MSBI 32400 – LAB 4
LARRY HELSETH, PHD AND
JASON EDELSTEIN

January 30, 2019

## Outline

**2**

- Install a reference proteome for local BLAST on VM
- Configure Linux environment to run BLAST
- Run searches
  - Output to text or console
  - Output to HTML
- Learn to run 2 way BLAST
- Identify some unknown sequences

MSBI 32400 Lab 4    1/30/2019

# Using wget

**3**

- See Chapter 6 of Vince Buffalo's book for a great discussion of **wget** and other commands for moving files.
  - Includes examples of how to use md5sum and other tools to check the integrity of large downloads

MSBI 32400 Lab 4    1/30/2019

# First let's get some reference data

**4**

- We'll use wget to retrieve the human reference proteome from UniProt (~7 MB file)
- Go to /data/ncbi-blast-2.7.1+/ and create a new directory called "db", then cd into that directory
- wget ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/reference_proteomes/Eukaryota/UP000005640_9606.fasta.gz
- Optionally browse to http://uniprot.org/ then click on Reference Proteomes, enter 9606, click Download and select All Proteins with Compressed format.
- gunzip UP000005640_9606.fasta.gz

MSBI 32400 Lab 4    1/30/2019

## Getting the reference data



MSBI 32400 Lab 4    1/30/2019

## How many proteins?

- Extract the fasta file using
  gunzip UP000005640_9606.fasta.gz
- grep the new fasta file for "^>" (lines that start with the greater than symbol, meaning they have FASTA annotation) and **pipe** that to the word count program **wc** with the -l (just show the number of lines)
- Include that information (and the wget command) in your Lab 4 README file

MSBI 32400 Lab 4    1/30/2019

## Let's add ncbi binaries to PATH

7

□ Tell your server where to find the BLAST binaries:
- ▫ PATH=$PATH:/data/ncbi-blast-2.7.1+/bin
- ▫ export PATH

MSBI 32400 Lab 4    1/30/2019

## "Hidden" files in Linux

8

□ If you type "ls -l" in your home directory you'll see some files. If you type "ls -la" you will start seeing things like ".bash_profile", ".bashrc", etc.
- ▫ Your terminal history is kept in ".bash_history"

□ The preceding "." hides the file or directory from ls unless you ask it to show all

□ Most are program or user profile files & directories

□ .bash_profile and .bashrc contains your environmental variables (loaded when you sign in)

MSBI 32400 Lab 4    1/30/2019

# BLAST lab

**9**

□ Need a hidden file for BLAST to work correctly

□ See **Recipe 11: Running Blast** from Safari on-line version of "**Managing Your Biological Data with Python**" by Allegra Via, Kristian Rother & Anna Tramontano, CRC Press 2014 (pp 431-435 print)

▪ (Files/Lab4)

□ Create a .ncbirc file in your <u>home</u> directory (/home/student)

▪ vi ~/.ncbirc (or nano or gedit)

MSBI 32400 Lab 4     1/30/2019

# Syntax for .ncbirc

**10**

□ ; Start the section for BLAST configuration
[BLAST]
; Specifies the path where BLAST databases are installed
BLASTDB=/data/ncbi-blast-2.7.1+/db

```
File  Edit  View  Search  Terminal  Help
[student@msbi32400lab5 ~]$ cat ~/.ncbirc
; Start the section for BLAST configuration
[BLAST]
; Specifies the path where BLAST databases are installed
BLASTDB=/data/ncbi-blast-2.7.1+/db
[student@msbi32400lab5 ~]$ █
```

MSBI 32400 Lab 4     1/30/2019

# Make a BLAST database

11

- □ cd /data/ncbi-blast-2.7.1+/db
- □ Use makeblastdb
  - ◻ Tell it what input file you're using (-in) and what type of database it is (-dbtype either prot or nucl)
- □ makeblastdb -in UP000005640_9606.fasta -parse_seqids -dbtype prot
- □ Look at the files in db (ls -la)

MSBI 32400 Lab 4    1/30/2019

# Let's download some files to test

12

From your /lab4/data folder:

- □ esearch -db protein -query "NP_000509" | efetch -format fasta > hbb.fasta
- □ esearch -db protein -query "NP_976312" | efetch -format fasta > myog.fasta
- □ Download "unknown.fasta" from Canvas under Files/Lab 4
- □ Download "unknown_fragment.fasta" from Canvas under Files/Lab 4

MSBI 32400 Lab 4    1/30/2019

# Running blastp (comparing proteins)

**13**

Check your PATH if you don't get this:

```
[student@MSBI32400Lab1 ~]$ blastp -h
USAGE
  blastp [-h] [-help] [-import_search_strategy filename]
    [-export_search_strategy filename] [-task task_name] [-db database_name]
    [-dbsize num_letters] [-gilist filename] [-seqidlist filename]
    [-negative_gilist filename] [-entrez_query entrez_query]
    [-db_soft_mask filtering_algorithm] [-db_hard_mask filtering_algorithm]
    [-subject subject_input_file] [-subject_loc range] [-query input_file]
    [-out output_file] [-evalue evalue] [-word_size int_value]
    [-gapopen open_penalty] [-gapextend extend_penalty]
    [-qcov_hsp_perc float_value] [-max_hsps int_value]
    [-xdrop_ungap float_value] [-xdrop_gap float_value]
    [-xdrop_gap_final float_value] [-searchsp int_value]
    [-sum_stats bool_value] [-seg SEG_options] [-soft_masking soft_masking]
    [-matrix matrix_name] [-threshold float_value] [-culling_limit int_value]
    [-best_hit_overhang float_value] [-best_hit_score_edge float_value]
    [-window_size int_value] [-lcase_masking] [-query_loc range]
    [-parse_deflines] [-outfmt format] [-show_gis]
    [-num_descriptions int_value] [-num_alignments int_value]
    [-line_length line_length] [-html] [-max_target_seqs num_sequences]
    [-num_threads int_value] [-ungapped] [-remote] [-comp_based_stats compo]
    [-use_sw_tback] [-version]

DESCRIPTION
  Protein-Protein BLAST 2.5.0+

Use '-help' to print detailed descriptions of command line arguments
[student@MSBI32400Lab1 ~]$
```

MSBI 32400 Lab 4    1/30/2019

# Let's BLAST some proteins

**14**

☐ blastp -query hbb.fasta -db /data/ncbi-blast-2.7.1+/db/UP000005640_9606.fasta

▫ This sends the output to console.  You can either pipe this to more or redirect to a file then view that with less or more

☐ blastp -query hbb.fasta -db /data/ncbi-blast-2.7.1+/db/UP000005640_9606.fasta -out mysearch.html –html

▫ firefox mysearch.html to view

▫ What is the E-value for the top hit?  (record in README)
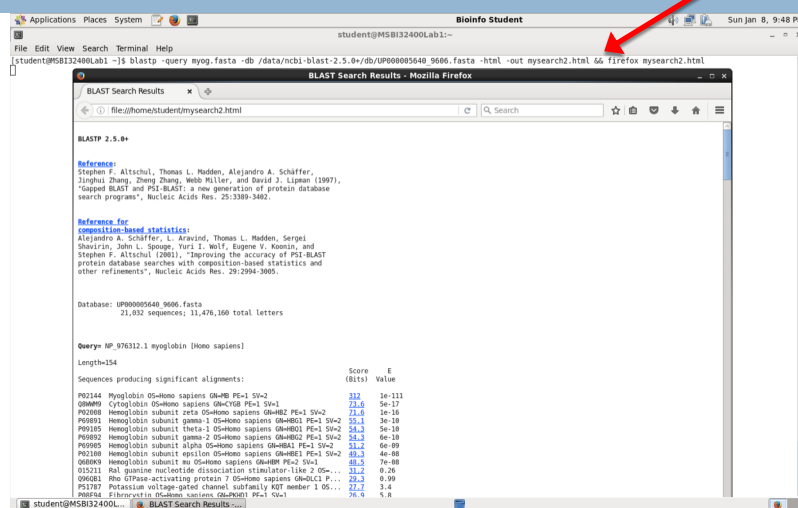
MSBI 32400 Lab 4    1/30/2019

# What is the unknown sequence?

**15**

- Assuming it's human, run blastp with unknown.fasta as your query against the UniProt database
- Record the top hit in the README you send to Jason

MSBI 32400 Lab 4    1/30/2019

# BLAST myog.fasta

Use '&&' to string together commands

**16**



MSBI 32400 Lab 4    1/30/2019

# BLAST two files



MSBI 32400 Lab 4    1/30/2019

# BLAST shows mismatches

□ Download unknown2.fasta from Files/Lab 4 on Canvas

□ Run blastp with this file as query

□ Look at the first hit in the HSP alignment.  Did it match all amino acids?  What is the base in the query and what is the base in the reference sequence?  Ignoring the first M in the sequence, what position is any potential mismatch?  (Hint- You've seen it before!)   Record what you found in the README file you send to Jason.

MSBI 32400 Lab 4    1/30/2019

## Repeat some of the searches on-line

19

- Upload unknown.fasta to
  https://blast.ncbi.nlm.nih.gov/Blast.cgi
- Run 16S ribosomal search against NR_119358.1
- Download and BLAST unknown_fragment.fasta
  - Look at the file by **head** to see if it's protein/nucleotide and use appropriate search engine
  - Assume it's human (9606)
  - What gene/protein is this from?  What is the E-value?

MSBI 32400 Lab 4    1/30/2019

## Homework

20

- Submit to Canvas or e-mail Jason (jasone@uchicago.edu) the README with the file information requested above before next class (2/13) with "**Lab #4**" in the subject line

MSBI 32400 Lab 4    1/30/2019