will retrieve the completed strategy:
https://plasmodb.org/plasmo/app/workspace/strategies/import/76a3cff6f01535ea