MCB 5472 Assignment #4:
Introduction to command line BLAST
February 12, 2014

## Assignment feedback

• Everyone is doing very well!
• Most people lose marks because they have not read the question close enough (e.g., not handing in pseudocode for Assign. #2)
• Check that your output files match your input and that they contain what you think they should

## Code hints:

• Filehandles should be in block capitols

```
open (INFILE, $ARGV[0]); # good
open (infile, $ARGV[0]); # less good
```

• Will work otherwise but not under "use warnings" and "use strict" pragmas, i.e., sloppy code

## Code hints:

• Align you code blocks properly
• Purpose: so you can intuitively see your code logic

```
foreach $word (@array){
    if ($word =~ /fox/){
        print "found the fox";
    }
}
```

## Code hints:

• When using pattern matching, you should always consider possible exceptions in your input file

```
$line =~ /^>/; # match fasta header
$line =~ />/; # matches any line with a ">" character
$line =~ tr/ACGT/TGCA/; # compliments high quality sequence
$line =~ tr/ACGTacgt/TGCAtgca/; # compliments high and low quality sequence
$line =~ /^[ACGT]/; # line starts with a nucleotide
$line !~ /^>/; # any line not a fasta header, accommodates degenerate bases: N, V, B, H, D, K, S, W, M, Y, R
```

## Code hints:

• "\s" means "any white space character"
• Includes:
```
" " # space
"\t" # tab
"\r" # return
"\n" # new line
```

## Code hints:

- Filehandles should be in block capitols

```
open (INFILE, $ARGV[0]); # good
open (infile, $ARGV[0]); # less good
```

- Will work otherwise but not under "use warnings" and "use strict" pragmas, i.e., sloppy code
- Keep your tabs aligned
- Matching /^>/ vs />/
- tr/ATCG/TAGC/ vs tr/ATCGatcg/TAGCtagc/
- "\s" matches: " ", "\t", "\r", "\n"
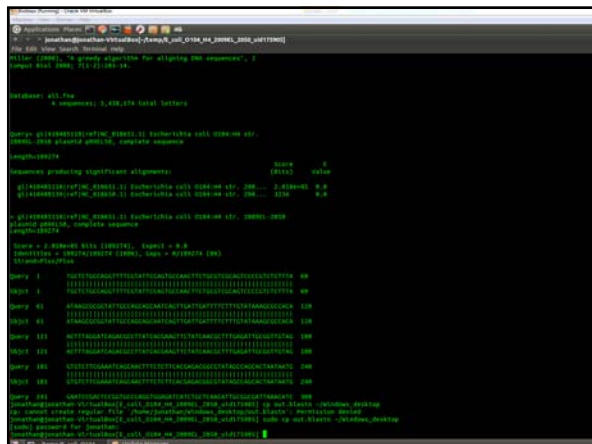
## Command line BLAST

- We will be using BLAST+
- Need to run on the Biotechnology Center server
  - Preinstalled on Biolinux so can be run locally
- Two parts to every BLAST
  1. Format the BLAST database
  2. Perform the BLAST itself

## Formatting a BLAST database

- `[jlklassen@bbcsrv3 ~]$ makeblastdb -in [name of input file] -dbtype [either 'nucl' or 'prot']`
  - e.g., `[jlklassen@bbcsrv3 ~]$ makeblastdb -in all.fna -dbtype nucl`
- Produces:
  - nucleotide: `[name].nhr`, `[name].nin`, `[name].nsq`
  - protein: `[name].phr`, `[name].pin`, `[name].psq`
  - where `[name]` is whatever was entered for the `makeblastdb -in` flag
- For help: `[jlklassen@bbcsrv3 ~]$ makeblastdb -help`

## Running BLAST

- `[jlklassen@bbcsrv3 ~]$ blastn – query [query file name] –db [database name]`
  - e.g., `[jlklassen@bbcsrv3 ~]$ blastn – query NC_018651.fna –db all.fna`
- For other BLAST flavors: replace `blastn` with `blastp`, `blastx`, `tblastn` or `tblastx`
- For help: `[jlklassen@bbcsrv3 ~]$ blastn –help`
- Multiple fasta file can be used as query
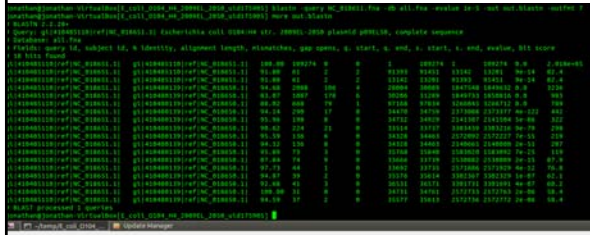  - multiple BLAST outputs in the same output file



## Other helpful BLAST options

- `–evalue [maximum evalue threshold]`
- `–out [output file name]`
- `–outfmt [0 for normal alignment format; 7 for easy to parse table format]`

## Tabular output format

- Separated by tabs ("\t")

```
@blastline = split "\t", $line;
```



## This week: Question #1

- Last week's complete genomes each had plasmids
- Are the plasmids from each organism with a complete genome homologous?
- Are the plasmids present in any of the draft genome sequences?
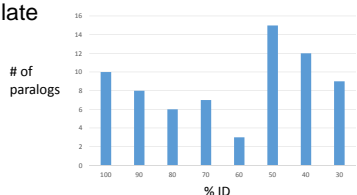
## Practically:

- Use BLASTn to compare plasmids with each other
- Use BLASTn to find homologous sequence to each plasmid type in the draft genomes
- Use your judgment to infer homologs – this is ultimately subjective and needs to be defended!
- YOU DO NOT NEED TO WRITE PERL SCRIPTS FOR THIS (unless you want to)

## This week: Question #2

- Find paralogous genes and proteins in the complete *Escherichia coli* O104:H4 str. 2009EL-2050 genome and its plasmids
- Compare number of gene and protein paralogs
- Tabulate paralog age estimated from their percent BLAST similarly

## Practically:

- Download the genes from NCBI
- BLAST all genes & proteins against each other using blastn and blastp (respectively)
- Round percent identity to the nearest 10% and tabulate



## Rounding in perl is not trivial

- `int`: truncates decimals to integers
  ```
  print int(1.4); # returns 1
  print int(1.6); # returns 1
  ```
- So to round:
  (1) Divide to convert to a decimal
  (2) Add 0.5
  (3) Apply `int`
  (4) Multiply to revert divide

  ```
  $number1 = 19;
  print int($number1/10+0.5)*10; # returns 20
  $number2 = 11;
  print int($number2/10+0.5)*10; # returns 10
  ```

## Discuss: how will we tabulate rounded %IDs

```
%tabulated_IDs = (
        100 => "",      90 => "",       80 => "",
        70  => "",      60 => "",       50 => "",
        40  => "",      30 => "",       20 => "",
        10  => "",       0 => "",
); # set up output hash
$rounded_ID = int($blast_table[2]/10+0.5)*10; # perform rounding
$tabulated_IDs{$rounded_ID}++; # tabulate in output hash
or:
$tabulated_IDs{$rounded_ID} = $tabulated_IDs{$rounded_ID} + 1;
# same thing as above
```

## To submit for next week

- Your conclusions from your results and your justification of them (esp. question #1)
- Your scripts and/or representative terminal commands
  - detailed enough that I can reproduce your results
- You don't have to submit input files or pseudocode