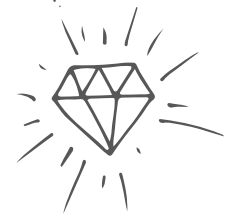# Map Proteins to OrthoMCL with Diamond BLAST- A Tutorial
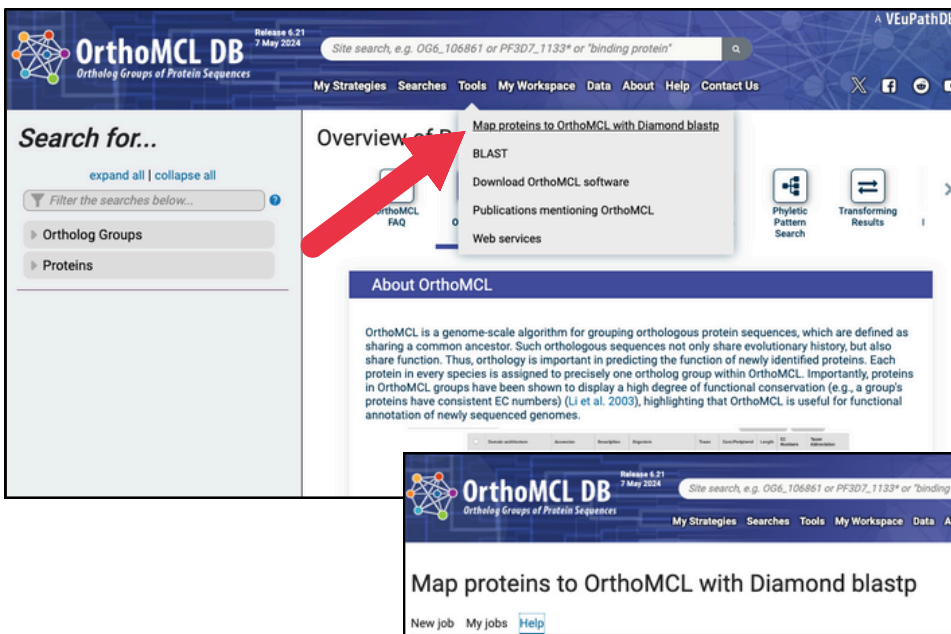
## Learning Objectives

- Understand the purpose of the protein mapping tool in OrthoMCL
- Learn how to prepare and upload data
- Explore the output and understand the Diamond job result page

1. Introduction
2. OrthoMCL is a genome-scale algorithm for grouping orthologous protein sequences. The OrthoMCL algorithm is employed on proteins from a set of 150 Core species to form Core and Residual ortholog groups.
3. Purpose of the protein mapping tool: The purpose of this tool is to allow users to map a set of proteins of interest, usually a complete proteome from an organism, to existing OrthoMCL groups. This tool uses Diamond, a newer computing alternative to BLAST, which is 10,000 times faster than BLAST while being only 0.1– 1% less sensitive.
4. Access the tool from the **Tools menu** in the header > Map proteins to OrthoMCL with Diamond BLASTP (red arrow below)
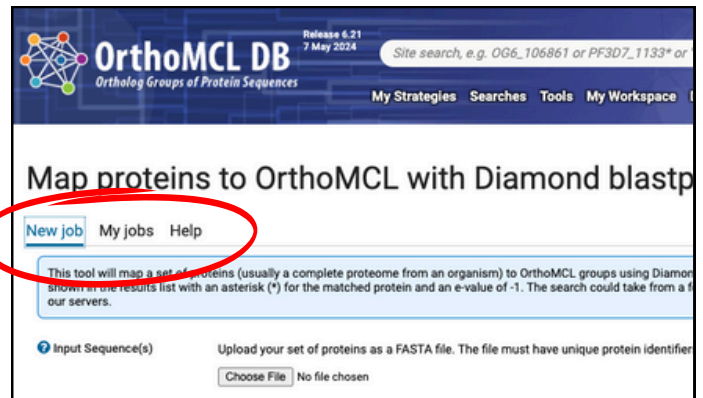


Note that you must be logged into your free OrthoMCL account to use this tool. Click on the person icon at top right (green arrow below) to log in or register.

**Layout of the Diamond BLASTP protein mapping tool**: There are three tabs (circled in red on right)

1. **New job**: Upload data here
2. **My jobs**: Table of all your previous jobs; these are saved in your account and persist between sessions
3. **Help**: Tips for using the tool



**Preparing your data:** Your set of proteins must be formatted as a plain text FASTA file.
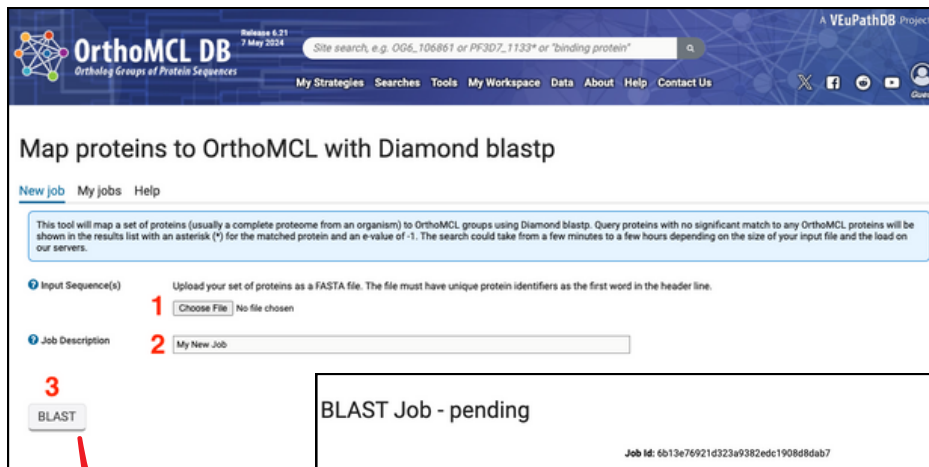
- Each protein in the FASTA file must have unique protein identifiers as the first word in its definition/header line.
- Header line must start with a greater than (>) symbol and end with a carriage return.

For example, the FASTA file on the right shows a set of predicted proteins from a newly sequenced Diatom species.



**Uploading data**: Do the following steps in the "New job" tab (refer to figure below)

1. **Input sequence**: Choose a FASTA-formatted data file with protein sequences from your computer
2. **Job description:** Add brief text describing your set of proteins
3. **BLAST**: Click on the button to start the job. You will see a message with a job ID assignment.

## Understanding the output:

The output page has two components
1. **The results table** (see below). This is a preview of the matching results for the first 100 sequences in your query file.
2. **A blue download button** at the top right (see red arrow below). The complete result can be downloaded as a tab delimited file with one best match for each query protein with the following columns:
    - Query_ID
    - Subject_ID
    - Orthogroup
    - Subject_description
    - Alignment_length
    - Percent_identity
    - e-value

Note: Unmatched query proteins are included in the results file without an OrthoMCL protein or group listed. For example, see red rectangle below.



Questions? Comments? Write to

help@veupathdb.org