# Orthology and Gene Ontology

## Learning objectives:

- Explore the orthology table on VEuPathDB gene pages

- Run and explore results of searches in OrthoMCL

- Leverage the phyletic pattern search

- Leverage the orthology transform tool

- Run and explore the results of a GO enrichment analysis

- Port GO enrichment results to Revigo

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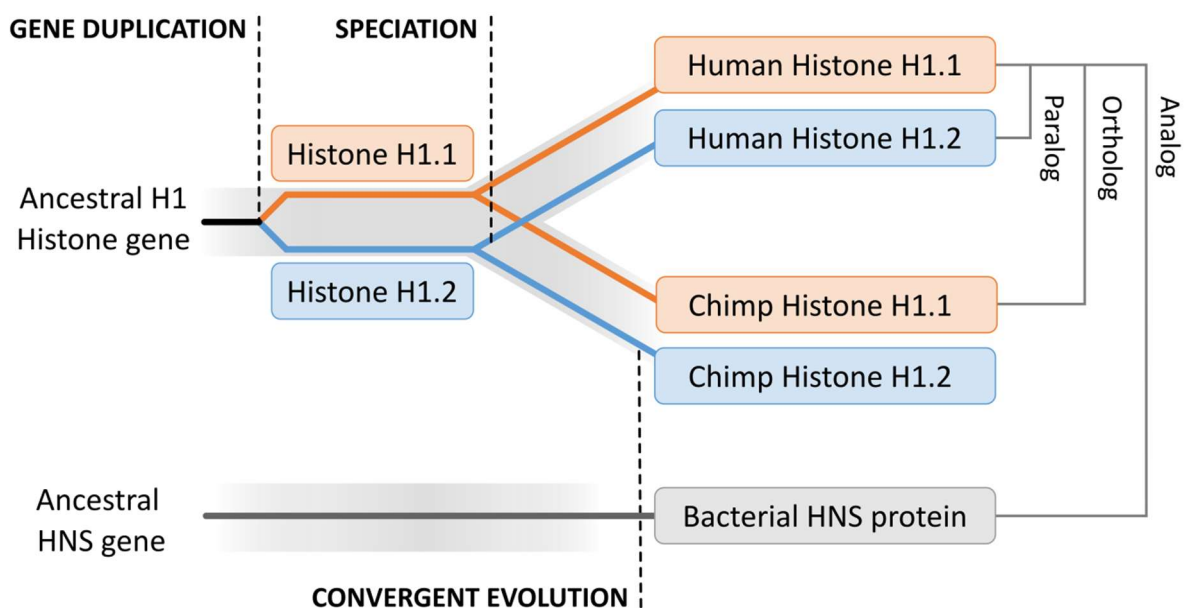Informatics Resources

# Table of Contents

# Introduction

## 1. Orthology and phyletics

Homologs are genes that share ancestry either by speciation (orthologs), gene duplication (paralogs), or gene transfer events (xenologs). Paralogs of a conserved gene may occur in a single species or strain. Conserved sequences in genomes can be used to infer evolutionary history (e.g., ribosomal sequences), their similarities and differences can be used to trace the divergence and evolution of organisms. Genes that share function by convergent evolution, but do not share ancestry are known as analogs.

Ortholog groups can also allow you to explore the potential functions of a gene, or group of genes across species. In pathogens like *Plasmodium falciparum*, ortholog groups might facilitate the identification of potential targets for drug or vaccine development.

For more detail on orthologs, paralogs and evolutionary genomics, read the review by Koonin[1].



*Figure 1. Gene phylogeny (orange and blue) within species phylogeny (grey). Top shows an ancestral gene duplication event, producing two paralogs of the Histone H1 gene, producing H1.1 and H1.2. This is followed by a speciation event leading to Chimpanzee and Human Orthologs of the two genes. Bottom shows a gene with separate evolutionary origin that has evolved similar function to H1 Histones through convergent evolution, HNS (histone-like nucleoid-structuring protein). HNS is a bacterial analog to H1 Histone. Figure adapted from* this image *by Thomas Shafee (2018).*
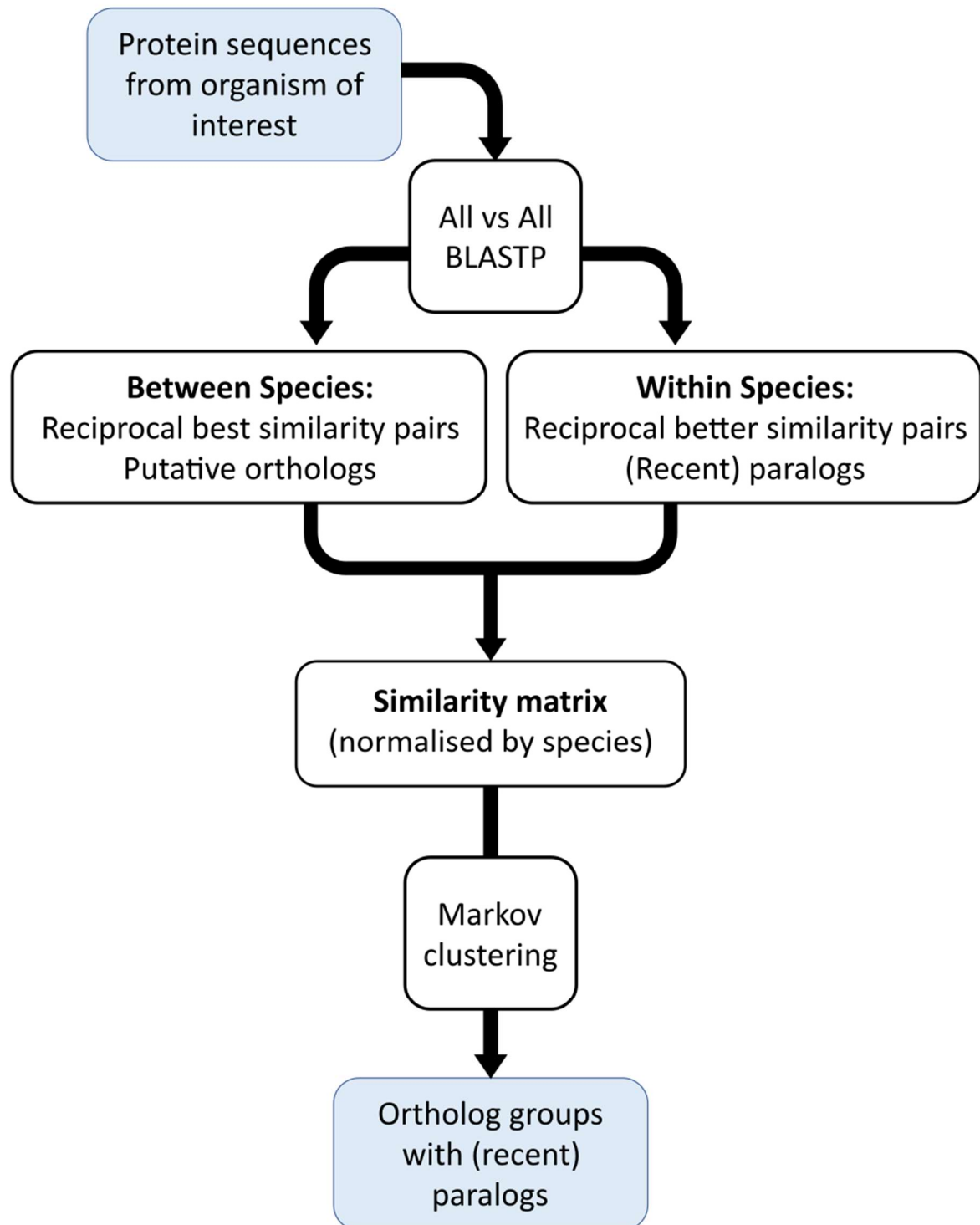
## 2. OrthoMCL

OrthoMCL is a genome-scale algorithm for grouping orthologous protein sequences which not only share evolutionary history, but also share function. Thus, ortholog prediction is important in predicting the function of newly identified proteins. Detection of orthologs has become more widespread with the rapid progress in genome sequencing and the discovery of protein sequences [2,3].

OrthoMCL provides a database of ortholog groups with high degrees of functional conservation (e.g., they have consistent EC numbers, which link genes to specific products in metabolic pathways), making it a useful tool for functional annotation of genomes [4,5].

OrthoMCL identifies shared protein groups between species and is also capable of representing species specific gene expansion families. To achieve this, the OrthoMCL algorithm starts with reciprocal best BLAST hits within each proteome as potential in-paralog/recent paralog pairs and reciprocal best hits across any two proteomes as potential ortholog pairs. Related proteins are interlinked in a similarity graph. Then, [MCL](link) [6] is invoked to split mega-clusters - clusters that are just too big and uninformative. MCL clustering is based on weights between each pair of proteins, which are normalised by species to account for evolutionary distance. If you want to know more about how the MCL algorithm works have a look at [this simplified explanation](link) [7].

The organism specific orthology information garnered from our OrthoMCL analysis of VEuPathDB organisms is presented on gene pages and integrated into an Orthology Phylogenetic Profile search. They are available for anyone to explore or use for their own investigations. The OrthoMCL.org site offers a deep look into all data associated with the OrthoMCL results for orthology groups and proteins.
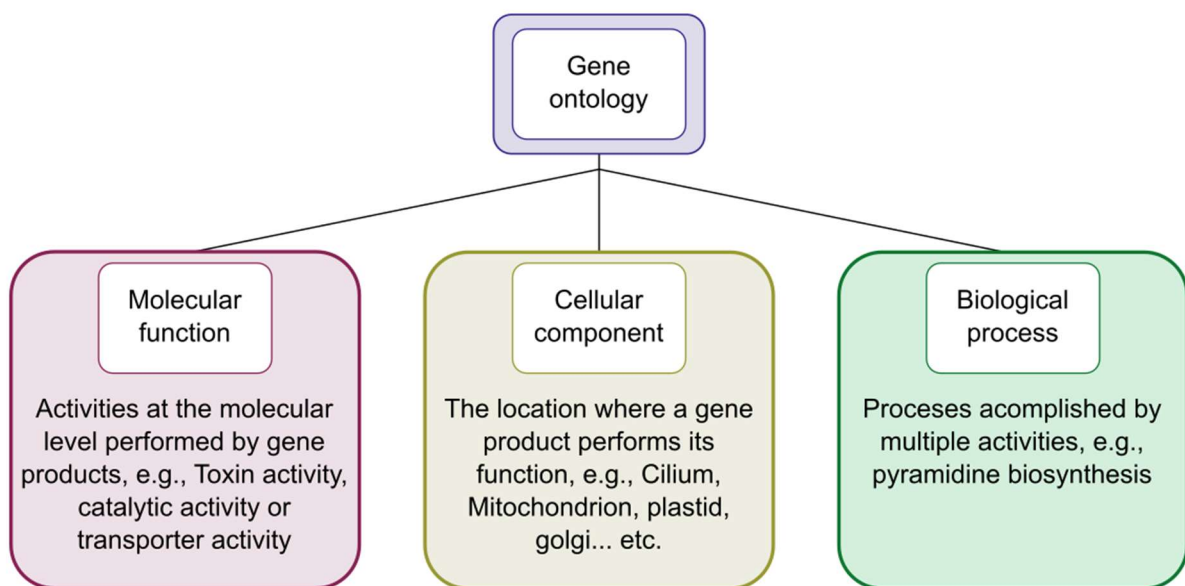
**Figure 2.** *OrthoMCL's workflow.*

## 3. Gene Ontology

Ontologies are a controlled vocabulary of terms and concepts with relationships between them. 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/



A gene can be assigned a GO term either manually (by an annotator or curator when they evaluate experimental evidence from a publication) or computationally (based on the GO terms of genes that share sequence or functional domains). The origin of the assignment is documented; some researchers believe that manually assigned functional annotations are more accurate than those that are electronically transferred since a researcher has reviewed the manually annotated assignments. GO terms can be used to test whether your set of genes are enriched for a molecular function, cellular component, or biological process.

**For example:** A researcher performs a proteomics experiment on a protein fraction collected during an antimalarial treatment and identifies 100 proteins in total. When they examine the GO terms assigned to the gene set corresponding to the proteome, they see that 25 genes are assigned GO:0016301, kinase activity. Out of 5000 genes in the genome, only 100 are assigned GO:0016301. There is an overrepresentation of GO:0016301 in the researcher's proteome which is 'enriched' for kinase activity.

A standard enrichment determination method employs Fisher's exact test, a statistical test that evaluates a 2x2 contingency table (in this case, the number of genes in my set versus number of genes from the same genome not in my set, and number of genes with GO term X versus number of genes without term X). This test produces a p-value between 0 and 1, where $p \leq 0.05$ is considered significant (that is, less than 5% probability that the enrichment is due to chance).

However, the test is performed for each of the 100s of GO terms, increasing the chances that a GO term will be incorrectly considered enriched (a false positive, or type I, error). Thus, the original p-value must be adjusted for so-called multiple hypothesis testing, resulting in an adjusted p-value such as the Benjamini-Hochberg false discovery rate (FDR) or Bonferroni adjusted p-value.

# Exercises

## 1. Using the orthology transform tool to identify apicoplast targeted genes in *Toxoplasma* and *Neospora*.

**Note:** For this exercise use http://veupathdb.org

**What is an apicoplast?**

The apicoplast likely became encased in four membranes via a double endosymbiotic event. The chloroplast arose by engulfment of a cyanobacteria by a plant/algae ancestor. An alga was then engulfed by the ancestor of all apicomplexans. Thus, an apicoplast organelle arose with four membranes.

1. Start by finding genes in *Plasmodium* that are predicted to target the apicoplast.

   **Hint:** Navigate to the Pfal 3D7 Subcellular Localization search for Apicoplast. You can filter the type of search by text query.



2. You can further expand your list of potentially Apicoplast targeted proteins by adding a GO terms search strategy for the term "apicoplast" or the GO ID: "GO:0020011" in *P. falciparum* 3D7.

   **Hint:** click on add step the go to the function prediction category and select the GO term search.

   a. Which Boolean operation did you use? Union or intersect?

Add a step to your search strategy

**Combine** with other Genes

Subcell Loc
499 Genes
Step 1    Step 2

**Transform** into related records

Subcell Loc
499 Genes
Step 1    Step 2

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

1. Choose *how* to combine with other Genes

○ 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
Phenotype
🔍 CRISPR Phenotype
Text
🔍 Text (product name, notes, etc.)

**Search for Genes by GO Term**

The results will be ◉ unioned with ▾ the results of Step 1.

Configure Search    Learn More    View Data Sets Used

❷ **Organism**

*1 selected, out of 622*
select only these | add these | clear these

3d ×

⊟ Apicomplexa
  ⊟ Aconoidasida
    ⊟ Haemosporida
      ⊟ Plasmodiidae
        ⊟ Plasmodium
          ⊟ Plasmodium falciparum
            ☑ Plasmodium falciparum 3D7 **[Reference]**

❷ **Evidence**

☑ Curated
☑ Computed
select all | clear all

❷ **Limit to GO Slim terms**

○ Yes
● No

● ❷ GO Term or GO ID

GO:0020011 : apicoplast : 7 ✕

**GO Term**
*370 Genes*

**Subcell Loc**
*499 Genes*

*635 Genes*

Step 1    Step 2

8

3. Add a step to your strategy that transforms the results with *Toxoplasma* and *Neospora* orthologs.

4. Although *Cryptosporidium* is an apicomplexan parasite it has lost its apicoplast! Can you use this fact to refine your results from the above search?

    a. First pull up all *Cryptosporidium* genes with the Genes by Taxonomy search and then transform these back to their *Toxoplasma* and *Neospora* orthologs for the subtraction to complete. Think about what kind of intersection you should be using!

**Hint:** try subtracting out any orthologs present in *Cryptosporidium*. You will need to use a nested strategy.

## How to make a nested strategy

Click edit > nested strategy



Then to view and edit the nested strategies, for instance, to add an ortholog transformation step for *Neospora* and *Toxoplasma*, click edit > view



This leaves you with apicoplast specific genes for *Toxoplasma* and *Neospora* that you could target in future research.

## 2. GO enrichment analysis

GO term enrichment analysis can be carried out at any stage in any search strategy to see what categories of genes are most common. This function is currently available on every organism site except for the VEuPathDB main site (this is due to the underlying complexities of the website).

Pick your favorite database (e.g. plasmodb, toxodb) and have a play around with GO terms or follow along with this example.

### Retrieving RNA-seq evidence

To run a GO enrichment analysis, you need a list of genes to test. This can be a list of gene IDs from your experimental results (upload them with the ID search) or a gene list resulting from a search you conducted on a VEuPathDB website. For this example, in ToxoDB, we will identify genes that are differentially regulated over time.

1. Navigate to the RNA-Seq searches and find the data set called "**Oocyst Time Series (M4)**" from Fritz et al. A quick way of getting to the RNA-Seq searches is to type 'rna' in the filter box on the left of the home page and click on the RNA-Seq Evidence link. See image below.



2. The RNA-Seq evidence page includes a list of all data sets that are loaded in the website. To quickly find a dataset, you can start typing key words in the "Filter Data Sets" box. For example, start typing the word "oocyst".

Identify Genes based on RNA-Seq Evidence

3. Once you find the data set of interest, choose the fold-change (FC) search. For this exercise, identify genes that are upregulated by 20-fold in days 4 and 10 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)

4. Click "Get Answer" to initiate the search. This will return a one-step search strategy.

   a. How many genes did you get?



TgM4 Oocyst RNA-Seq (fc)
873 Genes
Step 1

## Run the GO enrichment analysis

1. Click on the Analyze Results tab just above the list of genes (arrow in image below) to open the enrichment tools.

    a. Besides GO enrichment, what other analyses are available?



2. 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".

    • Organism = T. gondii ME49

    • Ontology = Cellular Component

    • Evidence = Computed and Curated

    • Limit to GO Slim terms? = NO

## Explore your results of your analysis

1.  What is the top enriched GO term from this analysis?

2.  Does this make sense for an enrichment analysis of the cellular component of your Oocyst expressed genes? Notice that the p-value with Benjamini or Bonferroni correction is very low.

3.  What do each of the columns in the analysis table represent?

    **Hint**: *move your mouse over the question mark next to each column header*



- *Fold enrichment* -The ratio of the proportion of genes in the list of interest with a specific GO term over the proportion of genes in the background with that term.

- *Odds ratio* -The odds of the GO term appearing in the gene list are the same as that for the background list.

- *P-value* –The null hypothesis or the probability of getting a result that is equal or greater than what was observed.

- *Benjamini-Hochburg false discovery rate* – A method for controlling false discovery rates for type 1 errors.

- *Bonferroni adjusted P-values* - A method for correcting significance  based  on multiple comparisons.

<span style="color:orange">Port your results to Revigo</span>

1. Click the Open in Revigo button to port the results to Revigo, the Reduce and Visualize Ontology tool.



2. Once at Revigo, you may need to scroll down to click Start Revigo to run the analysis with default parameters.  Revigo provides a scatterplot and table, an integrative map and a tree map to supplement the table provided in the VEuPathDB site. See the Revigo publication for more information.

3. Try rerunning the GO enrichment analysis, but this time select the Molecular Function ontology.
   a. What is the top enriched GO term?
   b. What is the p-value for the enrichment?
   c. Do you have more or less confidence than in the last search that this function is enriched in your gene set?
4. Click on the "Word Cloud" button above the Molecular Function analysis results. What type of analysis is this? What information can you (See image below).
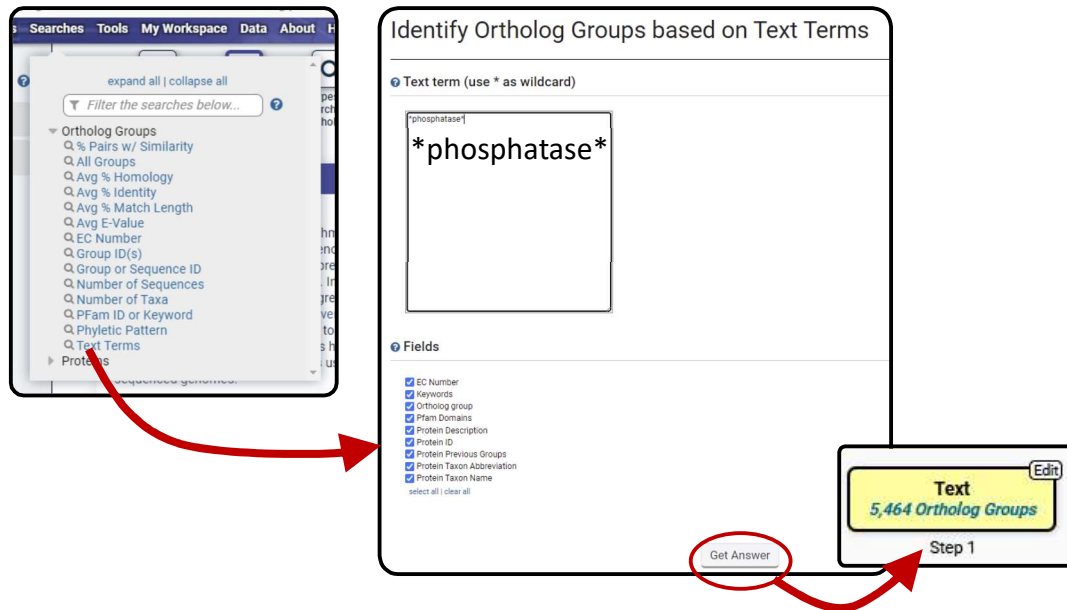
# 3. Search orthologs and filter by phyletics in OrthoMCL (optional)

**Note:** Use http://orthomcl.org for this exercise.

Find all plant proteins that are likely phosphatases that do not have orthologs outside of plants.

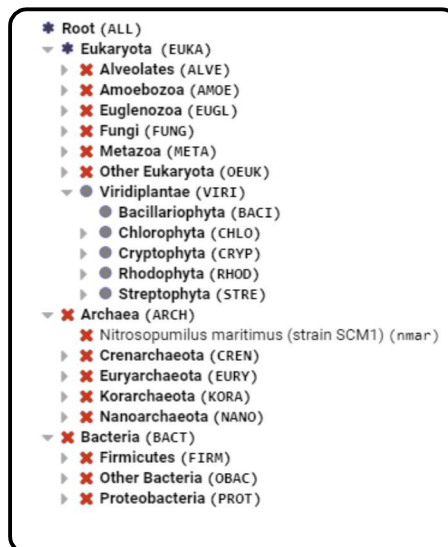1.  Use the text search to find OrthoMCL groups that contain the word "*phosphatase*". Note that the search should be run without the quotation marks but with the asterisks.

    

2.  Add a step and run a phyletic pattern search for groups that contain any plant protein but do not contain any other organism outside plants.
    **Hint:** make sure everything has a red x on it except for plants (Viridiplantae (VIRI)), which should be a grey circle.

3. Examine your results.

    a. How many groups were returned by the search?

    b. What is the distribution of plant proteins in each orthology group? (use the Add Columns tool to turn on the Viridiplantae column if it is not already on)



4. Next, you can run a multiple sequence alignment for OG6_112109.  Click on the group ID in your result table and navigate to the List of Proteins section of the group page.  The Clustal Omega tool is integrated into the table.  There are several formats available for the Clustal output, making it easy to take these results to other visualization programs.

## 4. Explore a specific OrthoMCL group - examining the cluster graph (optional)

**Note:** Use http://orthomcl.org for this exercise.

1. Visit the OrthoMCL group OG6_131670.  Use the site search to navigate to OG6_131670.

2. Examine the Phyletic Distribution.

    a. What is the phylogenetic distribution of the members of this group? The distribution is presented as a tree.  Expand the tree to view the distribution.



3. Navigate to the Cluster graph tab.  Modify the E-value cutoff slider.  What happens when you increase or decrease the E-value? Can you identify subclusters of orthologs?  The view of the graph can be changed using the Edge type options and the Node options.

# Resources

These are the operators you need to be aware of for these exercises. Don't forget to refer to the search strategies help sheets for more in-depth help!

| Operator | Combined search will contain: |
|---|---|
| 1 INTERSECT 2 | IDs common between both lists |
| 1 UNION 2 | All IDs from both lists |
| 1 MINUS 2 | IDs only in list 1 |
| 2 MINUS 1 | IDs only in list 2 |

**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

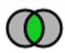**Fisher's Exact Test and the Hypergeometric Distribution (the M&M example)** - https://youtu.be/udyAvvaMjfM

**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

# References

1. Koonin, E. V. Orthologs, Paralogs, and Evolutionary Genomics. *Annu. Rev. Genet.* **39**, 309–338 (2005).

2. Glover, N. *et al.* Advances and Applications in the Quest for Orthologs. *Mol. Biol. Evol.* **36**, 2157–2164 (2019).

3. Linard, B. *et al.* Ten Years of Collaborative Progress in the Quest for Orthologs. *Mol. Biol. Evol.* **38**, 3033–3045 (2021).

4. Li, L., Stoeckert, C. J. & Roos, D. S. OrthoMCL: Identification of Ortholog Groups for Eukaryotic Genomes. *Genome Res.* **13**, 2178–2189 (2003).

5. Fischer, S. *et al.* Using OrthoMCL to Assign Proteins to OrthoMCL-DB Groups or to Cluster Proteomes Into New Ortholog Groups. *Curr. Protoc. Bioinforma.* **35**, (2011).

6. Van Dongen, S. Graph Clustering Via a Discrete Uncoupling Process. *SIAM J. Matrix Anal. Appl.* **30**, 121–141 (2008).

7. Jagota, A. Markov Clustering Algorithm. *Towards Data Science* https://towardsdatascience.com/markov-clustering-algorithm-577168dad475 (2020).

# Glossary

**Clusters of Orthologous Genes (COGs)**, phylogenetic classification of proteins encoded in complete genomes.

**Co-orthology**, recent descent and duplication

**EC numbers**, Enzyme Commission number used to classify enzymes based on the chemical reactions that they classify.

**Homologs**, genes that share ancestry either by speciation (orthologs), gene duplication (paralogs), or gene transfer events (xenologs)

**In-paralog**, recent duplication.

Markov Clustering Algorthm (MCL),

**Orthology** is the study of genes across species that are conserved over evolutionary time.

**Ortholog**, homologous genes resulting from gene duplication.

**Paralog**, homologous genes resulting from speciation.

**Phyletics**, or **phylogenetics**, is the study of patterns in evolutionary history across species.

**Reciprocal best hits**, where sequences from two different genomes find each other as the best scoring match

**Xenolog**, homologous genes resulting from a horizontal or lateral gene transfer event