

# Bioinformatics Introduction

Background and much more material is well presented in Prof Girke's class [GEN240B](#), in particular [Sequence Alignments](#)

## Running Analysis

- How to run BLAST on command line
- How to setup data files and process
- Development of workflows

## Sequence search tools - BLAST

- BLAST is by far the most taught tool in Bioinformatics. I am not going to rehash this completely in this class.
- See NCBI's [Introduction to BLAST](#)
- One of 7 Million pages by Googling "[blast introduction tutorial](#)"

## BLAST on Biocluster

There are multiple flavors of BLAST (implementations). Focus on the latest version from NCBI (2.9.0+). Default on the cluster is 2.2.30+

We will make links to two files which are ORFs from two yeast species. Try this in UNIX

```
# setup some files to do some searches
mkdir BLAST_demo
cd BLAST_demo
ln -s /bigdata/gen220/shared/data/C_glabrata_ORFs.* .
ln -s /bigdata/gen220/shared/data/S_cerevisiae_ORFs.* .
ln -s /bigdata/gen220/shared/data/Yeast_chr2_ORFs.fa
```

Now we have some files, set them up for running BLAST. Our question is, what ORFs are similar at the DNA level between these two species.

```
module load ncbi-blast/2.9.0+ # load the module on the biocluster
makeblastdb -dbtype nucl -in C_glabrata_ORFs.fasta
ls
# C_glabrata_ORFs.fasta      C_glabrata_ORFs.fasta.nhr
# C_glabrata_ORFs.fasta.nin C_glabrata_ORFs.fasta.nsq
# Yeast_chr2_ORFs.fa
#head -n 7 Yeast_chr2_ORFs.fa > YBL001C.cds # get 1st seq for an example

# we do this because I checked and that sequence takes up the first 7 lines
```

*# of the file*

```
blastn -query Yeast_chr2_ORFs.fa -db C_glabrata_ORFs.fasta
```

## BLAST Running

Change the output format to tab delimited with `-outfmt 6` or `-outfmt 7`

```
$ blastn -query YBL001C.cds -db C_glabrata_ORFs.fa \  
-evalue 0.001 -outfmt 7 -out YBL001C-vs-Cglabrata.BLASTN.tab
```

This will query the 1 sequence and produce a tab delimited file.

If you provide a multi-FASTA format file with many sequences, each one will be queried and all the results concatenated together.

```
$ blastn -query Yeast_chr2_ORFss.fa -db C_glabrata_ORFss.fa \  
-evalue 0.001 -outfmt 7 -out yeast_chr2-vs-Cglabrata.BLASTN.tab
```

## BLAST: what are the tools

- `makeblastdb` - index a database (required to do once before searching)
- `blastn` - DNA/RNA to DNA/RNA search
- `blastp` - protein to protein search
- `blastx` - translated query (DNA/RNA) against protein database
- `tblastn` - protein query against translated (DNA/RNA) database
- `tblastx` - translated query and database (both in DNA/RNA but search in protein space)
- `blastdbcmd` - retrieve a sequence from a blast formatted DB

##BLAST: what are the cmdline options?

All the tools have documented command line options. Use `-h` or `-help` to get detailed info. Sometimes with no arguments will print documentation, other times will not.

```
$ makeblastdb
```

USAGE

```
makeblastdb [-h] [-help] [-in input_file] [-input_type type] \  
-dbtype molecule_type [-title database_title] [-parse_seqids] \  
[-hash_index] [-mask_data mask_data_files] [-mask_id mask_algo_ids] \  
[-mask_desc mask_algo_descriptions] [-gi_mask] \  
[-gi_mask_name gi_based_mask_names] [-out database_name] \  
[-max_file_sz number_of_bytes] [-logfile File_Name] [-taxid TaxID] \  
[-taxid_map TaxIDMapFile] [-version]
```

DESCRIPTION

Application to create BLAST databases, version 2.2.30+

Use '-help' to print detailed descriptions of command line arguments

=====  
##BLAST: what are the cmdline options?

\$ blastn -h

USAGE

```
blastn [-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] [-evaluate evaluate] [-word_size int_value]
[-gapopen open_penalty] [-gapextend extend_penalty]
[-perc_identity float_value] [-qcov_hsp_perc float_value]
[-xdrop_ungap float_value] [-xdrop_gap float_value]
[-xdrop_gap_final float_value] [-searchsp int_value] [-max_hsps int_value]
[-sum_stats bool_value] [-penalty penalty] [-reward reward] [-no_greedy]
[-min_raw_gapped_score int_value] [-template_type type]
[-template_length int_value] [-dust DUST_options]
[-filtering_db filtering_database]
[-window_masker_taxid window_masker_taxid]
[-window_masker_db window_masker_db] [-soft_masking soft_masking]
[-ungapped] [-culling_limit int_value] [-best_hit_overhang float_value]
[-best_hit_score_edge float_value] [-window_size int_value]
[-off_diagonal_range int_value] [-use_index boolean] [-index_name string]
[-lcase_masking] [-query_loc range] [-strand strand] [-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] [-remote]
[-version]
```

DESCRIPTION

Nucleotide-Nucleotide BLAST 2.2.30+

Use '-help' to print detailed descriptions of command line arguments

## BLAST: some key arguments

- -query - query file name (required)
- -db - database file name (require)
- -evaluate - set the evaluate cutoff
- -max\_target\_seqs - max number of hit seqs to show
- -num\_alignments - max number of alignments to show

- -num\_threads - number of threads (parallel processing to run, 8 will be faster than 2)
- -outfmt - specify a simpler format than the text format, try '-outfmt 6' for tabular format
- -subject - instead of doing a DB search, search for alignments between query sequence and 1 to many subject sequences. Useful when want to just see the alignment of 2 sequences already picked out from other analyses

## BLAST: Putting it all together

This is a script. e.g. run\_blast.sh

```
#!/usr/bin/bash
#SBATCH -p short --nodes 1 --ntasks 4 --mem 2G --job-name=BLASTN
#SBATCH --output=blastn.%A.log
module load ncbi-blast/2.9.0+

CPUS=$SLURM_CPUS_ON_NODE
if [ ! $CPUS ]; then
    CPUS=1
fi
if [ ! -f C_glabrata_ORFs.fa.nhr ]; then
    makeblastdb -in C_glabrata_orfs.fa -dbtype nucl
fi
blastn -query yeast_chr2_ORFs.fa -db C_glabrata_ORFss.fa \
-evalue 1e-5 -outfmt 6 -out yeastORF-vs-CglabrataORF.BLASTN.tab -num_threads $CPUS
```

Now submit this script

```
$ sbatch run_blast.sh
$ squeue -u $USER # check on your submitted job
```

## Other types of search tools

- **HMMER**
  - Identify conserved domains in a protein
  - Sensitive searches for distant homologs
  - phmmer can be of comparable speed to BLASTP
  - HMMs are a way to not just match a single sequence but match a pattern
- **FASTA**
  - Another tool like BLAST
  - Doesn't require formatting the database

- FASTA/SSEARCH are more full length optimal alignments instead of individual scoring pairs, a single best alignment generated
- Global alignment also with ggsearch

## Other seq search tools

- **Exonerate**
- Another aligner useful for cDNA to genome alignment and protein to genome alignment
- splice-site aware
- output harder to parse but there is a GFF-flavor output and parsers in some toolkits
- **USEARCH / VSEARCH**
  - fast, near-exact search tool
  - useful in microbiome short-read
- **DIAMOND**
  - fast, near-exact short read search tool
  - translated BLASTX search option to search proteins against a short read database

## Retrieving Sequences from databases

Some of this will be mentioned later in Genome Assembly lecture. But here are some details about different ways to retrieve sequences is locate here.

## Already downloaded data

`/bigdata/gen220/shared/data/Afum` has some already downloaded datasets.

## Remote databases

The International Nucleotide Sequence Database Collaboration ([INSDC](#)) are the joint databases for sequence and related biomedical data deposition. These are critical central tools for archive of sequence data for the scientific community. Often when we say “deposited in GenBank” we mean this central repository which, but there are data for Sequence Reads, Assemblies, annotated genomes, Gene Expression, and Individual sequence records.

## Downloading FASTA databases from NCBI, Uniprot

On HPC there are already databases installed and indexed for BLAST searches.

```
#SBATCH -N 1 -n 16
module load db-ncbi
module load ncbi-blast
# loads the current ncbi folder as env variables
# $BLASTDB and $NCBI_DB
# after loading this you can run blast without specifying
# the location of the databases
blastp -db nr -query seqs.fasta -out seqs-nr.blastp -num_threads 16 -evalue 1e-5
```

### FTP / Web downloads

#### NCBI

The NCBI databases include several useful resources. Note these are now giant databases in some places so care is needed in whether you can download these to your own folder and if it already exists on cluster you should try to use those.

- [swissprot](#)
- [nr](#) - non-redundant protein seqs, very large ...
- [nt](#) - non-redundant nucleotide seqs, very large ...
- [env\\_nr](#) - environmental protein seqs
- [env\\_nt](#) - environmental nucl seqs

Another resource that is helpful is [Refseq](#) which are *somewhat verified* sequences from genomes.

To for example get all fungal refseq proteins use the `lftp` tool. Here is an interactive session:

```
lftp ftp://ftp.ncbi.nih.gov/refseq/release
lftp> cd fungi
lftp> mget fungi.*.protein.faa.gz
lftp> exit
pigz -dc *.faa.gz > refseq_fungi.faa
```

#### Uniprot

To Download uniprot\_swissprot database via ftp protocol. See <https://www.uniprot.org/downloads> for more download files available including uniref which is a set of clustered sequences at 100%, 90%, and 50% identity which can reduce the size of the total protein database but still leave representatives.

- [uniprot\\_swissprot](#)
- [uniref50](#)

The [UniRef50 database](#) for this as it isn't too big but useful for some relatively fast searching and more comprehensive than swissprot for taxonomic representation.

## Downloading from SRA

The easiest way to download from SRA is using the [sra\\_download.pl](#) script. This is already installed on the cluster in `/bigdata/stajichlab/shared/bin/sra_download.pl`.

To run this create a file that has a line for each SRA accession number. I have called it `sra.txt` here. This will download the fastq for all the data sets in the file and create a folder for each one.

```
#SBATCH -p short -N 1 -n 2 --mem 4gb
module load aspera
/bigdata/stajichlab/shared/bin/sra_download.pl --ascp --id $ASPERAKEY sra.txt
```

## Downloading sequence records from GenBank

You can use several tools to download accessions from genbank. It does require certain versions of perl are installed or conda.

```
module load perl/5.20.2
bp_download_query_genbank.pl --query 'AY295118.1'
>AY295118 Parmelia ernstiae voucher MAF 9805 tubulin gene, partial cds.
GAGGACATTCCTCCATAATGTGATACGTAGCTCACAGCTTTCAAGGCTTCAAACAACAAA
TATGTTCCCTCGTGCCGTA CTCTCGATCTCGAGCCTGGTACCATGGATGCTGTCCGCGCT
GGTCTTTTGGCCAGCTTTTCCGACCCGATAACTTCGTATTTGGTCAATCTGGTGCTGGT
AATAATTGGGCTAAGGGTCATTACACCGAGGGTGCAGAATTGGTGGACCAAGTCTCTCGAT
GTTGTGCGTCGAGAGGCTGAAGGATGCGACTGCCTCCAGGGCTTCCAGATCACGCACTCC
CTCGGTGGTGGAACTGGTGCTGGTATGGGTACGCTTTTGATCTCGAAAATCCGTGAGGAG
TTCCCAGATCGTATGATGGCTACATTCTCCGTGGTTCCTTACCAAAGGTATCCGACACT
GTTGTGGAGCCATACAACGCTACTCTCTCCGTGCATCAATTGGTCGAGAACTCGGATGAG
ACCTTCTGTATCGATAATGAGGTTGGTCAAGTGCATTTTTTTCACAGAGGCGCAAGGACT
GATATGTCAATCTAGGCGCTCTATGACATTTGCATGCGCACCTCAAGCTCTCCAACCCA
TCCTACGGGGATCTTAACCACCTTGTCTCCGCGGTCAATGTCTGGTGTACCACCTGCCTC
CGTTTCCCGGTCAACTCAATCCGACCTTCGAAAAC TAGCCGTCAACATGGTCCCATTT
CCCCGTCTACATTTCTTTCATGGTTGGCTTCGCACCTTTACCAGCCGAGGTGCTAACTCA
TTCCGTGCGGTACGCTACCAGAATTGACCCAACAAATGTACGAC
```

If you want to retrieve a number of sequences at a time you can specify a query. Below are the options for running the tool. If you want to retrieve data from protein database you need to specify the database with `--db` option.

```
bp_download_query_genbank --query "Neurospora[ORGN]" --db nucest -o Ncrassa_ESTs.fa --format
```

#### Other options

Provide ONE of:

```
--q --query query string OR
--queryfile profile file with query OR
--gi --gis --gifile file with list of GIs to download
```

#### Database type:

```
--d --db database (nucleotide [default], nucest, protein, )
--o --out --outfile output file (results are displayed on screen otherwise)
--f --format sequence file output format (fasta by default)
--v --verbose debugging output
```

#### Query options

```
--maxids maximum number of IDs to retrieve in a set (100 at a time by default)
--reldate
--maxdate maxdate for a record
--mindate minimum date for record
--datetype edat or mdat (entered or modified)
```

## Specialized fungal databases

### FungiDB

The [FungiDB](#) project provides access to a set of Fungal genomes loaded into this system. The resources for downloads are available at [this link](#) which includes current and previous releases. Data sets are organized by Abbreviations of genus + species and strain name. For example the Genome, CDS, Protein, and Transcripts associated with the *Neurospora crassa* OR74A strain are available from [this link](#).

### JGI

The JGI [Mycocosm](#) provides one of the largest collection of fungal genomes through the sequencing and annotation project. There are more than **1000 genomes available** through [several interfaces](#) hosted by the JGI. **Some scripts** for automation of downloads are needed to directly extract data from the site onto linux clusters. GLOBUS and other functionality do exist for dataset downloads as well.



## Ensembl

Ensembl provides a nearly complete set of public deposited genomes into GenBank organized by major domains. The [Ensembl Fungi](#) is a portal with access to thousands

## Saccharomyces or Candida Genome Database

See [SGD](#) and [CGD](#) for main site for genome browsers, comparative tools, and access to primary sequence data associated with these fungi.

The FTP site for yeast see [ftp://ftp.yeastgenome.org/sequence/S288C\\_reference](ftp://ftp.yeastgenome.org/sequence/S288C_reference) for example which has access to the yeast ORFs [proteins](#) and [coding sequence](#) as well as many other resources like upstream promotor files and other features. Data from multiple *Candida* species is available from <http://candidagenome.org/download/sequence/>

## Local databases

Once you have data files downloaded you can use these.

### FASTA Files

#### HMMER esl-sfetch

I find esl-sfetch one of the better tools for fasta file indexing. Database must be indexed first.

```
module load hmmer
esl-sfetch --index database.fasta
# fetch record based on single ID passed in on cmdline
esl-sfetch database.fasta accession > accession.fa
# fetch multiple records based on a list of IDs passed in a file (note the -f option)
esl-sfetch -f database.fasta list_of_ids > seqs.fa
# fetch list of IDs passed in on STDIN using a pipe and specifying the input file as '-'
cat ids_to_fetch | esl-sfetch -f database.fasta - > seqs.fa
```

#### cdbfasta

cdbfasta ([Constant database](#)) is a useful for indexing fasta and fastq files for retrieval by sequence ID.

```

module load cdbfasta
cdbfasta database.fasta

# to retrieve
echo "A" | cdbyan database.fasta.cidx > A.fa
cat list_of_ids | cdbyan database.fasta.cidx > retrieved.fa

```

### samtools (1.9 and later)

Samtools provides indexing and retrieval of FASTA

```

#SBATCH -p short -N 1 -n 4
module load samtools

```

```

# index file
samtools faidx DNA_sequences.fasta

```

To retrieve a sequence read after the file is indexed, where accession is the first text after the > in FASTA file, eg scaffold\_1 is the accession in the following:

```

>scaffold_1
TGCATGTCTAAGTATAAGCAATTATACCGTGAAACTGCGAATGGCTCATTAAATCAGTTATCGTTTATTTGATAGTACCTTACTACTTGGAA

```

To retrieve this sequence from the indexed file use this (specify DNA\_sequences.fasta if you used the uncompressed file).

```

module load samtools
samtools faidx DNA_sequences.fasta scaffold_1

```

### BLAST indexing

The BLAST indexing to setup a database for sequence alignment and searching also allows retrieval of sequences by identifier.

```

module load ncbi-blast
# index a nucleotide database
# to index a protein database change -dbtype from 'nucl' to 'prot'
makeblastdb -in sequences.fasta -dbtype nucl -parse_seqs

# to retrieve sequences
blastdbcmd -entry ACCESSION -db sequences.fasta -out ACCESSION.fasta

# use this database for sequence searches
# report the output as tab delimited format (outfmt 6)
blastn -query myquery.fasta -db sequences.fasta -out myquery-vs-seqs.BLASTN -outfmt 6 -eval

# Do a protein db search

```

```
blastp -query myquery.fasta -db protseqdb.fasta -out myquery-vs-seqs.BLASTP -outfmt 6 -eval
```

*# many other options for BLAST using blastx, tblastn, tblastx and many more options for run*

## DIAMOND indexing

**DIAMOND** is a rapid aligner for protein and translated searches which can operate on short sequence reads as well as assembled genomes.

DIAMOND does not provide a way to extract sequences back out from these indexed databases. Will report a tab delimited output file (m8 stands for the OLD NCBI -mformat 8 output which is tab delimited).

```
module load diamond
makedb --in my_protein_db.fasta -d mydb
diamond blastx -d mydb -q reads.fna -o hits.m8
```

## Short read aligner database indexing

Indexing DNA database for aligning short read DNA sequences against this database (usually a genome).

Indexing for *bwa* in order to setup searches:

```
module load bwa
bwa index database.fa
```

Indexing for *bowtie2*:

```
module load bowtie2
bowtie2-build database.fa database
```

Indexing for *gmap/gsnap*:

```
module load gmap
gmap_build -D genome_index -d genome_name database.fa
```

Indexing for *kallisto* (RNASeq analysis):

```
module load kallisto
kallisto index -i transcripts.idx transcripts.fasta
```

## FASTQ Files

### **cdbfasta**

**cdbfasta** ([Constant database](#)) is a useful for indexing fasta and fastq files for retrieval by sequence ID.

```

module load cdbfasta
cdbfasta -Q reads.fastq

# retrieve seqs by
echo "ACCESSION" | cdbank reads.fastq.cidx > fetched_read.fq
cat list_of_ids | cdbank reads.fastq.cidx > retrieved.fq

```

## samtools

Samtools provides indexing and retrieval of FASTQ Files.

If the file is compressed (.gz) it must be compressed with the bgzip tool - which is part of the htlib package. So if the file exists already as a compressed file you need to uncompress and recompress with bgzip.

```

#SBATCH -p short -N 1 -n 4
module load samtools
pigz -d READFILE.fq.gz
bgzip --threads 4 READFILE.fq

```

```

# now index
samtools fqidx READFILE.fq.gz

```

```

# you can also index an uncompressed file
samtools fqidx READFILE.fq

```

To retrieve a sequence read after the file is indexed, where accession is the first text after the @ in FASTQ file, eg ERR1309286.4 is the accession in the following:

```

@ERR1309286.4 H4:C3F32ACXX:2:1101:1849:2436/1
CTCTATTTCATCACGTTTCGAGAAGATCGCTACGCTTATCGAATTCAGATTATCATTGTCCGCTTCAACTTCTAGAGAACTGTGCATGAT
+
@CCFFFFFFGHHGHJJIIJIHJIIJJIIIIIIJIIIFJIIJJGIGIEHGIHIIGJIIIIJJJJJIHGF:BDDEEEDEEA>>CDDCDDDDDD

```

To retrieve this sequence from the indexed file use

```

module load samtools
samtools fqidx READFILE.fq.gz ERR1309286.4

```

## Databases

Lots of the data are in prim