

Finding Genes, Building Search Strategies and Visiting a Gene Page

- Finding a gene using text search.
For this exercise use <http://www.plasmodb.org>

- Find all possible kinases in *Plasmodium*.

Hint: use the keyword “kinase” (without quotations) in the “Gene Text Search” box.



- How many genes did you get? Does this include all genes with the word kinase in PlasmoDB? (*Hint* – the filter table is located between the strategy panel and the result table and shows the distribution of results across the organisms that you searched. Click on a number to only display results from a specific species or strain).

My Strategies: [New](#) [Opened \(1\)](#) [All \(1\)](#) [Basket](#) [Public Strategies \(26\)](#) [Help](#)

(Genes) Strategy: Text *

[Text](#) [Add Step](#)

3617 Genes Step 1

3617 Genes from Step 1
Strategy: Text

Click on a number in this table to limit/filter your results

All Results	Ortholog Groups	Plasmodium																					
		<i>P.bergherei</i>	<i>P.chabaudi</i>	<i>P.coatneyi</i>	<i>P.cynomolgi</i>	<i>P.falciparum</i> (nr Genes: 376)		<i>P.fragile</i>	<i>P.gaboni</i>	<i>P.gallinaceum</i>	<i>P.inui</i>	<i>P.knowlesi</i>	<i>P.malariae</i>	<i>P.ovale curtisi</i>	<i>P.reichenowi</i>	<i>P.relictum</i>	<i>P.vinckei</i> (nr Genes: 302)		<i>P.vivax</i> (nr Genes: 339)		<i>P.yoelli</i> (nr Genes: 472)		
		ANKA	chabaudi	Hackeri	strain B	3D7	IT	strain nilgiri	strain SY75	8A	San Antonio 1	strain H	UG01	GH01	CDC	SGS1-like	petteri strain CR	vinckei strain vinckei	P01	Sal-1	yoelli 17XNL	yoelli 17X	yoelli 17Y
3617	238	158	158	160	153	194	182	152	184	167	152	162	161	163	193	165	151	151	163	176	155	159	158

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

Genes: 3617 Transcripts: 3628 Show Only One Transcript Per Gene

First 1 2 3 4 5 Next Last [Advanced Paging](#) [Download](#) [Add to Basket](#) [Add Columns](#)

Gene ID	Transcript ID	Organism	Genomic Location (Gene)	Product Description	Found in	Score
PKNH_0407600	PKNH_0407600.1	<i>P. knowlesi</i> strain H	PKNH_04_v2:339,792..341,414(-)	protein kinase 7, putative	User Comments, GO Terms, Product, InterPro	23
PVX_002865	PVX_002865.1	<i>P. vivax</i> Sal-1	Pv_Sal1_chr04:336,029..337,650(-)	serine/threonine protein kinase, putative	User Comments, GO Terms, Product, InterPro	23

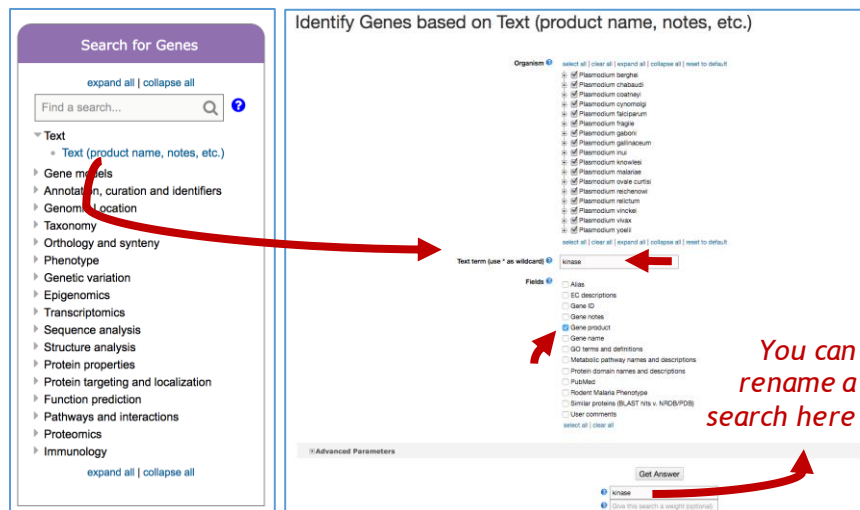
- Look closely at the sections of the result page. How many of those are in *P. falciparum*? How did you find this out?

- What happens if you search using the term **kinases** in the Gene Text Search box? How many results are returned?

b. Find only the kinases that specifically have the word “kinase” in the gene product name.

The search you ran in step 1a using the Gene Text Search box initiates a preconfigured search. Initiating the search from the full text search form - **Identify Genes based on Text**, allows you to configure the search yourself, choosing parameters that best meet your needs. Use the search form to search for genes that have the word kinase in their **gene product** name/description. Note that you can also revise the search from step 1a and configure the search parameters as described below.

- There are several ways to navigate to the **Identify Genes based on Text** page. Notice the sections of the search page. At the top are parameters and the Get Answer button followed by a search description and a list of datasets used by the search.



- How can you make sure to find your text term in plural form or in compound words like “kinases” or “6-phosphofruktokinase”. Adding a wildcard (wildcard = an asterisk (*), means any character) in your search term will broaden your search. Use the full text search, the specific page where you can define the fields to be searched (Fields = Gene Product).

Try kinase *kinase *kinase*

- Give each new search a name to help you keep track of the searches.
- How did you get to the Text Search page?
- How does limiting the number of fields searched affect your results?
- Did you remember to use the wild card?
- How many genes have the word kinase in their product names?

c. Combine the results of two text searches.

Find genes that were identified using the key word ***kinase*** but not the word **kinase**?

- Here we will build a search strategy that combines 2 of your searches. If you are not displaying the results of the ***kinase*** search (the strategy box will be highlighted in yellow), return to it by clicking on that step box in the

The screenshot shows a search strategy builder interface. At the top left, a yellow box labeled 'Step 1' contains the search '*kinase*' with 2871 Genes. To its right is a red 'Add Step' button. Below this, a search strategy list is shown with options like 'Gene', 'Genomic Segment', and 'SNP'. A second 'Add Step' dialog is open, titled 'Add Step 2 from existing strategy:'. It shows a search for 'kinase' with 2689 Genes. Below this, a section titled 'Combine Genes in Step 1 with Genes in Step 2:' offers several options: '1 Intersect 2', '1 Union 2', '1 Minus 2', '2 Minus 1', and '1 Relative to 2, using genomic colocation'. The '1 Minus 2' option is selected. A red callout box points to the search names 'kinase*' and 'kinase**'. Another red callout box points to the '1 Minus 2' option, asking 'Which operation will return genes from step 1 (*kinase*) but not step 2 (kinase)?'. At the bottom, a summary shows Step 1 (*kinase* 2871 Genes) and Step 2 (Copy of kinase 2689 Genes) combined to result in 182 Genes.

strategy panel. To add your **kinase** search to this strategy, click on “Add Step” and select “existing strategy”.

- Select the right strategy from your list of Gene Strategies and combine the strategies with the correct operation. Notice that there is an extra asterisk at the end of an unsaved strategy name. The list of available searches will have an * at the end of the name.
- Do the results make sense? Do all the product names contain the word **kinase**? From the result page look at the table of gene IDs returned by the search. The Product Description column contains the gene product name.

2. Combining text search results with results from other searches

a. Find kinase genes that are likely secreted.

In exercise 1b. you identified genes that have the word **kinase** somewhere in their gene product name (searching ***kinase*** in gene product field). Grow your search strategy by adding a step that returns genes whose protein products are predicted to have a signal peptide. In this search you are querying the results of our genome-wide analysis that used the SignalP

program to predict the presence and location of signal peptide cleavage sites in amino acid sequences.

<http://www.cbs.dtu.dk/services/SignalP/>

Focus your Strategies section on the ***kinase*** search and click Add Step. For the second search choose **Identify Genes based on Protein targeting and localization, Predicted Signal Peptide**

- How did you combine the search results?
- How many kinases are predicted to have a signal peptide?

Operator	:	Combined Result will contain:
<input type="radio"/> 1 INTERSECT 2	:	IDs in common between the two lists
<input checked="" type="radio"/> 1 UNION 2	:	IDs from list 1 and list 2
<input type="radio"/> 1 MINUS 2	:	IDs unique to 1
<input type="radio"/> 2 MINUS 1	:	IDs unique to 2
<input type="radio"/> 1 Relative to 2	:	IDs whose features are near each other (colocated) in the genome

The screenshot shows the 'Add Step' dialog box in a bioinformatics software. The 'Add Step' dialog has a search for 'Predicted Signal Peptide' and a list of genes from the Plasmodium species. The 'Combine Genes in Step 1 with Genes in Step 2' section shows the '1 INTERSECT 2' option selected. A red box highlights the 'Intersect' option with the text 'Which operation will return genes that are in both search result'.

Below the dialog box, the search strategy is shown as a flowchart:

- Step 1: *kinase* 2871 Genes
- Step 2: Signal Pep 20722 Genes
- Result: 154 Genes

b. Now that you have a list of possible secreted kinases, expand this strategy even further.

There is no wrong answer here!!

- From a biological standpoint what else would be interesting to know about these kinases? Add more searches to grow this strategy. Open the categories under Identify Genes By: on the home page and explore the types of searches that are available. You can reduce (or expand) your result set by adding searches that are based on many types of data.
- For example, how many of the secreted kinases also have transmembrane domains?

c. In the above example, how can you define kinases that have either a secretory signal peptide AND/OR a transmembrane domain(s)?

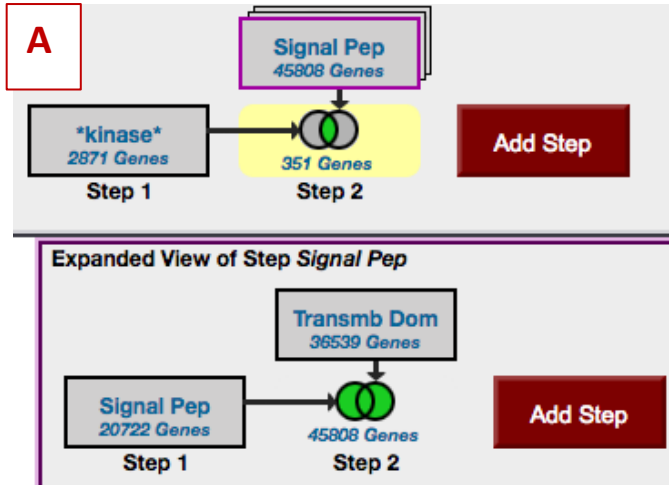
Hint: to do this properly you will have to employ the “Nested Strategy” feature. Nesting a strategy allows you to control the order in which your result sets are combined. Think about the difference between two mathematical equations.

The image displays a bioinformatics tool interface with a nested search strategy. It consists of three main parts:

- Strategy Overview:** A flowchart showing two steps. Step 1 is a search for "*kinase*" resulting in 2871 Genes. Step 2 is a search for "Signal Pep" resulting in 20722 Genes. The intersection of these two searches is highlighted in yellow and labeled "154 Genes".
- Expanded View of Step Signal Pep:** A detailed view of the Step 2 search. It shows a search box containing "Signal Pep" (20722 Genes) and an "Add Step" button. Below this, a list of organisms is displayed, including *Plasmodium berghei*, *Plasmodium chabaudi*, *Plasmodium coactroyi*, *Plasmodium cynomolgi*, *Plasmodium falciparum*, *Plasmodium fragile*, *Plasmodium gaboni*, *Plasmodium gallinaceum*, *Plasmodium inui*, *Plasmodium knowlesi*, *Plasmodium malariae*, *Plasmodium ovale*, *Plasmodium rechenowi*, *Plasmodium relictum*, *Plasmodium SGS1-like*, *Plasmodium vinckei*, *Plasmodium vivax*, and *Plasmodium yoelii*. Search parameters are listed: Minimum SignalP-NN Confusion Score: 3, Minimum SignalP-NN D-Score: 0.5, Minimum SignalP-NN Signal Probability: 0.5, and Matches any of all advanced parameters: any. The results are 20722 Genes.
- Search Results Window:** A window titled "STEP 2: Signal Pep" showing the same list of organisms and search parameters as the expanded view. A red circle highlights the "Make Nested Strategy" button in the window's title bar.

Red arrows indicate the flow of information: one arrow points from the "Signal Pep" box in the Strategy Overview to the "STEP 2: Signal Pep" window, and another arrow points from the "Make Nested Strategy" button in the window back to the intersection of the two steps in the Strategy Overview.

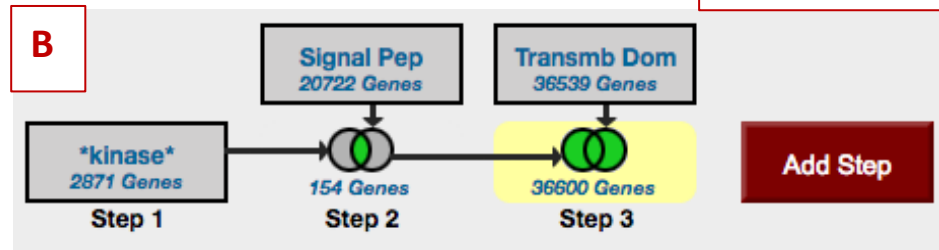
Equation without nesting: $2 \times 3 + 5 = 11$
 Equation with nesting: $2 \times (3 + 5) = 16$



Strategy Logic:

Every gene returned by Strategy A will be a kinase. These kinases with have a signal peptide OR a TM domain OR both. (SP and/or TM; either or both)

Strategy B returns kinases that have a signal peptide as well as TM domain containing genes.



3. Finding a gene by BLAST Similarity.

Note: For this exercise start with <http://www.toxodb.org>

Imagine that you generated an insertion mutant in *Toxoplasma* that is providing you with some of the most interesting results in your career! You sequence the flanking region and you are only able to get sequence from one side of the insertion (the sequence shown below). You immediately go to ToxoDB to find any information about this sequence. What do you do?

```
aaaggagagaaagataaaaatatacaaaggtcccagagacacgatagtgttactgacaa  
catacagaatcaggtcgagcaatggaagaaccaagcaccggcgccaagagattgaaactcgc  
ttggattgccgtagcgtttatgagttgatagcttggtctctaaaaaacaaggctgaaaa  
atggaaaaaaatgtctccaat
```

- Sequence is also available from this URL:
<http://tinyurl.com/ex1blast>
- Try using the BLAST search with this sequence

The image displays three screenshots from the ToxoDB website, illustrating the navigation path to the BLAST search tool. Each screenshot has a green header bar.

- Search for Genes:** Shows a search bar and a list of categories. The 'Sequence analysis' category is expanded, and 'BLAST' is highlighted with a red arrow.
- Search for Other Data Types:** Shows a search bar and a list of categories including Popset Isolate Sequences, RFLP Genotype Isolates, Genomic Sequences, Genomic Segments, SNPs, ESTs, ORFs, Metabolic Pathways, and Compounds.
- Tools:** Shows a list of tools. 'BLAST' is highlighted with a red arrow, and a red arrow points from the 'BLAST' link in the 'Search for Genes' screenshot to this 'BLAST' tool.

- Which blast program should you use? (hint: try different Blast programs, just keep in mind that you have a nucleotide sequence so you have to use an appropriate BLAST program).

1. Choose your target data type. What type of sequence in the database do you want to match your sequence to?
2. Choose the BLAST program to use.
3. Choose the target organism. What genome do you want to match your sequence to?

Target Data Type Transcripts
 Proteins
 Genome
 EST
 ORF
 PopSet

BLAST Program blastn
 blastp
 blastx
 tblastn
 tblastx

Target Organism

- Cyclospora
- Eimeria
- Hammondia
- Neospora
- Sarcocystis
- Toxoplasma
 - Toxoplasma gondii ARI
 - Toxoplasma gondii FOU
 - Toxoplasma gondii GAB2-2007-GAL-DOM2
 - Toxoplasma gondii GT1
 - Toxoplasma gondii MAS
 - Toxoplasma gondii ME49
 - Toxoplasma gondii RUB
 - Toxoplasma gondii TgCatPRC2
 - Toxoplasma gondii VAND
 - Toxoplasma gondii VEG
 - Toxoplasma gondii p89

Input Sequence

Note: only one input sequence allowed.
maximum allowed sequence length is 31K bases.

Expectation value

Maximum descriptions/alignments (V=B)

Low complexity filter

Note on BLAST programs:

- blastp compares an amino acid sequence against a protein sequence database;
 - blastn compares a nucleotide sequence against a nucleotide sequence database;
 - blastx compares the six-frame conceptual translation products of a nucleotide sequence (both strands) against a protein sequence database;
 - tblastn compares a protein sequence against a nucleotide sequence database dynamically translated in all six reading frames (both strands);
 - tblastx compares the six-frame translations of a nucleotide sequence against the six-frame translations of a nucleotide sequence database.
- Are you getting any results from blastx? tblastn? What about blastn?
 - What is your gene? (hint: after running a blastn against *Toxoplasma* ME49 (Target organism) genomic sequence (Target Data Type), click on the “link to the genome browser”. In the genome browser zoom out to see what gene is in the area).

4. Viewing data on a gene page.

Note: For this exercise use <http://plasmodb.org/>

a. Find the gene page for cysteine-tRNA ligase (PF3D7_1015200).

- How did you navigate to this gene? What other ways could you get there?
- Examine the information at the top of the gene page:
 - What is the gene name?
 - What chromosome is this gene on?
- Explore the “shortcuts” section at the top of the gene page – try clicking on the magnifying glass. This option opens up a preview of various sections of the gene page for quick access.

Add to basket Add to favorites Download Gene

PF3D7_1015200 cysteine--tRNA ligase, putative

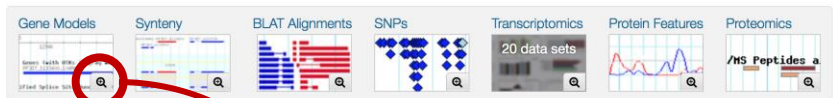
Name: CysRS
Type: protein coding
Chromosome: 10
Location: PF3D7_10_v3:614,872..617,736(-)

Species: Plasmodium falciparum
Strain: 3D7
Status: Curated Reference Strain

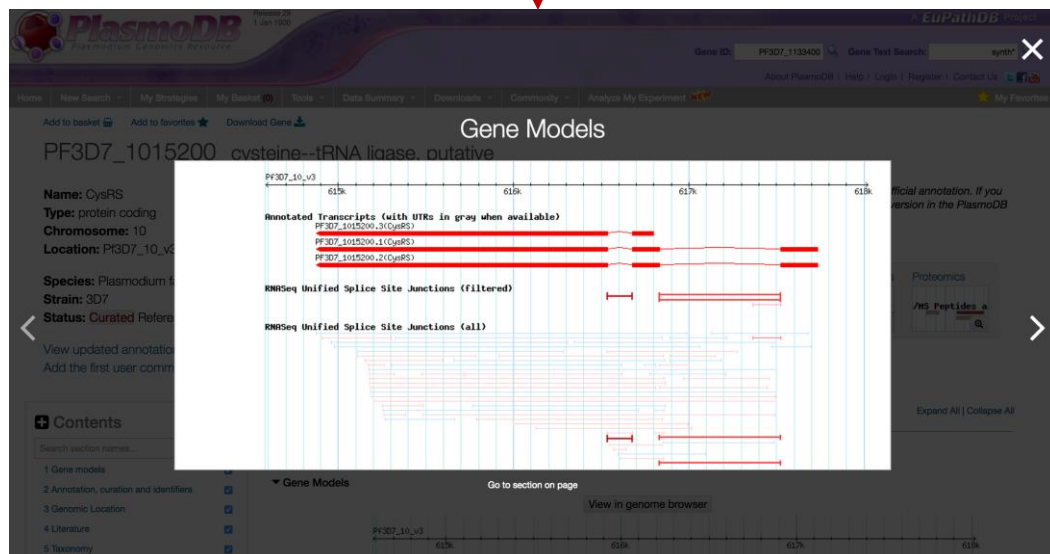
View updated annotation at GeneDB
Add the first user comment

GeneDB curates, researches and improves this genome, and will incorporate appropriate User Comments into the official annotation. If you wish to publish whole genome or large-scale analyses, please contact the primary investigator or use the published version in the PlasmoDB version 5.3 download folder.

Shortcuts



Also see PF3D7_1015200 in the Genome Browser or Protein Browser



- Examine the “Gene Models” section of the gene page.
 - How many exons does this gene have?
 - How many transcripts does this gene encode?
 - What direction are the transcripts relative to the chromosome?

- What does the “splice site junctions” information mean?

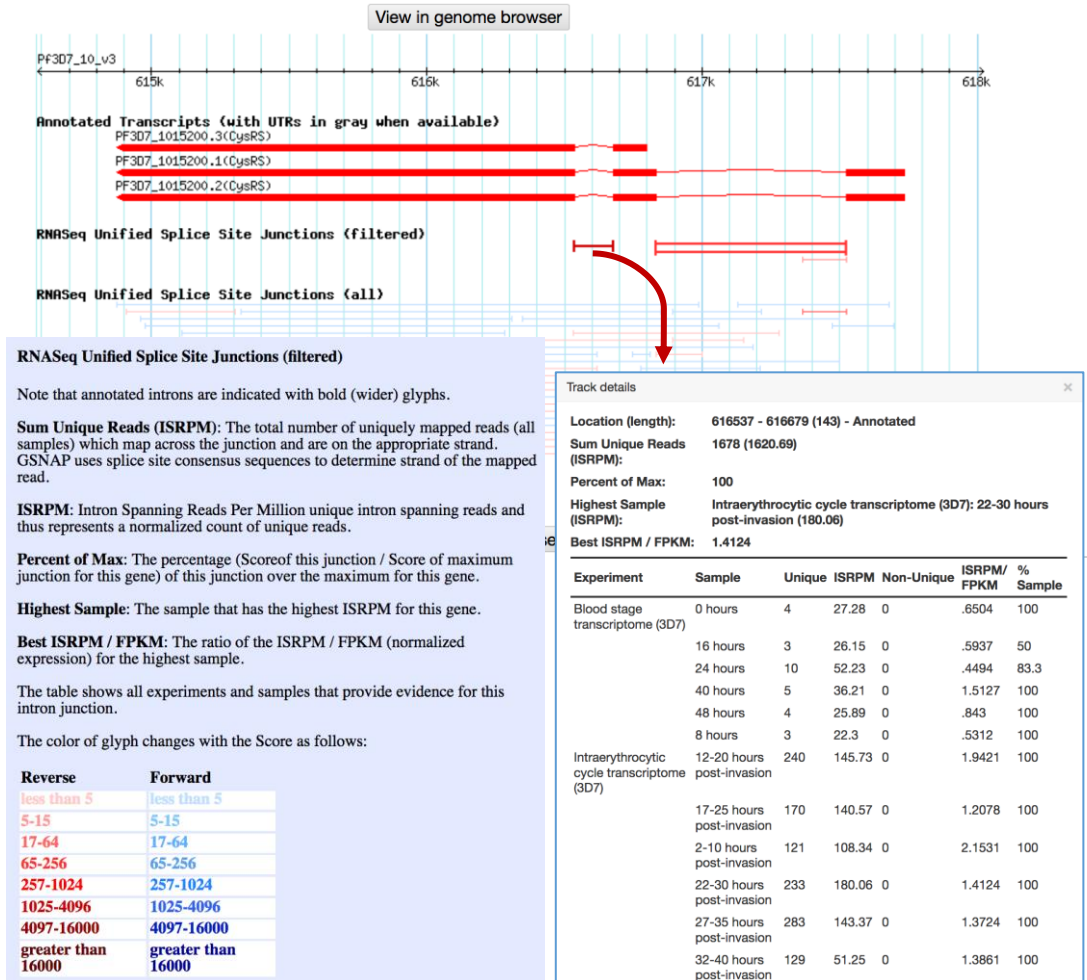
1 Gene models

Expand All | Collapse All

Exons in Gene 5

Transcripts 3

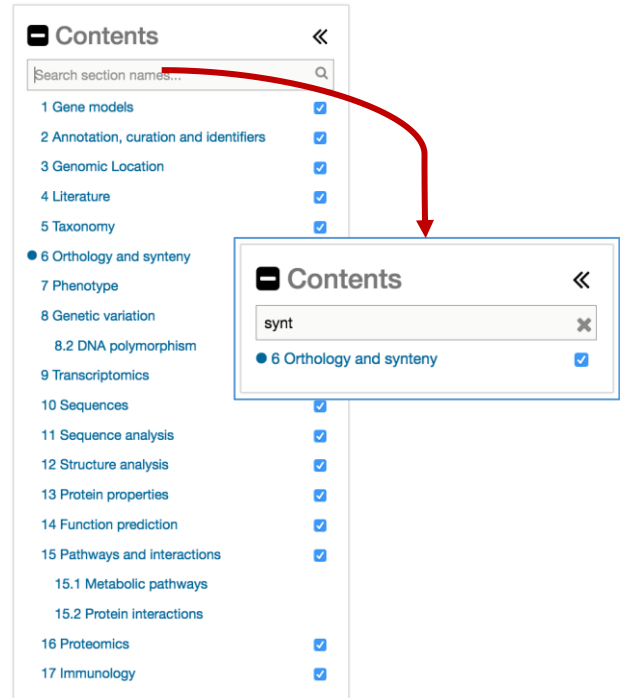
▼ Gene Models



- How many nucleotides is the largest transcript? (hint: examine the transcripts table underneath the gene models).

b. What does the synteny of this gene look like? How did you find/navigate to this section? (hint: you can use the “Contents” menu on the left side of the gene page to find/navigate to the different sections).

- Is synteny (chromosome organization) in this region maintained in other species? Hint: compare gene organization between the different species in the synteny section.
- What does the shading between genes indicate?
- What does synteny look like across the entire chromosome? To do this:
 - Click on the “View in Genome Browser” button right under the synteny section on the gene page.



View in genome browser

- Zoom out to the entire chromosome. There are a few ways to do this. For example, drag your cursor across the entire chromosome in the Overview panel and then select “zoom” from the popup menu (this may take a minute to load).
- For each genome notice that there are two lines: one called genes and the other contig. Which genome is composed of the most fragments? Are there any other interesting observations you can support by looking at synteny over large genomic regions?

c. Does this gene contain Single Nucleotide Polymorphisms (SNPs)?

In gene pages, SNPs are represented in a section called “Genetic variation”. This section includes an isolate alignment tool that shows SNPs between chosen isolates and a DNA polymorphism browser textual and graphical SNP representation.

- Examine the DNA polymorphism section.
 - What is the total number of SNPs in the gene?
 - How many SNPs impact the predicted protein sequence?
 - Is this likely to define the full spectrum of sequence variation in this gene?
 - What do the different color diamonds in the browser view signify? (Hint: move your cursor over a diamond – without clicking - to get more information in a popup).
- Compare Specific isolates to each other:

- Using the isolate alignment tool, run an alignment between several isolates: 303.1, 383.1, 7G8, GB4, N011-A, O222-A, PS097, PS206_E11, RV_3635, RV_3675
- What do Ns indicate?

```

Pf3D7_10_v3 600512 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATATA
303.1 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATATA
383.1 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATANN
7G8_2 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATATA
GB4 600511 NNNNNNNNNN NNNNNNNNNN NNNNNNNNNN NNNNNNNNNN NNNNNNNNNN NNNNNNNNNN NNNNNNNNNN NNNNNNNNNN
N011-A 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATATA
O222-A 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATATA
PS097 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATATA
PS206_E11 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGNNNNNT TCTATATAAT TATATATATA
RV_3635 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATATA
RV_3675 600511 AAATATGTTT AATAAGTTGA AATTTTGTAA TTTATGAAAA TATTTTTTTC TTAGGAACTA TCTATATAAT TATATATATA

```

d. Is this gene expressed at the protein and/or transcript level?

Look at the gene page sections entitled “Proteomics” and “Transcriptomics”.

- What kinds of data in PlasmDB provide evidence for protein expression? (Hint, view the Mass Spec.-based Expression Evidence table).
- Is this gene expressed at the protein level in salivary gland sporozoites?
- Does it contain any post-translational modifications?
- Can you quickly link to the data set record for proteomics experiments?

Mass Spec.-based Expression Evidence [Data sets](#)

post Showing 3 of 15 rows

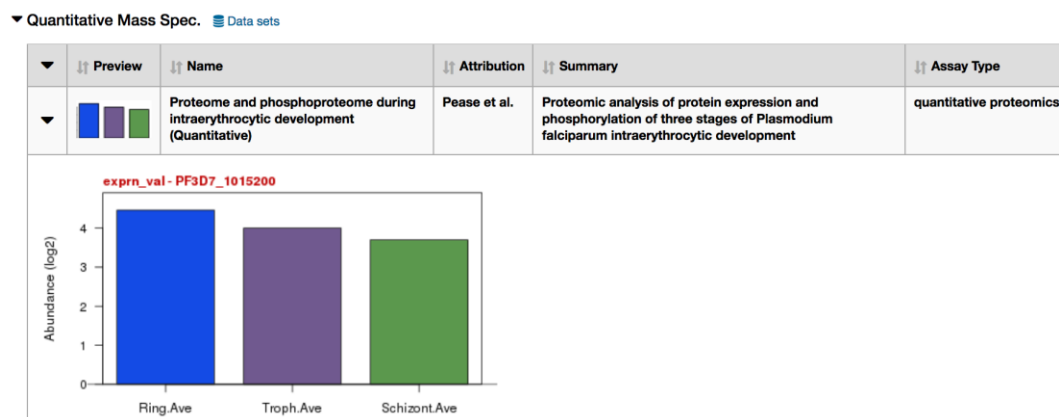
Transcript ID(s)	Experiment	Sample	Sequences	Spectra
PF3D7_1015200.1, PF3D7_1015200.2, PF3D7_1015200.3	Parasite rupture from erythrocyte (D10)	schizont, 42 h post-infection	15	26
PF3D7_1015200.1, PF3D7_1015200.2, PF3D7_1015200.3	Parasite rupture from erythrocyte (D10)	schizont, 48 h post-infection	17	32
PF3D7_1015200.1, PF3D7_1015200.2, PF3D7_1015200.3	Sporozoite Maturation Proteome (NF54)	salivary gland sporozoites 18-22 days post-infection	12	12

is
are

- How abundant is this protein? How confident you of this analysis? Abundance can

be estimated by counting the number of spectra supporting a peptide. Where do you find information about the number of spectra?

- Is the protein more abundant in the ring or schizont life cycle stage? Hint: Scroll down to open the quantitative Mass Spec track called **Proteome and phosphoproteome during intraerythrocytic development (Quantitative)**.



- Does the proteomic data agree with the available transcriptomic data? (Hint, navigate to the transcriptomic section – remember you can use the contents table on the left side of your screen).
- Find the RNAseq experiment by Otto et al. Where is this gene most highly expressed? How did you find this experiment? (Hint, you can search the transcriptomic table with key words).

10 Transcriptomics

10 Transcript Expression Data sets

otto Showing 1 of 22 rows

Preview	Name	Summary	Attribution	Assay Type
	Blood stage transcriptome (3D7)	Transcriptome analysis of <i>P. falciparum</i> 3D7 at seven time points during the intraerythrocytic developmental cycle	Otto et al.	RNA-seq

fpkm - PF3D7_1015200

Time (hours)	FPKM
0	~45
8	~45
16	~45
24	~115
32	~55
40	~25
48	~35

Attribution	Assay Type
Otto et al.	RNA-seq

percentile - PF3D7_1017300

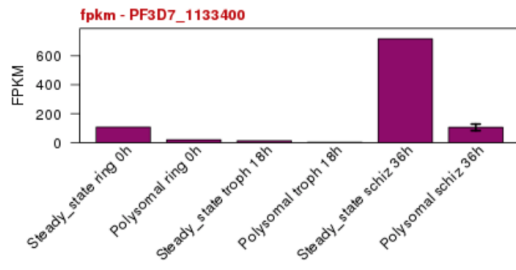
Time (hours)	Percentile
0	~65
8	~65
16	~65
24	~65
32	~65
40	~65
48	~65

- How does the RNAseq data compare with the microarray data?
- What does the polysomal RNAseq data look like?

▼ Transcript Expression [Data sets](#)

poly Showing 1 of 22 rows

▶	↕ Preview	↕ Name	↕ Summary	↕ Attribution	↕ Assay Type
▼		Polysomal and steady-state asexual stage transcriptomes	Transcriptome and translome of the <i>P. falciparum</i> asexual cell cycle.	Bunnik et al.	RNA-seq



▶ Coverage

▶ Data table