## SYSB 3036 W01: Introduction to Unix

Frank Aylward[1]

[1]Virginia Tech

Frank Aylward
Virginia Tech

PROTOCOL STATUS

**Working**

We use this protocol in our group and it is working

---

**start**

1   A Unix or Unix-like operating system generally has a wide variety of build-in functions that are extremely useful for navigating between folders, exploring the contents of files, and getting various summary statistics that are useful before undertaking a bioinformatic analysis.

Here we will explore some of the more useful Unix commands that you will find useful throughout this course and future assignments. Make sure you get comfortable using these commands, since they will make your life a lot easier later on.

If you are ever curious about how to use a command, you can type "man" right before the command name and hit enter ("man" is short for "manual" in this case, so it will give you the command's manual).

**mkdir**

2   **mkdir**

The "mkdir" command will make an empty directory (folder) with a specified name. Often we may wish to organize our work into separate folders, so it's nice to be able to easily create new ones.
Let's start by creating a folder called "test_project"

**mkdir test_project**

**ls**

3   **ls**

Now that we've created a new folder, it would be nice to know that it exists.
To do this we can use the "ls" command, which will list all folders and files present in our current folder.

**ls**

or to get a list-like output with human-readable values

**ls -lh**

You should see the folder "test_project" there.
It's generally a good idea to run "ls" after every step of a tutorial just to see what new files may have been created.

**autocomplete in the command line**

4   One nice thing about the command line is that it will try to autocomplete the name of any file or folder you type in if you hit the "tab" button.

---

So, for example, if you type "ls test_" and then hit a tab, it should autocomplete to "ls test_project".

This is handy since it is often tedious to always type exact file names in the command line, especially since even one typo will results in an error message.

If there is any ambiguity about how to autocomplete, the command line will autocomplete up to the first ambiguity. So, for example, if you had two folders called "test_project_1" and "test_project_2", you could type "te" and then hit tab, and the command line would autocomplete to "test_project_". Try creating another folder called "test_project_2" and giving this a try.

**writing to output**

5 **cd and pwd**

Now let's say we want to move into the folder test_project so we can do something there. To navigate between folders we can use the "cd" command.

**cd test_project**

and to move back out:

**cd ..**

Or to move back to our home folder immediately, just type

**cd**

The text before your cursor in the command line should tell you the location of the folder in which you are currently located (aka, your PATH). However, sometimes full PATHs can get long and may be truncated. If ever want the command line to print out what your full PATH is, use "pwd".

**pwd**

**wget**

6 **wget**

One command that is very useful is the "wget" command. This is a standard Unix command that will let you download a file directly to whatever folder your command line is situated in. There are many options, but all that is required is the full URL and an internet connection.

This is useful if you want to download genome files from a public server such as the National Center for Biotechnology Information (NCBI). On NCBI the URL for the genome of the bacterium Yersinia pestis is:

**ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/009/585/GCF_000009585.1_ASM958v1/GCF_000009585.1_ASM958v1 _genomic.fna.gz**

So if I want to download this file all I need to type is:

**wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/009/585/GCF_000009585.1_ASM958v1/GCF_000009585.1_ASM958v1 _genomic.fna.gz**

**gunzip**

7 **gunzip**

You will note that the file we just downloaded had the strange ending ".gz". That is because this is a compressed gzip file that we will need to uncompress in order to inspect further.
Do this this we can use the "gunzip" command:

**gunzip GCF_000009585.1_ASM958v1_genomic.fna**

And if you run "ls" now you should see the same file without the .gz ending.

The new file has a ".fna" ending, which stands for FASTA Nucleic Acid. This contains the raw genomic information for the Yersinia pestis genome.

## head and tail

8  **head** and **tail**

Sometimes files can be too large and unweildy to open up directly, but we will still want to take a look inside to see what the file format looks like. We can use the "head" and "tail" commands to quickly look at the first and last few lines of the file, respectively.

**head GCF_000009585.1_ASM958v1_genomic.fna**

and

**tail GCF_000009585.1_ASM958v1_genomic.fna**

Based on this you can get an idea of what FASTA format looks like. Generally, there are header lines that start with a ">" symbol and have the name of the sequence, and subseqent lines have the actual DNA sequence (ATGC letters).

## grep

9  **grep**

Grep is perhaps the most popular and widely used Unix command. It is essentially a robust search tool that will find occurrences of a pattern in a given file.

Given the genome file we downloaded, maybe we want to know how many lines start with ">", since this will tell us how many distinct DNA sequences there are in the file.
To do this we can type:

**grep "^>" GCF_000009585.1_ASM958v1_genomic.fna**

Here, the quotes surround the pattern we are searching for, and the "^" symbol specifies that we only want ">" symbols that start at the beginning of a line.

## wc

10  wc

The "wc" command gives us the size of a given file. By default it gives us the newline count, the word count, and the byte count. If we want these statistics for our genome file, we can simply type:

**wc GCF_000009585.1_ASM958v1_genomic.fna**

## pipes

11  Pipes, or the "|" symbol

One very useful aspect of Unix commands is that we can "pipe" the output of one command directly into another. This can simplify workflows substantially, since we don't always have to collect the output of one command and put it into another.

For example, let's combine two commands we already have experience with. Let's say we want to use "grep" to find all lines in our genome file that start with a ">", and then let's ask "wc" to count the occurrences for us:

**grep "^>" GCF_000009585.1_ASM958v1_genomic.fna | wc**

This output is simple enough now, but you can imagine if you had thousands of header lines in a file you would not want to count them by hand.

**12** Writing to output with the ">" symbol

Lastly, if we want to write the output of a command to a specific file, we can do so with the ">" symbol (this is unrelated to the use of this symbol in FASTA files).

For example, perhaps we want to record only the FASTA header lines and put them in a file called "fasta_headers.txt". We can do this with:

`grep "^>" GCF_000009585.1_ASM958v1_genomic.fna > fasta_headers.txt`

Now if you use "ls" and "head" you should see the new file has been created, and the contents are what we would expect.