

## Unix 101 Part 1 - files and directories

Concept: Terminal lets you access the Unix system that lies beneath the fancy Mac GUI.

### Before we begin

1. Download the file  
[https://research.nhgri.nih.gov/Training/unix/downloads/UnixClass\\_120516.tar.gz](https://research.nhgri.nih.gov/Training/unix/downloads/UnixClass_120516.tar.gz)
2. Double-click on the file to uncompress it. This will create a directory (folder) called *UnixClass*
3. Drag the *UnixClass* directory to the Desktop folder on your Mac.

List directory contents and change working directory	
Run Terminal.app	Find using Spotlight (/Applications/Utilities)
ls	what's here?
pwd	where am I?
cd Desktop	change directory
ls	
cd UnixClass	tab to complete names
ls	
pwd	
cd ..	change to parent directory
pwd	
ls /usr/bin	many of the utilities we're using live here
pwd	we haven't moved
ls UnixClass	list files in directory
ls -l UnixClass	up arrow to modify last command
cd !\$	useful shortcut: last word of previous cmd
ls -ltr	list files in order from oldest to newest; note we can often combine single-letter options
man ls	manual (help) pages for most Unix commands

Note that most punctuation characters have a specific meaning and purpose in the Unix shell. For this reason, it is generally a good idea to avoid using punctuation, including spaces, in filenames. The "safe" characters are letters, numbers, and \_ (underscore), - (dash) and . (period). Anything else will likely create a conflict with one program or another.

Also note that Unix doesn't care about the last few characters of your file name. Whether you name a file "mydata.txt" or "mydata.csv" has no effect on how the programs work. You must tell programs to do the right thing with them, by specifying the column separator, etc. In general, Unix programs are set up to work well with tab-delimited data, as we will use in the examples in this class.

## Unix 101 Part 2 - viewing, counting and searching text files

Concepts: Redirecting output, piping the output of one command to the input of another.

Tasks: See what's inside files, even very large ones, and start to be able to summarize and search through text files.

<b>View file contents</b>	
<code>cat hgmd_brca1.txt</code>	dump file contents to screen
<code>less hgmd_brca1.txt</code>	show one screen at a time
<code>&lt;down&gt;</code> or <code>&lt;enter&gt;</code>	line forward
<code>&lt;up&gt;</code>	line back
<code>f</code> or <code>&lt;space&gt;</code>	page forward
<code>b</code>	page back
<code>g</code>	go to start
<code>G</code>	go to end
<code>q</code>	quit
<code>less -S hgmd_brca1.txt</code>	show without wrapping text
<code>&lt;right&gt;</code>	page right
<code>&lt;left&gt;</code>	page left
<code>q</code>	Note content: some Ovarian cancer, some Breast cancer, some both
<b>Interrupt a running program</b>	
<code>cat</code>	no file; sits and waits
<code>ctrl-C</code>	to interrupt (cancel) running command
<b>View first or last lines of file</b>	
<code>head hgmd_brca1.txt</code>	first 10 lines
<code>head -3 hgmd_brca1.txt</code>	first 3 lines
<code>tail hgmd_brca1.txt</code>	last 10 lines
<code>tail -1 hgmd_brca1.txt</code>	last line
<b>Count characters, words and lines</b>	
<code>wc hgmd_brca1.txt</code>	line, word, character count
<code>wc -l hgmd_brca1.txt</code>	just count lines
<b>Search for patterns, redirect output to a file</b>	
<code>grep Ovar hgmd_brca1.txt</code>	gets some of them (how many?)
<code>grep Ovar hgmd_brca1.txt &gt; ovarian.txt</code>	redirect output into file
<code>wc -l ovarian.txt</code>	seems too small.
<code>less -S ovarian.txt</code>	What's wrong?
<code>grep ovar hgmd_brca1.txt &gt; ovarian.txt</code>	gets the others (grep is case-sensitive)
<code>wc -l ovarian.txt</code>	note that we overwrote file, no warning!
<code>grep -i ovar hgmd_brca1.txt &gt; ovarian.txt</code>	gets all (-i option makes case-insensitive)
<code>wc -l ovarian.txt</code>	

<code>grep -v -i ovar hgmd_brca1.txt &gt; no_ovarian.txt</code>	gets inverse: lines not containing ovar
<code>wc -l no_ovarian.txt</code>	we're creating a lot of files...
<b>Pipe output from one program to input of another</b>	
<code>grep -i ovar hgmd_brca1.txt   wc -l</code>	pipeline sends output of first command as input to second
<code>grep "Hum Mut" hgmd_brca1.txt   wc</code>	quotes needed when search with spaces or other punctuation used by the Unix shell
<code>grep X hgmd_brca1.txt   less -S</code>	pipe to less lets us see what we're getting
<code>grep X hgmd_brca1.txt   grep -i ovar   wc -l</code>	pipe grep to grep to create "and" queries; can have many pipes

## Hands-On Exercises #1

Note: Answer exercises 1.1 and 1.2 using the file `hgmd_brca1.txt`.

- 1.1. Use head and tail to extract lines 2-4 from `hgmd_brca1.txt`
- 1.2. How many BRCA1 variants were published in 2003? How many of these involve both ovarian and breast cancer?
- 1.3. Use man to learn what the `-w` option does in the `grep` command. Try it.

## Unix 101 Part 3 - Wildcards in file names and in grep

Concepts: Wildcards let you apply the same command to multiple files at the same time. Wildcards in `grep` are different and allow you to find matches to patterns, not just specific strings.

<b>Wildcards for file names</b>	
<code>ls *.txt</code>	show all files that end with '.txt'
<code>ls r*</code>	show all files that begin with 'r'
<code>wc -l *.txt</code>	count lines in all files that end with '.txt'
<b>Wildcards in grep patterns</b>	
<code>grep W...X hgmd_brca1.txt   less -S</code>	in <code>grep</code> , '.' is wildcard for single character
<code>grep "W.*X" hgmd_brca1.txt   less -S</code>	use '*' to indicate "zero or more". Need quotes to prevent * from being interpreted by Unix shell.
<code>grep "W[0-9]*X" hgmd_brca1.txt   less -S</code>	square brackets define a "character class", a set of characters that will match a given position. This is more precise than previous command.

## Unix 101 Part 4 - sorting and more counting

Concepts: One can sort alphabetically or numerically. Sorting not only puts things in the desired order; it can also identify unique combinations of values.

<b>Sort alphabetically</b>		
<code>less wgEncodeRegTfbs.txt</code>		new BED file from UCSC; transcription factor binding sites
<code>sort wgEncodeRegTfbs.txt   less -S</code>		
<code>wc -l wgEncodeRegTfbs.txt</code>		it's a big file
<code>head -1000 wgEncodeRegTfbs.txt &gt; tfbs.txt</code>		let's take a smaller sample for today
<code>sort tfbs.txt   less</code>		now it doesn't take too long to sort
<code>sort tfbs.txt   head</code>		use 'head' instead of 'less' to keep a view of recent output in the terminal window
<code>sort -k5,5 tfbs.txt   head</code>		can sort on the score column; but this is still an alphabetical sort
<b>Sort numerically</b>		
<code>sort -n -k5,5 tfbs.txt   head</code>		-n sorts numerically
<code>sort -r -n -k5,5 tfbs.txt   head</code>		and -r reverses the order of the sorting
<code>sort -k5,5nr tfbs.txt   head</code>		another way to say that
<b>Sort on multiple columns</b>		
<code>sort -k1,1 -k2,2n tfbs.txt   head</code>		can sort on multiple columns; -n and -r can be applied to individual columns
<b>Count unique (distinct) values</b>		
		how many different TFs are in there? how many sites for each?
<code>cut -f4 tfbs.txt   less</code>		extract single column
<code>cut -f4 tfbs.txt   sort -u   less</code>		sort -u gives only unique (distinct) values
<code>cut -f4 tfbs.txt   sort -u   wc -l</code>		count them
<code>cut -f4 tfbs.txt   sort   uniq   wc -l</code>		this is just like sort -u
<code>cut -f4 tfbs.txt   uniq   wc -l</code>		<b>don't do this:</b> sort is needed before uniq
<b>Count number of repetitions of each value</b>		
<code>cut -f4 tfbs.txt   sort   uniq -c   less</code>		uniq -c counts the number of sequential repetitions
<code>cut -f4 tfbs.txt   sort   uniq -c   sort -nr   head</code>		show which TFs have the most binding sites
<b>Count unique combinations of values, make subtotals</b>		
		which TFs have sites on the fewest different chromosomes?
<code>cut -f4,1 tfbs.txt   less</code>		can extract multiple columns, but can't specify order
<code>cut -f1,4 tfbs.txt   head</code>		keep the first 10 lines of output for reference
<code>cut -f1,4 tfbs.txt   sort -u   head</code>		First, get distinct chrom/TF pairs...

<code>cut -f1,4 tfbs.txt   sort -u   cut -f2   head</code>	...extract column that we want to "subtotal" on
<code>cut -f1,4 tfbs.txt   sort -u   cut -f2   sort   uniq -c   head</code>	...then re-sort to lump TFs together and count number of chroms for each TF
<code>cut -f1,4 tfbs.txt   sort -u   cut -f2   sort   uniq -c   sort -n   head</code>	sort again to find TFs with fewest chroms

## Unix 101 Part 5 - Working with compressed (.gz) files

Concepts: compression is very useful for large files, and many files that you download from the internet are compressed, to save download time as well as space. You can uncompress the files, or use pipes to work with them while they are compressed.

<code>ls -l refFlat.txt</code>	check out file sizes
<code>gzip refFlat.txt</code>	compress file; replaces original file
<code>ls -l refFlat.txt.gz</code>	this file is compressed
<code>less refFlat.txt.gz</code>	can't work with compressed data directly
<code>gunzip refFlat.txt.gz</code>	uncompress file
<code>ls -l</code>	compressed version is gone
<code>gzip refFlat.txt</code>	compress file again
<code>gunzip -c refFlat.txt.gz   less</code>	can uncompress to pipe, saving disk space
<code>zless refFlat.txt.gz</code>	shorthand for the same thing

## Unix 101 Part 6 – Working with text files

We caution against using Word or Excel to edit text files that you later want to manipulate on the command line using the Unix commands that you've learned in this class. While both programs provide options to save files in various text formats, the resulting .txt or .csv files are not always compatible with Unix. A better option is to use a text editor like BBedit (Mac) to edit any files that will be also accessed under Unix. A list of recommended free text editors is provided at the end of this document.

## Hands-On Exercises #2

Note: Answer the following exercises using the file `tfbs.txt`.

- 2.1. Are all TF binding sites on the "+" (forward) strand?
- 2.2. What are the highest scoring binding sites for c-Myc? the lowest scoring?
- 2.3. How many binding sites for NFkB have the highest score? the second highest?
- 2.4. Going back to `hgmd_brca1.txt`, count how many variants for each type of cancer.

## Unix 101 Extra Topics

These are some basic concepts that didn't fit in the three-hour session. But you can explore on your own, now that you are more comfortable working at the command line.

### Biowulf HPC Cluster at NIH

Biowulf support

<https://hpc.nih.gov/>

Everything you need to know about the NIH Biowulf HPC Cluster, including how to get an account, user guides, training, list of installed applications, etc.

Online class: Introduction to Biowulf

[https://hpc.nih.gov/training/intro\\_biowulf/](https://hpc.nih.gov/training/intro_biowulf/)

The 'Introduction to Biowulf' class is an online, self-paced class. New Biowulf users are encouraged to work through the entire class, and experienced Biowulf users can view specific videos to brush up on a particular section.

### Copy, move, rename and delete files and directories

Tasks: list, copy, move and rename files and directories, just like you can do with the Finder.

So why bother? These work on any Unix machine, including remote servers and clusters.

<code>cp hgmd_brca1.txt my.txt</code>	copy file
<code>ls -l</code>	
<code>mv my.txt brca1.txt</code>	rename (move) file
<code>rm brca1.txt</code>	delete (remove) file <b>Note: there is no undo, and no "recycle bin"</b>
<code>mkdir brca1</code>	create directory
<code>ls -l</code>	notice 'd' for directory, permission flags
<code>cp hgmd_brca1.txt brca1</code>	copy file into different dir
<code>ls -l brca1</code>	notice: same file name, new timestamp
<code>rm brca1</code>	not right command for directory
<code>rmdir brca1</code>	can't do this; dir not empty
<code>rm -r brca1</code>	remove recursively (all contents) <b>Warning: this can be quite dangerous; you can permanently delete a large number of files. USE WITH CAUTION.</b>
<code>mkdir brca1</code>	
<code>cp hgmd_brca1.txt brca1/hgmd</code>	copy using different name
<code>mv brca1/hgmd .</code>	move file here (current directory)
<code>rmdir brca1</code>	remove empty dir
<code>rm hgmd</code>	
<code>rm "refFlat 2.txt"</code>	must use quotes if filename includes spaces or other punctuation

## Accessing Remote Computers

Concepts: Unix commands work exactly the same on Linux servers and HPC clusters (like biowulf) as they do on your Mac. To access these computers over the network, we use ssh (secure shell) and scp (secure copy).

Log in to a remote computer	
<code>ssh studXX@biowulf.nih.gov</code>	"studXX" is your username on the remote computer. If username is the same on local and remote, can leave out "username@"
Password:	Enter password (case sensitive)
<code>pwd</code>	All commands that we know work here, e.g. Where are we?
<code>ls -l</code>	What's here?
<code>exit</code>	Log out
Copy files to remote server	
<code>scp UnixClass.tar.gz studXX@biowulf.nih.gov:</code>	Note colon after hostname.
Password:	Asks for password on remote (won't mention this again, happens every time we connect)
<code>ssh studXX@biowulf.nih.gov</code>	
<code>ls</code>	The file is here, in your home directory, because we didn't provide a destination directory name
<code>tar -xzf UnixClass.tar.gz</code>	Untar gzipped archive
<code>ls -l</code>	created a new directory
<code>ls -l UnixClass</code>	Familiar files
<code>exit</code>	
Copy files from remote server	
<code>scp studXX@biowulf.nih.gov:UnixClass/hgmd_brca1.txt .</code>	Note dot to indicate that the destination is the current directory; always have to provide a destination for scp

## Text editors

If you need to edit text files on your desktop machine, we recommend the following:

BBedit (Mac)

<https://www.barebones.com/products/bbedit/>

Notepad (Windows)

Installed on NHGRI Windows machines

Notepad++ (Windows)

<http://notepad-plus-plus.org/>

Visual Studio Code (Mac and Windows)

<https://code.visualstudio.com/>