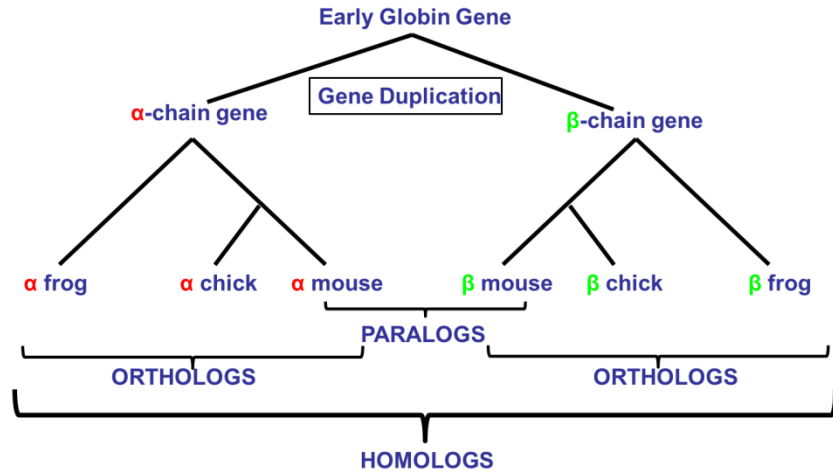


## Orthology and Phyletic Patterns

### Homology



#### 1. Getting to OrthoMCL from EuPathDB databases

Note: For this exercise use <http://cryptodb.org> and <http://orthomcl.org/>

- Go to the gene page for the *Cryptosporidium parvum* gene with the ID: cgd7\_2290
- What information on the gene page can you use to guess a function for this gene? It is annotated as a hypothetical protein! Hint: look at the orthologs table and the domains in the protein features graph. You may also want to visit some of the external links.
- Scroll down to the table labeled “Orthologs and Paralogs within CryptoDB”. Does this gene have orthologs in other *Cryptosporidium* species? What about other organisms? (hint: click on the link below the table that takes you to OrthoMCL).

Gene	Organism	Product	is syntenic	has comments
Cvel_467	Chromera velia CCMP2878	rRNA-processing protein FCF1 homolog, putative	no	no
Chro.70261	Cryptosporidium hominis TU502	hypothetical protein	yes	no
CMU_034340	Cryptosporidium muris RN66	hypothetical protein, conserved	yes	no
GNI_088410	Gregarina niphandrodes Unknown strain	rRNA-processing Fcf1-like protein	no	no
Vbra_6876	Vitrella brassicaformis CCMP3155	rRNA-processing protein FCF1 homolog, putative	no	no

View the group (OG5\_127679) containing this gene (cgd7\_2290) in the OrthoMCL database

- Does this protein have orthologs in other organisms? Does it have any orthologs in bacteria or archaea? (Hint: mouse over the colorful boxes in the table to reveal the full species and phylum names – see image below).



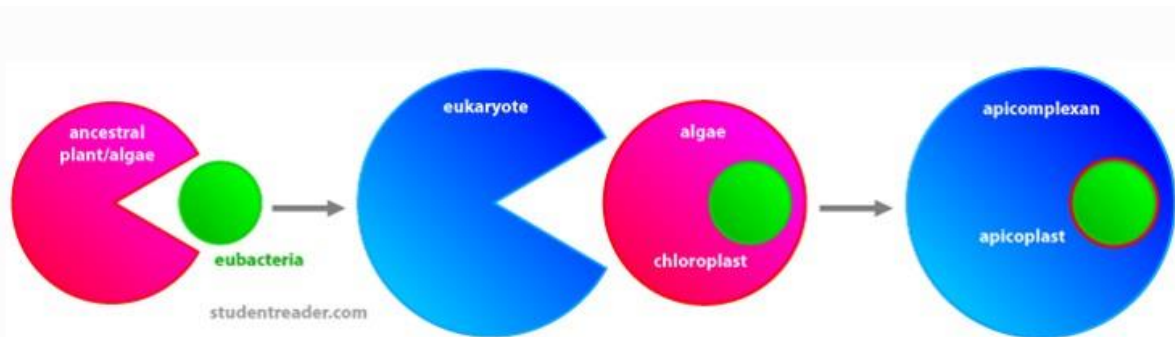
category). To specify a phyletic pattern click on the icon next to the taxonomic group or species to include or exclude it.

- a. How many protein groups do not contain orthologs from eukaryotes?
- b. Find all groups that contain orthologs from at least one species of *Cryptosporidium* and *Giardia* but not from bacteria or archaea.

All EuPathDB sites also have a phyletic pattern search that uses OrthoMCL data under Genes -> Evolution -> Orthology Phylogenetic Profile. This search is very useful to identify genes in your organism of interest that are restricted in their profile. For example, you frequently want to identify genes that are conserved among organisms in your genus but not present in the host as these genes may make good drug targets or vaccine candidates. Optional: go to your favorite EuPathDB site and run this search to identify all genes that are not present in human or mouse.

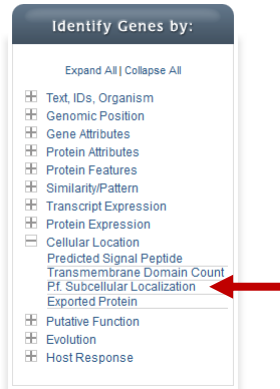
### 3. Using the orthology transform tool to identify apicoplast targeted genes in *Toxoplasma* and *Neospora*.

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



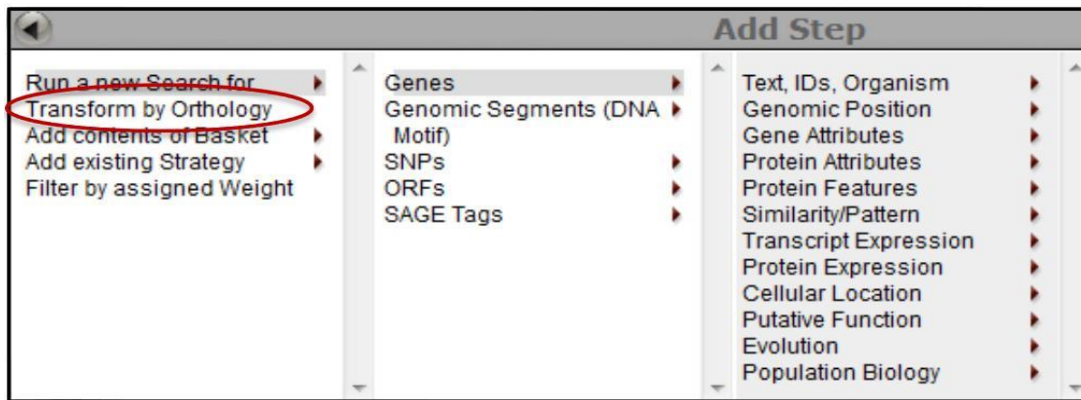
The apicoplast likely became encased in four membranes via a double endosymbiotic event. The chloroplast arose by engulfment of a cyanobacteria by a plant/algae ancestor. An algae was then engulfed by the ancestor of all apicomplexans. Thus an apicoplast organelle arose with four membranes.

- a. Start by finding genes in *Plasmodium* that are predicted to target to the apicoplast. Hint: click on “Cellular Location” then on “P.f. Subcellular Localization”.



b. Transform the results of the above search to their *Toxoplasma* orthologs.

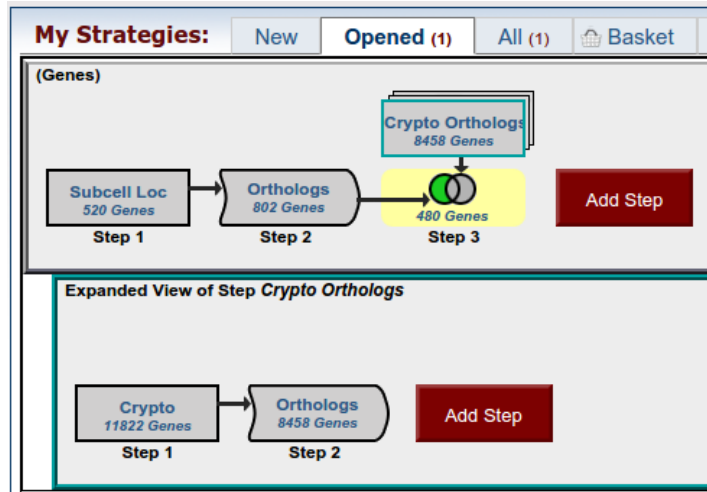
Hint: add a step, then select “Transform by Orthology”. On the search page, select all



*Toxoplasma* and *Neospora*.

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

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



4. Combining searches in OrthoMCL (Use <http://orthomcl.org> for this exercise).

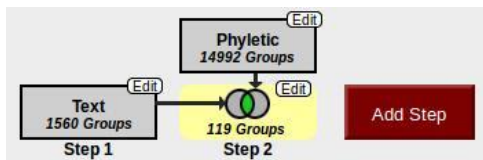
Find all plant proteins that are likely phosphatases that do not have orthologs outside of plants.

- a. 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).

The screenshot shows the OrthoMCL DB website interface. The header includes the logo 'OrthoMCL DB' (Version 5, 10 May 13) and 'A EuPathDB Project'. There are two search bars: 'Groups Quick Search:' containing '\*phosphatase\*' (circled in red) and 'Sequences Quick Search:' containing 'synth\*'. Below the search bars is a navigation menu with items like Home, New Search, My Strategies, My Basket (0), Tools, Data Summary, Downloads, and Community. The main content area displays search results for '\*phosphatase\*', organized by taxonomic groups. The groups are listed on the left, and corresponding taxon IDs are listed on the right. Plant groups (Viridiplantae) are highlighted in grey, while other groups are in red. The groups shown include Aconoidasida (ACON), Haemosporida (HAEM), Piroplasmida (PIRO), Amoebozoa (AMOE), Euglenozoa (EUGL), Viridiplantae (VIRI), Streptophyta (STRE), Chlorophyta (CHLO), Rhodophyta (RHOD), Cryptophyta (CRYP), Bacillariophyta (BACI), Fungi (FUNG), Microsporidia (MICR), and Basidiomycota (BASI). The taxon IDs listed include pber, pcha, pfal, pkno, pviv, pyoe, bbov, tann, tpar, ddis, ehis, edis, einv, lbra, linf, lmaj, lmex, tbru, tbrg, tcon, tcru, tviv, atha, osat, ppat, rcom, micr, crei, otau, vcar, cmer, gthe, tpse, ecun, ebie, eint, cneo, cneg, lbic, pchr.

- b. 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).

- c. How many groups did you return? Explore the multiple sequence alignments from some of these groups. (Hint: click on a group ID and open the MSA tab).
5. Exploring a specific OrthoMCL group - examining the cluster graph. (Use <http://orthomcl.org> for this exercise).
- Visit the orthomcl group OG5\_127676. You can either type the ID in the group quick search option at the top of the page or follow this link: [http://orthomcl.org/group/OG5\\_127676](http://orthomcl.org/group/OG5_127676)
  - Examine the "Sequences & Statistics" tab: Based on the EC description and the product descriptions of the members of this group, what kind of a protein does this group represent? What is the phylogenetic distribution of the members of this group?
  - Examine the "PFam Domains (graphic)" tab: How many PFam domains are represented in this group? What is the most common one? Which one is the least common one?



- Examine the "Cluster Graph" tab: Modify the E-value cutoff slider. What happens when you increase the E-value? What happens when you decrease the E-value? Can you identify subclusters?

Group: OG5\_150204 (10 sequences)

Sequences & Statistics | PFam domains (graphic) | PFam domains (details) | MSA | Cluster graph

Phyletic Distribution Hide

Legend: 0 no o

MUSCLE (3.7) multiple sequence alignment

```

otau est:Ext_fgenehsl_pg_Chr_06  -----MSRF-----HSDGFLA---RTVDDDDL-----
osat NP_001052931  MAAAAAATVEAVGVAGRRRR-----SGSVALGDLRREASERASASAGAGGRERE
osat NP_001047178  -----NSTRSKSVFPAGGGAATVPLAVLLRREVSEKTAAE-----
osat NP_001050291  -----REVQGEV-----SSVFLAVLLKRELCKQVE-----
rcm 30I70.m013899  -----MNAARE-----HQTVELSVLLKRELNERVE-----
atha NP_564504  -----MSTKGE-----HHTVFLSVLLKRESANEKID-----
ppat e_gwl.29.62.1  -----MVPLASLAVELKNEIVE-----
atha NP_001031252  -----NASREKRRNHNHDEKLVFLAALISRETKAAKME-----
atha NP_177008  -----NASREKRRNHNHDEKLVFLAALISRETKAAKME-----
rcm 30066.m000733  -----NASGEGRCR-----DNLVFLAALISREKNEKME-----
                                     . * .
otau est:Ext_fgenehsl_pg_Chr_06  -----ASIVASRKRKQGEDFLSVPALTWTFRATGEDAFV-----MFLGYLFDHG
osat NP_001052931  RRP5VAAGQACRAKRGEDFALLKPCACERLPAGGA-----PFSFALFDHG

```

Edge Options: E-value Cutoff Max E-Value: 1E-140

Node Options: Show Nodes By: EC Numbers, PFam Domains