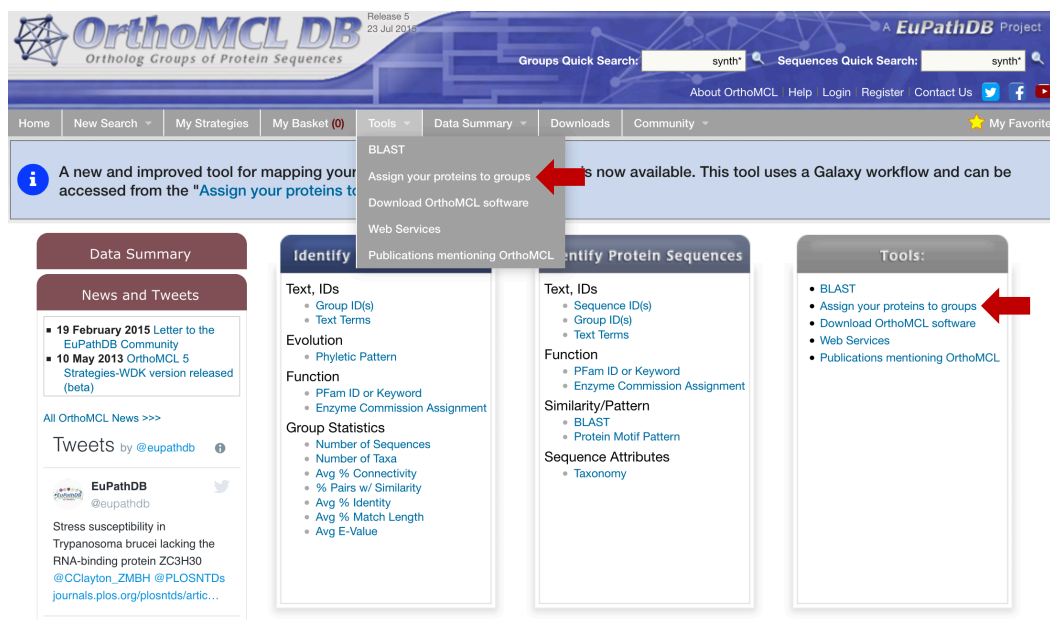# Map your proteins to OrthoMCL groups

This tool allows you to map proteins from a FASTA file to OrthoMCL groups.  The tool has been implemented as a workflow in the VEuPathDB Galaxy workspace.  To use this tool, you can follow these steps:

1.  Click on the "Tools" item in the grey menu then select "Assign your proteins to groups" or click on the "Assign your proteins to groups" link from the right-hand side of the main OrthoMCL page.



2.  The next page provides some instruction and a link to the VEuPathDB Galaxy server.  Click on the VEuPathDB Galaxy Server link to access this service. **Note:** to use this service you will have to create a VEuPathDB account.

3. Once in Galaxy, you can import your protein FASTA file by clicking on the "Get Data" option in the left pannel and selecting an import option, such as "Upload file from my computer".



4. Your imported file will appear as a step in the right hand history pannel. The color of the step indicates its status. The step is grey when it is in queue, yellow when in process and green when completed.
5. Once your protein FASTA file has been successfully uploaded into Galaxy, select the "Workflow to map your proteins to OrthoMCL groups" from the middle section.

6. This will import the workflow into your workspace and allow you to run it on your protein file.  In most cases you will not need to change any of the default parameters.  Simply click on the "Run workflow" button at the bottom of the middle section. Note that there is a few second lag between clicking on the button and the workflow starting to run – please be patient.

## Running workflow "imported: Map your proteins to OrthoMCL groups"

[Expand All] [Collapse]

Assign your set of proteins to OrthoMCL groups. This workflow uses BLASTP and the OrthoMCL algorithm to (1) map proteins to OrthoMCL groups using BLAST similarity and (2) form paralog groups from proteins with no significant similarity to any OrthoMCL proteins. The workflow produces a file with the mapping from protein ID to group ID, along with similarity metrics. It also produces a file of paralog groups. The latter file is only valid if your input proteins all belong to a single proteome. This workflow might take 24 hours or more to run, depending on the size of the job queue.

Step 1: OrthoMCL Clean FASTA file (version 1.0.0)
1

**Proteome FASTA file:**
1: MyProteins

**Maximum allowed number of input sequences**
100000 ✎

**Action:**
Hide output 'output'.

Step 2: NCBI BLAST+ makeblastdb (version 0.3.0)
7

Step 3: NCBI BLAST+ blastp (version 0.3.0)
3

Step 4: NCBI BLAST+ blastp (version 0.3.0)
2

Step 5: OrthoMCL Reformat Blast (version 1.0.0)
5

Step 6: OrthoMCL Reformat Blast (version 1.0.0)
4

Step 7: OrthoMCL Map Proteome to Groups (version 1.0.0)
6

Step 8: MCL (version 14.137)
7

☐ Send results to a new history

[Run workflow]

7.  Workflow steps will queue up in the right-hand panel.  The entire workflow may take up to 24 hours or more to run depending on the size.



8.  Once the workflow has completed, you can work with the output files.  Click on the name of the step in your history in the right-hand panel to expand it.  This provides additional options including the option to download file.