

Unix, Perl and BioPerl

Session 1: Introduction to Unix for Bioinformatics

Exercise 1: BLASTing ESTs against a RefSeq database

Goal: Learn the most common Unix commands while manipulating sequence files and “identifying” some rat ESTs by searching RefSeq, an annotated database, with BLAST.

Some commands written on multiple lines should be entered as a one-line commands.

See <http://jura.wi.mit.edu/bio/education/bioinfo2006/unix-perl/> for course page

#	To do / To answer , Commands , and Comments
0	<p>Mac OS X: Open X11 (In the dock, or Applications >> Utilities >> X11) or the terminal on your computer. If you’re running Mac OS X, you’re in Unix now.</p> <p>Windows: Open Cygwin (preferred) or SSH</p> <p>See the course page for links to more information about these applications.</p>
1	<p>Open your home account on hebrides. <code>ssh -l username hebrides.wi.mit.edu</code></p> <p>replacing ‘username’ with yours. The switch ‘-l’ is the 12th letter (not “one”). You will be prompted for your password. If it’s the first time for connecting, you may get a message about a “RSA key fingerprint” and then asking, “Are you sure you want to continue connecting?” You can answer “yes”.</p> <p>Also if it’s the first time connecting, change your password.</p> <p>To change the password whenever you want: <code>passwd</code></p>
2	<p>What is the full path to your home directory? <code>pwd</code></p> <p>“print working directory” show the full path to your current location</p>
3	<p>What files are in your home directory? <code>ls</code></p> <p>“list” all files (except hidden files)</p>

4	<p>What files (including hidden files) are in your home directory? How big are they?</p> <pre>ls -a ls -al (the letter l)</pre> <p>The -a option will also show files starting with a dot; the -l option shows a long listing</p>
5	<p>What's in these files?</p> <pre>more myfile</pre> <p>“myfile” is replaced by any name you got with the "ls" command</p>
6	<p>Create a directory called "unix_class":</p> <pre>mkdir unix_class</pre> <p>“make directory”</p>
7	<p>Go to the "unix_class" directory:</p> <pre>cd unix_class</pre> <p>“change directory”</p>
8	<p>Make directories called “rat-ests” and “dbs”, and go to the “dbs” directory:</p> <pre>mkdir rat-ests mkdir dbs</pre>
9	<p>Change to the ‘dbs’ directory:</p> <pre>cd dbs</pre>
10	<p>(Optional) To get a preview of what you’ll be downloading, open a web browser to view the NCBI’s FTP site. Point it to the rat RefSeq sequence “database”:</p> <pre>ftp://ftp.ncbi.nih.gov/refseq/R_norvegicus/</pre> <p>You can view the “README” file and note the “mRNA_Prot” directory. Note that this isn’t really a database but rather a set of multiple sequence files.</p>
11	<p>Download the README file, renaming it as refseq_README:</p> <pre>wget ftp://ftp.ncbi.nih.gov/refseq/R_norvegicus/README -O refseq_README</pre> <p>This is a one-line command. The -O refseq_README (the “O” is the letter, not a zero), which is optional, lets you the rename the downloaded file. In any cases of poor network connections to NCBI, README can be copied from /home/george/unix/dbs/ :</p> <pre>cp /home/george/unix/dbs/README ./refseq_README</pre>
12	<p>Download the “rat.rna.fna.gz” file and rename it to rat.fna.gz. This is a fasta-format file of all rat RefSeq gene sequences:</p> <pre>wget ftp://ftp.ncbi.nih.gov/refseq/R_norvegicus/mRNA_Prot/</pre>

	<pre>rat.rna.fna.gz -O rat.fna.gz</pre> <p>This is a one-line command. In any cases of poor network connections to NCBI, <code>rat.fna.gz</code> can be copied from <code>/home/george/unix/dbs/:</code></p> <pre>cp /home/george/unix/dbs/rat.fna.gz .</pre>
13	<p>Check to make sure you downloaded what you wanted to get:</p> <pre>ls</pre> <p>You should have 'refseq_README' and 'rat.fna.gz'</p>
14	<p>Look at the README one screenful at a time:</p> <pre>more refseq_README</pre> <p>Hit the space bar to advance to the next screenful or 'q' to quit 'more'</p>
15	<p>Unzip the sequence file. What's the file called now?</p> <pre>gunzip rat.fna.gz</pre> <pre>ls</pre> <p>It's generally assumed that a file ending in .gz needs to be unzipped; the opposite is gzip.</p>
16	<p>How big is the sequence file?</p> <pre>ls -l</pre>
17	<p>How would you list files in order of modification time? (Consult the man pages for ls, using the space bar to advance and "q" to quit)</p> <pre>man ls</pre> <p>Extra credit: How would you list in reverse order of modification time (from oldest to newest)?</p>
18	<p>Look at rat.fna to check that it's in fasta format</p> <pre>more rat.fna</pre> <p>Fasta format requires a one-line header (starting with ">") followed by sequence</p>
19	<p>What are the arguments to use for "grep"?</p> <pre>grep</pre> <p>"general regular expression parser" – very useful!</p>
20	<p>Use grep to print all the header lines into a file called rat.headers:</p> <pre>grep ">" rat.fna > rat.headers</pre> <p>">" marks the beginning of a fasta-format sequence</p>
21	<p>Check out the new file to be sure it looks okay at the beginning and the end</p>

	<pre>head rat.headers tail rat.headers</pre> <p>Add the option <code>-n</code> to print “n” lines with “head” or “tail”</p>
22	<p>How many sequences are in the sequence file?</p> <pre>grep ">" rat.fna wc -l</pre> <p><code>wc</code> (“word count”), with the <code>-l</code> option (for “lines”), prints the number of lines.</p>
23	<p>Make a BLAST database of the <code>rat.fna</code> sequence file using the “formatdb” command.</p> <pre>formatdb -i rat.fna -p F -o T</pre> <p>“formatdb -“ prints all options. The options used here are the minimal/usual ones. The “-o T” option indexes the sequences in such a way that, given accession(s) you can extract sequence(s) from the big file using ‘fastacmd’.</p>
24	<p>What files have been created? What does the log file say?</p> <pre>ls more formatdb.log</pre> <p><code>formatdb.log</code> will show any <code>formatdb</code> errors. You can safely ignore any errors that look like <code>CoreLib [002.003] FileOpen(".formatdbrc","r") failed</code></p> <p>How many sequences were indexed?</p>
25	<p>Change to the ‘rat-ests’ directory</p> <pre>cd ../rat-ests</pre> <p>Remember that ‘..’ means up one level in the directory tree</p>
26	<p>Get a file of ESTs from <code>/home/george/rat-ests</code> and place it into the directory “rat-ests”</p> <pre>cp /home/george/rat-ests/* .</pre> <p>Note the final “.” which indicates that you want to copy the files into the current directory.</p>
27	<p>Extract the first sequence and place it into a file by itself.</p> <pre>head ests.fa head -8 ests.fa > est1.fa</pre>
28	<p>Run BLAST on the single sequence, using an expect cutoff of 0.05, printing text output (only the best 5 hits):</p> <pre>blastall -i est1.fa -d ../dbs/rat.fna -p blastn -e 0.05 -v 5 -b 5 -o est1_blast.txt</pre> <p>This is a one-line command. The command “blastall” shows and describes all options. What do the options mean in the command above?</p>
29	<p>Remembering that you use the <code>↑</code> to get back to the previous command, run a similar BLAST</p>

	<p>search but with a default expect value cutoff and generate tab-delimited output</p> <pre>blastall -i est1.fa -d ../dbs/rat.fna -p blastn -v 5 -b 5 -o est1_blast_tab.txt -m 9</pre> <p>This is similar to the previous command, but with “-m 9” added. Running BLAST on est1.fa with “-m 8” instead of “-m 9” performs the exact same search but creates an output file without the header line.</p>
30	<p>Open est1_blast.txt in pico (a simple text editor), using ^X (control-x) to exit.</p> <pre>pico est1_blast.txt</pre>
31	<p>Extract a sequence (ex: NM_199463) from the BLAST database:</p> <pre>fastacmd -s NM_199463 -d ../dbs/rat.fna</pre> <p>This only works if you had used the “-o T” option when indexing rat.fna with formatdb</p>
32	<p>Make a file with the five sequence IDs from est1_blast.txt, and extract these sequences from rat.fna</p> <pre>tail -5 est1_blast_tab.txt cut -f2 > list1.txt fastacmd -i list1.txt -d ../dbs/rat.fna > myseqs1.fa</pre> <p>The command ‘cut’ can cut out a field (by number, after the ‘-f’) from a tab-delimited file. Accessions (ex: NM_133594) or GIs (ex: 19424297) can be used (or both).</p>
33	<p>BLAST the set of ESTs with standard text output</p> <pre>blastall -i ests.fa -d ../dbs/rat.fna -p blastn -e 0.05 -v 5 -b 5 -o est_blast.txt</pre> <p>This is a one-line command. Note that BLAST is very fast when searching a database smaller than the default “nt” database</p>
34	<p>BLAST the set of ESTs with tab-delimited output:</p> <pre>blastall -i ests.fa -d ../dbs/rat.fna -p blastn -e 0.05 -v 5 -b 5 -T F -m 9 -o est_blast_tab.txt</pre>
35	<p>Any questions?</p>
36	<p>Logout from hebrides and return to your desktop terminal:</p> <pre>logout</pre> <p>Make sure the “command prompt” no longer shows something like “username@hebrides”.</p>
37	<p>Copy the BLAST output files from hebrides to your laptop using ‘scp’ (secure copy):</p> <pre>scp username@hebrides.wi.mit.edu:/home/username/unix_class/rat-ests/* .</pre>

	replacing username with yours. This is a one-line command.
38	Look at the files you downloaded with Unix commands or with Mac/Windows software. Note that the file(s) will be in the directory where you ran the last command, probably C:/cygwin/home/student (Windows) or /Users/student (Mac)
39	Delete any of your files from the laptop. Thanks!

Notes:

1 – If, when trying to use the ‘pico’ text editor, you get a message that hebrides or barra doesn’t know anything about your terminal, use the command

```
setenv TERM vt100
```

and try again with pico.