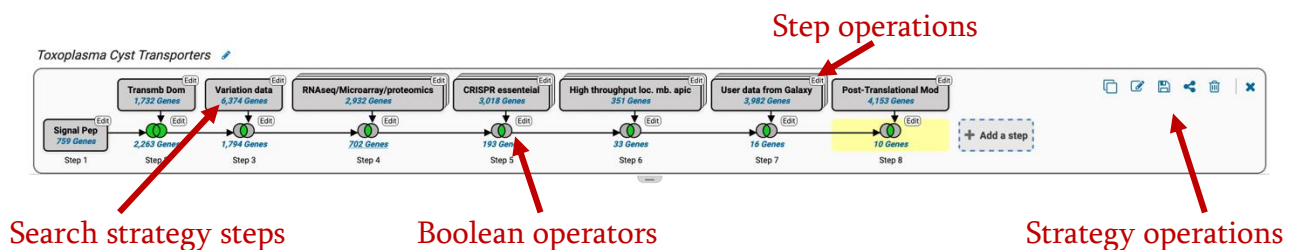# Strategies 2

# Data Integration Through Search Strategies

This exercise illustrates how to combine search results from different data types and how to effectively explore the results. **Specific objectives include**:

1. Understanding search strategy functions including adding/revising/deleting steps, copying search strategies, and saving and sharing strategies.
2. Interacting with gene results and adding columns
3. Navigating transcriptomic searches
4. Exploring proteomics data
5. Exploring subcellular localization data
6. Exploring genome wide CRISPR data
7. Exploring variation data
8. Leveraging orthology searches
9. Running enrichment analyses

## Search strategies

Search strategies in VEuPathDB resources allow you to combine results from different datatype searches using Boolean operators (e.g. Intersect, union, minus). Search strategies enable you to develop *in silico* experiments based on data from the species of interest of from



other species (or strains) by leveraging orthology.

## Getting started with your first search strategy

There are a few things to consider before developing a search strategy:

1. What is your question? Or what are you trying to find out? (overall strategy)
2. Can you break down your question into smaller components? (strategy steps)

3.  What data or analyses can be used to answer the various components of your main question?
4.  How will you combine the different components of your question? Ie. Which Boolean operators.
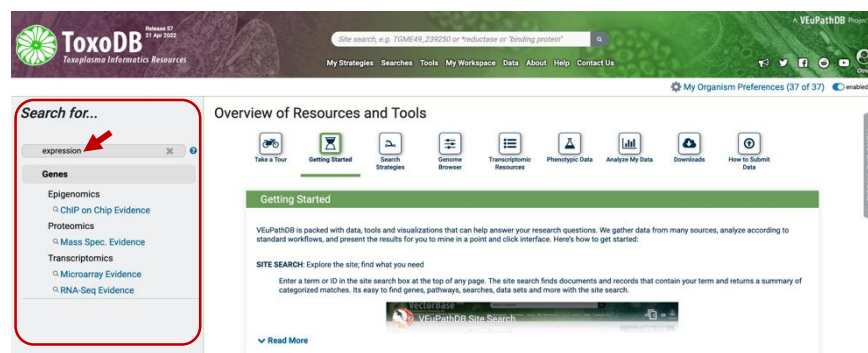
## Example question

**Big question**: I would like to identify bradyzoite/tissue cyst specific therapeutic targets.

**Let's break it down:**
1.  How do I identify genes whose expression is upregulated in bradyzoites?
    a.  Does upregulated mean a gene is not expressed in other stages?
    b.  How do you remove genes that are expressed in other stages?
2.  Should I exclude expression from other stages? How can I do this?
3.  How do I identify genes that have a specific type of variation?
4.  How do I leverage orthology to define phyletic pattern?
5.  What about essentiality? How do I find the genes that are important for parasite fitness?

## Running your first search

1.  Explore the data available in ToxoDB. What data can tell you about expression timing of genes? Expand the menu on the left-hand side of the home page and look for datatypes that would tell you about expression. Hint: try filtering the searches with a key work like "expression".



2.  Explore the RNA-Seq evidence data. Are there any experiments that tell you about bradyzoite expression? Try filtering the datasets using a keywords like "bradyzoite" or "differentiation".

3. Select the fold change search for the "Feline enterocyte, tachyzoite, bradyzoite stage transcriptome (Hehl, Ramakrishnan et al.)".



4. Configure this search to identify genes that are upregulated by 2-fold in tissue cysts compared to all other stages (use average expression values).

5. Add a step and combine these results with results from another experiments containing bradyzoite samples. For example, try selecting the "Stage-specific RNA-sequencing in Toxoplasma gondii (Waldman et al.)". Configure this search to identify genes that are upregulated by 2-fold when comparing 48hr bradyzoites to 24hr tachyzoites.
   o Did you use an intersect or a union operator? How would your results change if you us one or the other? Is one better than the other in this case?



6. How about combining this with data from a microarray experiment? Add a step, go to the microarray data section and select, for example, the "Bradyzoite Differentiation (3-day time series)(Pru) (Buchholz, Fritz and Boothroyd et al.)" experiment. Configure the search to identify genes that are upregulated by 2 fold between time 0 and the other time points.

## Identify Genes based on T. gondii ME49 Bradyzoite Differentiation (3-day time series)(Pru) Microarray (fold change)



7. Add a step and combine any genes that have mass spec evidence from the "Mouse brain bradyzoite proteomics time course (Garfoot et al.)" experiment.

# Identify Genes based on Mass Spec. Evidence

## ❓ Experiments and Samples

*5 selected, out of 101*

select only these | add these | clear these

| brady | ✕ | ❓ |

☐ **Toxoplasma gondii**
  ☐ **Toxoplasma gondii ME49**
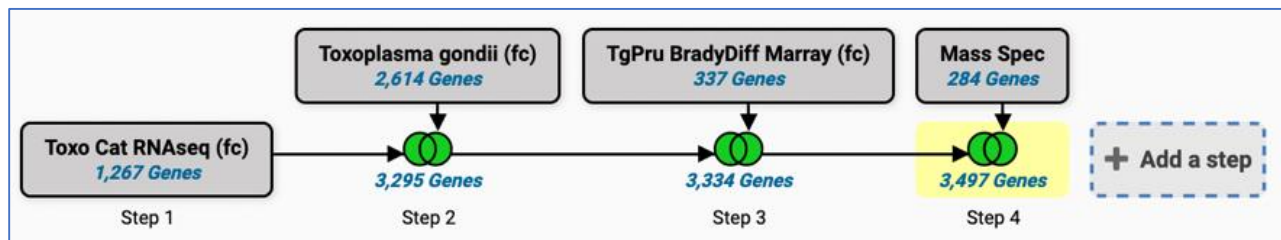    ☑ **Mouse brain bradyzoite proteomics time course (Garfoot et al.)**
      ☑ Bradyzoites 21 days
      ☑ Bradyzoites 28 days
      ☑ Bradyzoites 3 months
      ☑ Bradyzoites 4 months
      ☑ Bradyzoites 5 months

select only these | add these | clear these

## ❓ Minimum Number of Unique Peptide Sequences

| 10 | ⬅ |

| | Toxoplasma gondii (fc) **2,614 Genes** | TgPru BradyDiff Marray (fc) **337 Genes** | Mass Spec **284 Genes** | |
|---|---|---|---|---|
| **Toxo Cat RNAseq (fc)** *1,267 Genes* | 🟢 *3,295 Genes* | 🟢 *3,334 Genes* | 🟢 *3,497 Genes* | ➕ Add a step |
| Step 1 | Step 2 | Step 3 | Step 4 | |

8. Now let's exclude any gene that is highly expressed in tachyzoite and sexual stages. To do this, add a step and select the **percentile** search for the "Feline enterocyte, tachyzoite, bradyzoite stage transcriptome (Hehl, Ramakrishnan et al.)" dataset. Configure the search to exclude any gene that is expressed in all stages except tissue cyst at 70 or higher percentile.

   ❓ **Experiment**

   ◉ Feline enterocyte, tachyzoite, bradyzoite stage transcriptome toxo Transcriptomes of enteroepithelial stages - Sense
   ○ Feline enterocyte, tachyzoite, bradyzoite stage transcriptome toxo Transcriptomes of enteroepithelial stages - Antisense

   ❓ **Samples**

   ☑ EES1
   ☑ EES2
   ☑ EES3
   ☑ EES4
   ☑ EES5
   ☑ Tachyzoites
   ☐ Tissue cysts
   select all | clear all

   ❓ **Minimum expression percentile**

   | 70 |

   ❓ **Maximum expression percentile**
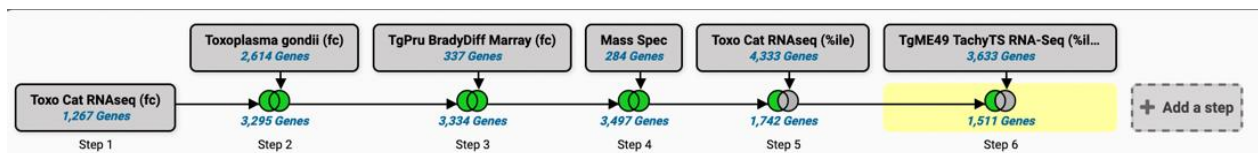
   | 100 |

   ❓ **Matches Any or All Selected Samples?**

   | any ⌄ |

   o What does the "Matches Any or All Selected Samples?" parameter do? Which option is more stringent, any or all?
   o Which Boolean operator did you use?

9. You can exclude additional genes that show expression in stages you are not interested in from other experiments using the above method. Try this with another experiment – for example the "Tachyzoite Transcriptome Time Series (ME49) (Gregory)".



10. Now that we have a list of genes that are upregulated in bradyzoites and are likely not highly expressed in other stages, let's find out which of these have more that 10 non-synonymous SNPs. To do this, add a step and find the search for genes by SNP characteristics.



- o Configure the SNP search to find genes to select all samples aligned to *T. gondii* ME49.

6

**Search for Genes by SNP Characteristics**

The results will be ⬤ intersected with ⌄ the results of Step 6.

**❓ Organism**

| Toxoplasma gondii ME49 | ⟵ |

**❓ Set of Samples**

65 Set of Samples Total

expand all | collapse all

🔍 ❓ Find a variable

📊 Collection year
☰ Country
📊 obsolete_average mapping coverage
⟶ ☰ Data Set
📊 proportion mapped reads
▸ Sample
▸ Sample source
▸ Organism under investigation

65 of 65 Set of Samples selected  Data Set ✕

**Data Set**

⬤ Keep checked values at top

65 (100%) of 65 Set of Samples have data for this variable

| | | | Remaining Set of Samples | | | Set of Samples | | Distribution ❓ | % ❓ |
|---|---|---|---|---|---|---|---|---|---|
| ☑ | ☰ | Data Set | 65 | (100%) | | 65 | (100%) | | |
| ☑ | | Aligned genomic sequence reads - RH Strain | 1 | (2%) | | 1 | (2%) | ▌ | (100%) |
| ☑ | | Aligned genomic sequence reads - White Paper Strains | 62 | (95%) | | 62 | (95%) | ▬▬▬▬▬ | (100%) |
| ☑ | | Toxoplasma gondii ME49 Genome Sequence and Annotation | 1 | (2%) | | 1 | (2%) | ▌ | (100%) |
| ☑ | | Toxoplasma gondii strain CZ clone H3 aligned genome sequence | 1 | (2%) | | 1 | (2%) | ▌ | (100%) |

- o Next set the percent isolates with a SNP call to 60, the SNP type to non-synonymous and the number of SNPs of this type to >/= 10.

**❓ Read frequency threshold**

| 80% ⌄ |

**❓ Minor allele frequency >=**

| 0 |

**❓ Percent isolates with a base call >=**

| 60 |

**❓ SNP Class**

| Non-Synonymous ⌄ |

**❓ Number of SNPs of above class >=**

| 10 |

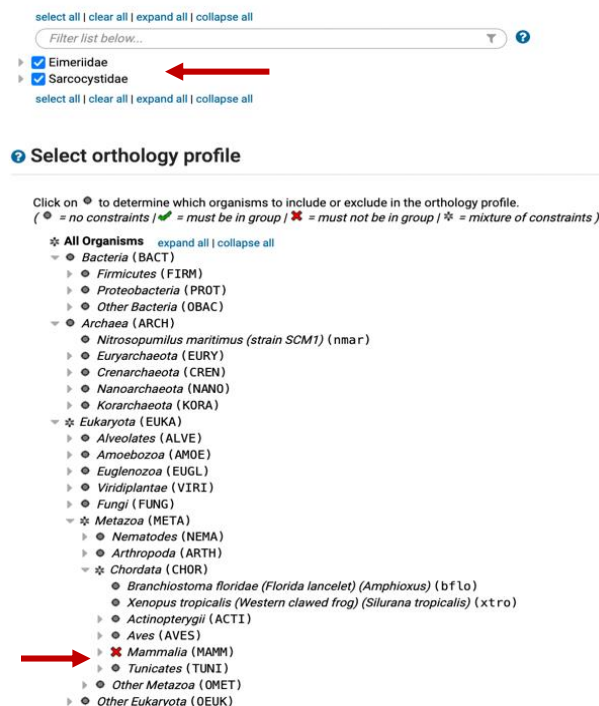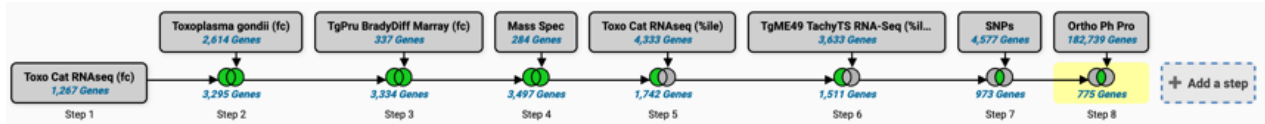| **Toxo Cat RNAseq (%ile)** 4,333 Genes | **TgME49 TachyTS RNA-Seq (%ile)** 3,633 Genes | **SNPs** 4,577 Genes | ┊ ➕ Add a step ┊ |
|---|---|---|---|
| ⬤ 1,742 Genes | ⬤ 1,511 Genes | ⬤ 973 Genes | |
| Step 5 | Step 6 | Step 7 | |

11. Now let's determine how many of these genes do not have orthologs in mammals. Add a step and find the search called "Orthology Phylogenetic Profile".



- o There are different ways to configure this search depending on which Boolean operator you use. If you use the intersect operator, then configure the search to return all genes in ToxoDB that do not have orthologs in mammals.

12. As a final step let's determine which of these genes are essential based on the genome wide CRISPR screen from the Lourido lab. Add a step and find the CRISPR phenotype search. Set the phenotype score < = to -2.

**Search for Genes by CRISPR Phenotype**

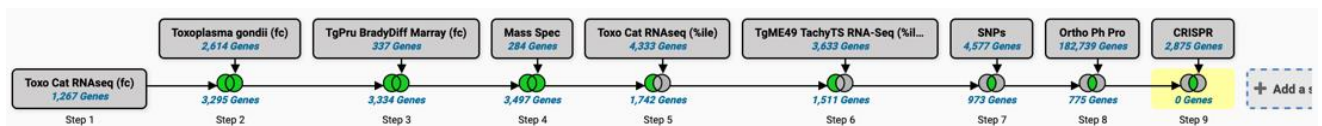The results will be [○ intersected with ∨] the results of Step 8.

**❷ Phenotype Score >=**

-6.89

**❷ Phenotype Score <=**

-2  ⟵

Run Step



- o How many results did you get? Is this surprising? Why do you think you got 0 results?
- o How can you get over the problem observed above? Is there a tool that would allow you to convert *T. gondii* GT1 genes to *T. gondii* ME49 genes?

13. Hover over the CRISPR step and click on the edit icon. In the popup click on the "orthologs" option and select ME49 from the list of organisms to transform to.
- o Did this improve the results?

14. Explore the genes in your result list. Are there any interesting genes that you might purpsue further in the lab?

15. How many hypothetical genes are in your results? A quick way to find out is to click on the graph icon in the Product Description column heading. This generates and



interactive word cloud. Hover over the word hypothetical to see the number.

16. What can you do to figure out what some of these hypothetical genes do? Here are a couple of suggestions:

   o Add a column for InterPro domain descriptions. Click on add column, search for InterPro and add the appropriate column. Did this give you an idea for



possible functions for some of the hypothetical genes?

o   Add another column for the hyper_LOPIT subcellular localization data.  Add the column called "Predicted Location (TAGM-MAP)". Did this reveal some possible clues about the some of the other hypotheticals?

17. Do your results contain an enrichment of certain functions? Try some of the analysis tools available in the analysis tab.



Link to strategy:

https://toxodb.org/toxo/app/workspace/strategies/import/ce05740b75a965d9