

# **Introduction to Mac and Linux Command Line for Bioinformatics**

**INBRE/COBRE Joint Bioinformatics core**



**IDeA Network  
of Biomedical Research  
Excellence**



**Pacific Center for Emerging  
Infectious Diseases  
Research**

## Introduction to the Mac and Linux Command Line

INBRE/COBRE Bioinformatics Workshops Series

Mahdi Belcaid<sup>1</sup>

<sup>1</sup>INBRE/COBRE Joint Bioinformatics Core

December 7, 2010

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

### Notes

- 1 Introduction
- 2 The Essentials  
First Steps
- 3 Concepts for Digging Deeper  
Files Summary Commands  
Regular Expressions  
grep, sed and awk: a Brief Overview  
Running Tools from the Blast Suite

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

### Notes

BLAST through command line for 454 data

Control: More ways than through the GUI  
Speed: Certain commands require many GUI step to do equivalent  
Availability of applications: Bioinfo tools and powerful system applications  
Because (some) bioinformaticians are Lazy!

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

### Notes

## Why Mac and Linux?

- Mac is Unix based and Linux is Unix-like
  - They are both "Unix" compatible and possess a comparable set of core system-applications
- A wide range of applications available for both
- Platform of choice for bioinformaticians
  - More and more biologists embracing it (Mac more than Linux)
- There are ways to emulate Linux on Windows
  - Easier to just run Linux from a Virtual Machine

Mahdi Belcaid Introduction to the Mac and Linux Command Line

Notes

---

---

---

---

---

## What You Should Know

- During this workshop, we will interact with the OS only through the Command Line Interface (CLI)
- The CLI is interactive, i.e. takes commands, executes specific tasks and returns (or not) an output
  - Each command accomplishes a specific task
- Case-sensitive: Remember that commands, files, folders, behaviors, . . . , are all case sensitive

Mahdi Belcaid Introduction to the Mac and Linux Command Line

Notes

---

---

---

---

---

## Tips for a "Happy Interaction"

- History: The shell keeps a history of a set of typed commands, uses arrow to access them or to type **history** to see the full list
- Proprietary formats: Avoid working with annotated files. (Ex: Word, Excel, pdf, etc...)
- Avoid Spaces: Rename any files/folder that contains spaces in its name.
- To view the contents of a file, use the **less** program

**less <file\_name>**

- To quit the **less** viewer, type "**q**"
- Aborting commands: You can press **Ctrl+c** to stop a running command

Mahdi Belcaid Introduction to the Mac and Linux Command Line

Notes

tab button fills in command

"history"

## I.E.1 Some Shortcuts

- Command completion: Use tab to complete a partially typed command, a file or a directory name
- **Ctrl+A** Takes you back to the beginning of the line
- **Ctrl+E** Takes you back to the End of the line
- **Ctrl+delete** delete one word back
- **Ctrl+→** jump one word to the right
- **Ctrl+←** jump one word to the left
- To get help, use "-help" option or **man** command
  - type "q" to exit **man** or **less**
- More useful commands on the cheat sheet

### <sup>1</sup>I.E. Interactive Example

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

## Unix Command Line Structure

- Typical structure<sup>2</sup>:

command [options] [arguments]

- Options alter the behavior of the command
  - ...and in most cases, the output
- Arguments are the operands of the command

Example<sup>3</sup>:

`mahdi@laptop$ grep -i ACGT myInutFile.fasta`

- Command line structure across Unix-like systems is similar

<sup>2</sup>Code in red boxes is not executable

<sup>3</sup>Code in green boxes is executable

## Unix-like File System Hierarchy

- One of the major differences between Windows and Linux/OS X
- Unix-like systems follow a "Quasi-standardized" hierarchical structure
- The root of the hierarchy is: /
  - Similar to C: or D: under Windows
- All files, folders, applications, drives, etc..., are mounted (placed by the OS) under /

Windows path    C:\Folder\subfolder\file.txt  
Unix-like path   /Folder/subfolder/file.txt

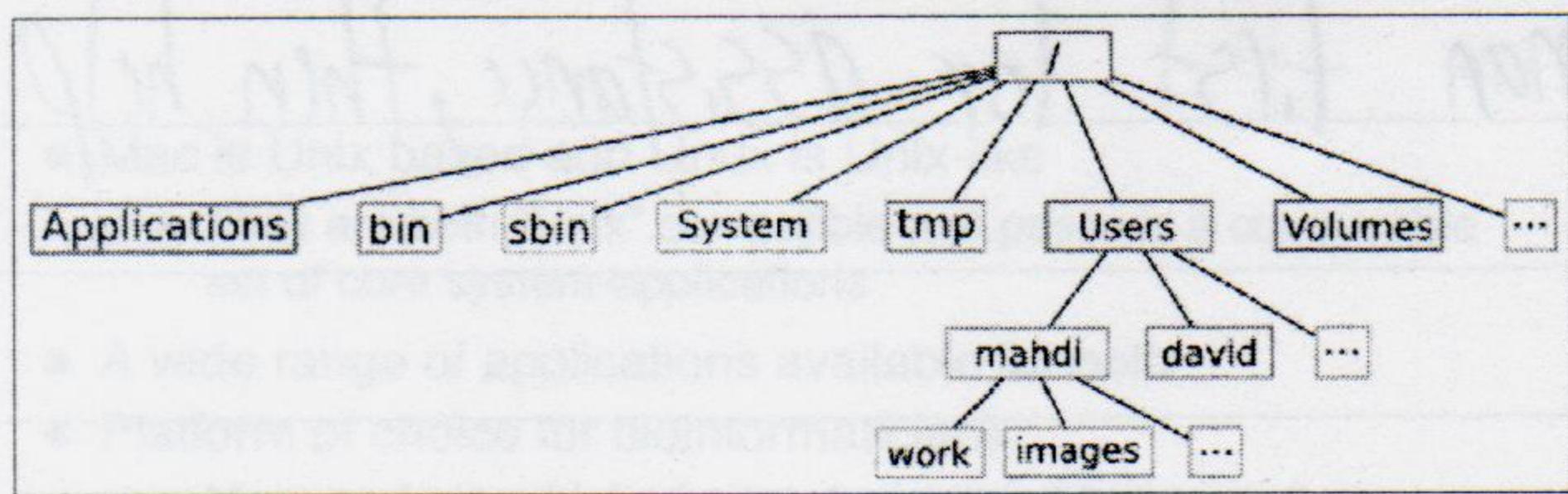
## Notes

man first for assistance, then help

## Notes

## Notes

## OS X File System Structure



/Applications	User apps
/bin	Essential programs (date, ls, ...)
/sbin	Essential system admin utils
/System	Critical OS files (...)
/tmp	Temporary files, caches, ...
/Users	All user homes on the machine
/Volumes	Mounted devices & volumes (hard disks, DVD's, DMG mounts, ...)
...	Other directory partitions (/dev, /etc, /usr, ...)

Mahdi Belcaid Introduction to the Mac and Linux Command Line

Notes

bin = binary      binary = program  
sbin = admin, system binary

tmp: deletes on restart

## Gathering Information

- List the content of a directory:

```
ls [directory_path]
```

```
mahdi@laptop$ ls /home/mahdi
Applications  Dropbox  Music  Photos
work
```

- To know where you are, use the **pwd** command (print working directory)

```
mahdi@laptop$ pwd
/Users/mahdi/
```

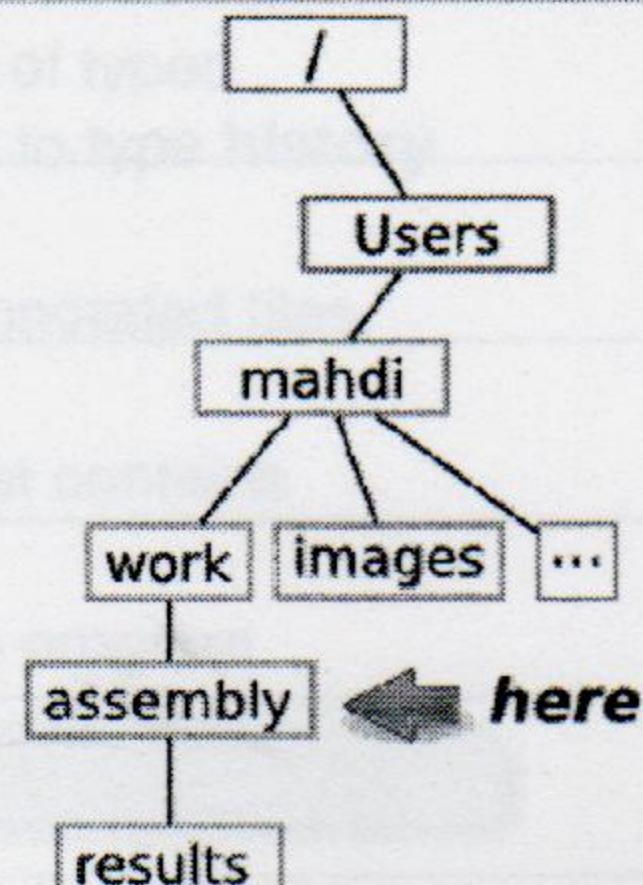
Mahdi Belcaid Introduction to the Mac and Linux Command Line

Notes

## Navigating the Hierarchy

- To move from one place to another place on the file system, type arrow in front of them, or move to see the full path
- Open every file with a word processor (e.g. Word, Excel, pdf, etc.)
- Avoid Spaces: Rename any file that has spaces in its name.
- To view the contents of a file, use the **cat** command

- You are in:  
`/Users/mahdi/work/assembly`



Notes

Introduction  
The Essentials  
Concepts for Digging Deeper

First Steps

## Navigating the Hierarchy: `cd`

- Changing directory: Going up and down the hierarchy

```
cd [directory_path]
```

- `directory_path`: The *absolute* or *relative* path

- ① *Absolute path*: Starts at the root and spell full path
- ② *Relative path*: with regards to your current location use .. to specify a parent directory

Notes

---



---



---



---



---

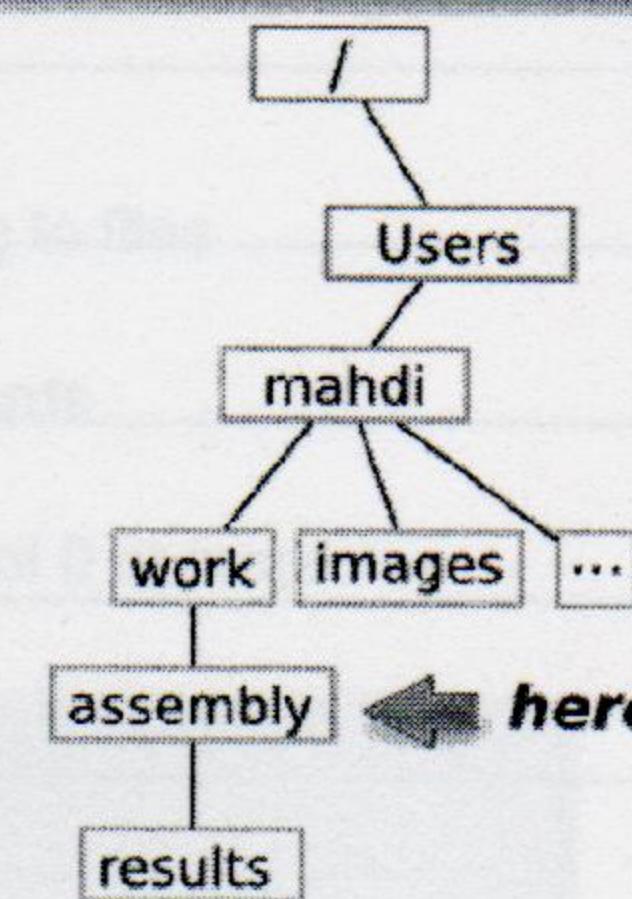
Mahdi Belcaid | Introduction to the Mac and Linux Command Line

Introduction  
The Essentials  
Concepts for Digging Deeper

First Steps

## Navigating the Hierarchy: Going-down

- You are in:  
`/Users/mahdi/work/assembly`



Notes

---



---



---



---



---

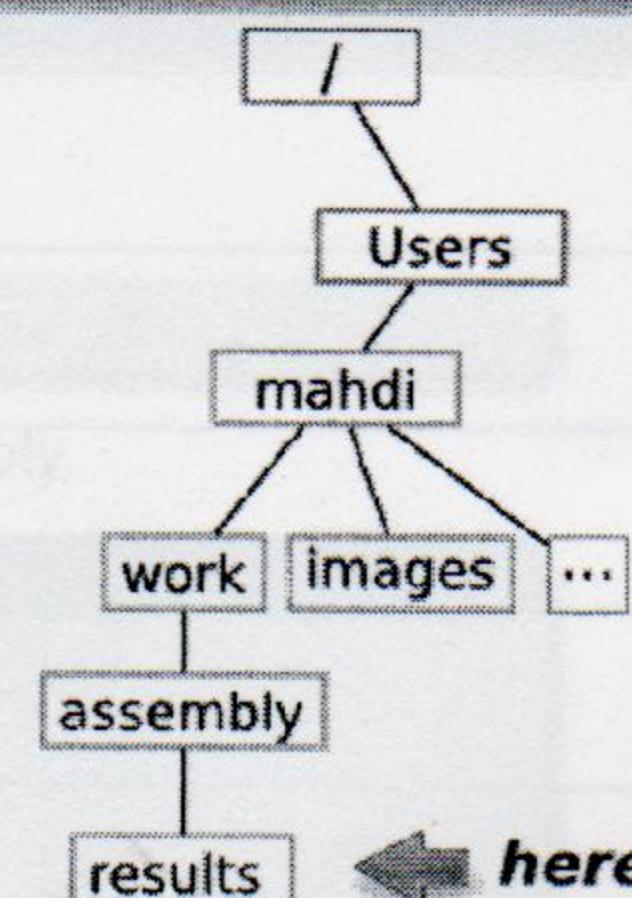
Mahdi Belcaid | Introduction to the Mac and Linux Command Line

Introduction  
The Essentials  
Concepts for Digging Deeper

First Steps

## Navigating the Hierarchy: Going-down

- You are in:  
`/Users/mahdi/work/assembly`
- You need to go into the directory `results`, located in `assembly`



Notes

---



---



---



---



---

Mahdi Belcaid | Introduction to the Mac and Linux Command Line

Notes

```
mahdi@laptop$ pwd  
/Users/mahdi/work/assembly
```

- Using the *absolute* path:

```
....$ cd /Users/mahdi/work/assembly/results
```

- Using the *relative* path:

```
mahdi@laptop$ cd results
```



Notes

• List the contents of a directory

- You are in:

```
/Users/mahdi/work/assembly/results
```



Notes

- You are in:

```
/Users/mahdi/work/assembly/results
```

- To go into the directory

```
/Users/mahdi (user's home  
directory)
```



Introduction  
The Essentials  
Concepts for Digging Deeper

First Steps

## Going-Up: Example

```
mahdi@laptop$ pwd
/Users/mahdi/work/assembly/results
```

- Using the *absolute path*:

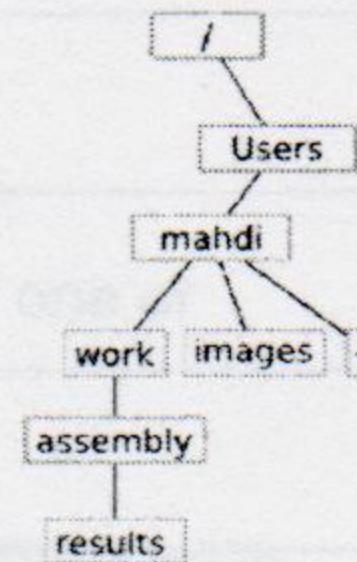
```
mahdi@laptop$ cd /Users/mahdi
```

- Using *relative path*:

```
mahdi@laptop$ cd ../../..
```

Or use a shortcut

```
mahdi@laptop$ cd
```



Notes

---



---



---



---



---

Introduction  
The Essentials  
Concepts for Digging Deeper

First Steps

## CLI wildcards: \*

- They are used by the CLI when referring to files
  - These are not regular expressions!
- There are various command line wild-cards
  - We will only look at: \*
- \* is interpreted by the CLI as any string of 0 or more characters.

```
file* Any file name that starts with the prefix "file".
Ex: file, file128, file_A, ...
*.txt Any file that ends with the extension ".txt"
```

wildcards

Notes

---



---



---



---



---

Mahdi Belcaid | Introduction to the Mac and Linux Command Line

file names started with . are invisible

Introduction  
The Essentials  
Concepts for Digging Deeper

First Steps

## Creating or Deleting Directories mkdir and rmdir

- You can create or delete directory with:

```
mkdir dir_name rmdir dir_name
```

- For **rmdir**, the directory needs to be empty.

### Example

```
1 mahdi@laptop$ mkdir /Users/mahdi/work/assembly/results/tempDir1
2 mahdi@laptop$ mkdir /Users/mahdi/work/assembly/results/tempDir1/tempDir2
3 mahdi@laptop$ rmdir /Users/mahdi/work/assembly/results/tempDir1
4 rmdir: failed to remove '/Users/mahdi/work/assembly/results/tempDir1': Directory not empty
```

Notes

---



---



---



---



---

Mahdi Belcaid | Introduction to the Mac and Linux Command Line

rmdir → file needs to be empty

## Copying Files: cp

- Done through the **cp** command:

```
cp file1 file2
```

```
cp file1 file2 file3 ...filen dir
```

```
cp -r dir1 dir2
```

- Paths of arguments can be relative or absolute

-r → recursive, needed for copying directory w/ contents

Mahdi Belcaid Introduction to the Mac and Linux Command Line

## Deleting Files

- Done through the **rm** command:

```
rm file1 file2 file3 ...filen
```

- You can remove directories with **rm**
  - Consult the **man** page of **rm**

```
mahdi@laptop$ rm -r dirPath
```

- Be extra careful with **rm**... Especially when using it in **-r** (recursive) mode or with the "\*" wild-card
  - No way to undo **rm**

Mahdi Belcaid Introduction to the Mac and Linux Command Line

## Exercise 1

- Download the following zip archive:  
<http://dl.dropbox.com/u/28270/work.zip>
- Move the **work.zip** file to your home directory
  - Example in my case it is: /Users/mahdi/
- unzip the **work.zip** archive

```
mahdi@laptop$ unzip work.zip
```

- This will create a new directory **work** in your home
- Create a directory in **work** called **exercise\_1**
- Copy all the fasta files (.fasta) in **work/temp** to **exercise\_1**
- Delete the directory **temp**

## Notes

## Notes

## Notes

## Counting File Characteristics: **wc**

- Generates the word, line and character counts for one or more inputs
- Three useful options:

```
wc -l [file1, ..., filen] #counts the number of lines  
wc -m [file1, ..., filen] #counts the number of characters  
wc -w [file1, ..., filen] #counts the number of words
```

Notes

## Sorting a File: **sort**

- The command **sort** sorts the content of a file passed as a parameter
- Default behavior is to sort alphabetically, in ascending order
- Most used options are:

```
-n Sort using numerical order  
-r Sort in reverse(descending) order
```

Notes

## Removing Redundant Lines: **uniq**

- Deletes the duplicated lines in a file
- File needs to be sorted
- Most used options are

```
-c prefix the unique lines by the number of occurrences  
-d print only the repeated lines
```

commonly used

Notes

Uniq tells you one of each thing you have

## Output Redirection: >

- The ">" is used to redirect **standard output** (what gets printed to the screen)
- You can send the output to another program (**pipe**), or **redirect** it to a file
- Redirection is selected through the ">" character

```
program [option] [arguments] > results_file.txt
```

```
mahdi@laptop$ ls /Users/mahdi/work/viruses > listing.txt
```

- listing.txt contains all the names of files and directories located in /Users/mahdi/work/viruses

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

Notes

To one test command to the terminal and know what happened

Terminal output goes here

Copy and paste command to terminal

Ctrl + C to stop command

Ctrl + V to paste command

Ctrl + S to save command

Ctrl + D to exit command

## Unix Pipes: |

- A pipe is a way of chaining two or more commands in Unix
  - The standard output of the first command becomes the input of the next command
  - The output is put in a "virtual" files for the subsequent program to use
- Example: Since **uniq** requires sorted files, we can easily pipe the output of **sort** into **uniq**

```
mahdi@laptop$ sort genes.ids | uniq
```

- There is no limit on the number of pipes you can do

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

Notes

P chain pipe used when you don't need to save output

## I.E. Using Pipes

- Counting the number of unique blast hits:

/Users/mahdi/work/blasts/contigs\_rePLICATE\_all.csv  
contains the blast results of a contigs file against nt

/Users/mahdi/work/blasts/genes.ids  
contains all the hits from the blasts files (8th column)

```
mahdi@laptop$ sort genes.ids | uniq | wc -l
```

This counts the number of unique ids in the file genes.ids

Notes

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

## Exercise 2

- How many fasta files are there in the exercise\_1 directory you created earlier?
- Print a sorted listing of the assembly/libraries directory into a file named "temp"
  - Make sure the file is sorted alphabetically
- How many distinct libraries are in assembly/library/data

ls | wc -l

### Notes

sort ./\*libraries/\* > temp.txt

## Inclusion, Exclusion Comparison: comm

- comm is a way of comparing 2 files
  - Same logic as a Venn diagram
- The syntax of the program is:

```
comm [-1, -2, -3] file_1 file_2
-1 Supresses lines unique to file_1
-2 Supresses lines unique to file_2
-3 Supresse lines common to both files
```

### Notes

## I.E. Comparing two files

- Directory /Users/mahdi/work/examples/comm contains two files
  - file\_1 and file\_2

```
mahdi@laptop$ comm -13 file_1 file_2
```

Prints lines unique to file 2

```
mahdi@laptop$ comm -3 file_1 file_2
```

Prints lines that are not shared by file\_1 and file\_2

### Notes

<p>Introduction</p> <p>The Essentials</p> <p>Concepts for Digging Deeper</p>	<p>Files Summary Commands</p> <p>Regular Expressions</p> <p>grep, sed and awk: a Brief Overview</p> <p>Running Tools from the Blast Suite</p>
--	---

## Exercise 3

- Compare the singletons rate between two alternative assemblies  
assembly1\_singletons.ids and assembly2\_singletons.ids in the results folder
    - How many ids are singletons in assembly1\_singletons.ids?
    - How many ids are considered singletons in both assemblies?

## Notes

# What is a Regular Expression?

- "A regular expression, abbreviated as regex or regexp, is a string that describes or matches a set of strings, according to certain syntax rules." *Wikipedia*
  - A regex is a pattern used to give a concise description of a set, without listing all its elements
  - An exact string or one literal character is a trivial example of a regular expression
    - Describes exactly one element
    - Example: A, or ACGT, etc...

## Notes

Mahdi Belcaid | Introduction to the Mac and Linux Command Line

# How Does a Regular Expression Work?

- Uses meta-characters, meta-symbols to better describe a pattern
  - Each meta-(character|symbol) has a special meaning
  - Meaning is expanded during search phase to represent a class
  - Suppose we need to search for exact occurrences of poly-tails in a file
    - We start with A

## Notes

Mahdi Belcaid | Introduction to the Mac and Linux Command Line

## Using Anchors

- To enforce matching to either the beginning or the end of a string, anchors are used:

^ Forces matching at beginning of line  
\$ Forces matching at end of line

^A Searches for all the lines that start with A  
A\$ Searches for all the lines that end with A

- However, we know that a poly-A's contain more than one A.
  - How do we define the multiplicity of a character?

## Using Quantifiers

- The minimum, as well as the maximum number of matches can be set using quantifiers
- Quantifiers apply to the object they follow
- Quantifiers can be explicit:

A{5} find exact occurrences of AAAAA (5 As)  
A{5, 10} Minimum 5 consecutive As, maximum 10  
A{5, } 5 consecutive As with no maximum

- Or they can be implicit (using a meta-character)

A\* match A, 0 or more times  
A+ match A, 1 or more times. Example: A, AA, AAA  
equivalent to A{1,}  
A? match A, 0 or one time

## Using Alternation

- We know that poly-tails can be comprised of Ts as well
  - How can we define alternative patterns?
- To describe alternative options to match, use the | meta-character

A|T match an A or a T  
A{5,}|T{5,} match 5 consecutive As or Ts

- The regular expression engine tries to match any occurrence of at least 5 consecutive As with no maximum
- If it fails, it tries to match any occurrence of at least 5 consecutive As with no maximum

or

options

## Notes

## Notes

## Notes

## Other Useful Meta-Characters: ":"

- The . represents any character
  - Depending on the OS, the context, the applications, this previous statement might be false
  - It is safe to assume that any alphanumeric character + space can be represented by .
- Considering the DNA alphabet

A.T will match AAT, ACT, AGT, ATT, or A T

## Notes

to line edit in principle edit mode of Drush... editing of a file

to save file in current directory a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

## Other Useful Meta-Characters:

Character Classes

- [] are used to represent a class of characters
- Match one character of the class

[ACGT] represent either A or C or G or T

[A-M] represent either, A, B, ..., through M

[0-9] represent any one digit character

- ." and "[]" can be combined with any of the previous quantifiers

[0-9]+ a number of one or more digits

## Notes

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

## Quick Review

- To describe a pattern, you can use a combination of the meta-characters|symbols)
- You need at least one literal, or a character class, such as ." or "[0-9]"
- Use quantification to define multiplicity of the occurrence
- Separate patterns with alternation symbol "|"

### Examples

AGT*	Matches AG, AGT, AGTT, AGTTT ...
N{4,}	Matches any occurrence of four or more consecutive Ns
^>.+	Matches id lines from a fasta files

## Notes

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

to use both command a file is now saved in current directory

0 of  
more times

Examples

AGT*	Matches AG, AGT, AGTT, AGTTT ...
N{4,}	Matches any occurrence of four or more consecutive Ns
^>.+	Matches id lines from a fasta files

FASTA start with these

Starts with > with one or more characters

## grep: An Intro

- Very powerful utility program used to search input for lines matching a given pattern
- The syntax of grep is the following:

```
grep [OPTIONS] PATTERN [File_1, ..., File_n]
```

- To find all the lines containing the pattern ACGT in the file sequences.fasta

```
mahdi@laptop$ grep ACGT sequences.fasta
ATGCTAGTGATCGTAGCTGTGAACGT
ACTGTGTAGACGTATGTGTAGTCGTGAT
```

- Patterns in grep are case sensitive

Mahdi Belcaid | Introduction to the Mac and Linux Command Line

## Notes

grep - kind of like find

## More with grep

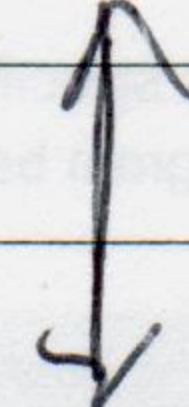
- The pattern used in grep can be a regular expression
- Use the -E option for more flexibility

```
mahdi@laptop$ grep -E "^.+." input.fasta
>tempFT13I length=103 xy=2275_2189 region=3 run=R_2010_09_22_21_17_58
>tempHH7RR length=74 xy=3960_3134 region=3 run=R_2010_09_22_21_17_58
>tempGS0CL length=178 xy=2856_0539 region=3 run=R_2010_09_22_21_17_58
>tempPS18I length=93 xy=2409_1839 region=3 run=R_2010_09_22_21_17_58
```

Mahdi Belcaid | Introduction to the Mac and Linux Command Line

## Notes

grep > kbg.fasta | wc -l



grep -c > tbg.fasta



Often used to count # of sequences

## Useful grep Options

- Some of the most used options

```
-c    count the number of occurrences of the pattern
-i    ignore case of the pattern
-n    print the line number containing hit
-v    print lines that DO NOT contain pattern
-A N   print N lines after the pattern
-B N   print N lines before the pattern
```

- Many other options...

Mahdi Belcaid | Introduction to the Mac and Linux Command Line



## sed: An Introduction

find & replace

- **sed** (Stream editor) is used to perform a set of editing operations on a file
  - Reads the input, one line at a time, and executes the edit operation specified
  - We will only be looking at **sed** for making substitutions
  - Substitute a pattern occurrence with a string
- By default the modified file is printed to the standard output (Screen)

## Notes

## Mahdi Belcaid Introduction to the Mac and Linux Command Line

## sed Basic Format

- The basic command line structure is:

```
sed [options] 's/PATTERN/STRING/[g]' file
```

### Example

```
mahdi@laptop$ sed 's/ACGGCA/TGCCGT/' sequences.fasta
```

- ① Replaces the first occurrence of 'ACGGCA' (if any) on each line, with the string 'TGCCGT'
- ② Outputs the contents of sequences.fasta to stdout

## Notes

## Mahdi Belcaid Introduction to the Mac and Linux Command Line

## Useful sed Options

- By default, only the first string is replaced
  - You can change the behavior by running **sed** in greedy mode
- Also, you can instruct **sed** to make the changes directly to the input file

```
mahdi@laptop$ sed 's/ACGGCA/TGCCGT/g' sequences.fasta
```

- **sed** makes the modification directly to sequences.fasta, but...
  - Creates a backup of the original file with the ".backup" extension

## Notes

↗ with no g, it will replace only the first instance per line

s/abc//g ← replaces abc w/nothing

## Regular Expressions With sed

- **sed**'s strength comes from the fact that you can use it with regular expressions
- To instruct **sed** to expect an extended regular expression, the "-E" option is used
  - or "-r" under Linux

```
...$ sed 's/\.[0-9]+//g' accession.fasta >  
accession_no_version.fasta
```

- Remove the trailing version number for a NCBI accession.
  - Replaced with empty string "//"
  - We need to escape wild-cards chars to give them their literal meaning

Used to  
get duplicate accession numbers to have same tag

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

## Notes

## Exercise 5

- In the file:  
`/Users/mahdi/work/blasts/genes.ids`
- Remove the dd identifier in the beginning of the name.  
(gb1, db1, ...)
- Remove the version from NCBI accessions names
- Remove spaces at the beginning or end of line

ref |      sp |      embl

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

## Notes

this character makes  
you take literal meaning  
↓ of bar

Sed -E 's/.+\|//g' genes.ids  
Sed -E 's/\.[0-9]+//g' genes.ids

## Working with awk

- **awk** is an excellent tool for working with column delimited files
  - Ex: Comma delimited, tab delimited, etc...
- Complete, but often complicated, reporting language
- We will only look at some of the useful functionalities
- Basic command-line structure:

```
awk [options] 'program text...program text' input_file
```

## Notes

Introduction  
The Essentials  
Concepts for Digging Deeper

Files Summary Commands  
Regular Expressions  
grep, sed and awk: a Brief Overview  
Running Tools from the Blast Suite

## How awk Addresses Columns

- By default, **awk** splits columns based on `\*` (one or more spaces)
- Behavior can be changed in the program text
- Each column is assigned an "automatic" variable name
  - ① Column 1 is assigned variable name `$1`
  - ② Column 2 is assigned variable name `$2`
  - ③ Column n is assigned variable name `$n`
  - ④ The whole line is assigned variable name `$0`
- In the program text, to print a column, use the following syntax:

```
print column_a"SEPARATOR"column_n,...
```

Mahdi Belcaid | Introduction to the Mac and Linux Command Line

### Notes

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...  
...  
...

...  
...  
...<

## I.E. Tying awk with Other Programs

- Extracting gene information from NCBI genbank file  
`/Users/mahdi/work/proteins/protein.gbk`
- Extract gene name from  
`textt/gene="gene name"`
- Ignore `/gene_synonym`
- Remove formatting (= and ")

~~awk~~ /t = tab as delimiter

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

### Notes

awk to awk / tabbed entries into lines

tab spacing will be kept at new

new entries will be added to lines

the new entries will be added to lines

## Conditionals in awk

- We will make a slight exception to our "No Programming" rule
- We can code conditionals in awk

```
awk '{ if (<condition>) {execute this block}
else{ execute this other block} }' input_file
```

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

### Notes

the new entries will be added to lines

## I.E. Extracting Specific Hits from a Blast Report

- Conditionals are useful if you want to extract a subset to the data
  - Split a file into different files
- Example: to extract all the hits that are above 1e-05

```
awk '{if($3>1e-05){print $1,$8}}'
contigs_viruses_all.csv
```

### Notes

if statement e must have number before it otherwise it thinks that it is a string.

the new entries will be added to lines

Mahdi Belcaid

Introduction to the Mac and Linux Command Line

## Blast Databases

- Blast requires specifically formatted database to run
- **nt, nr, envxx** are available, preformatted, from the NCBI ftp site
- For this example, we will format our own database:
- The basic syntax for formatting a blastdb is:

```
formatdb -i <fasta_input> -p <F|T> -o  
-i <input> Input file to be formatted (Mostly Fasta)  
-p F|T The file type. protein(T) or nucleotide(F)  
-o Created indexes for fastacmd
```

## Notes

## Running Blast from the Command Line

- The basic Blast command line is:

```
blastall -p <program_name> -i <input_file> -d <database>  
-o output_name [other options]
```

- Blast has many params that are not available on public web interfaces
- For more details, see **man** page or type:

```
mahdi@laptop$ blastall -help
```

## Notes

## Some Blast Parameters

- There is no consensus on the most used blast options
  - Depends on the task
- The following params are very useful

```
-m INT changes the output type (-m 8 for table)  
-e eval set the e-value threshold  
(usually 1e-05 for nucleotide blast.)  
-F F|T filter sequence for low complexity  
(T by default)  
-W INT changes the word size, for small input sequence
```

## Notes

## Parsing Sequences from a Blastdb

- **fastacmd** is a program from the blastall suite
- Used to extract sequences from a blast formatted database
  - Extract one or more sequences from the DB
- Syntax of the program is:

```
fastacmd -d <database> -s <ids> [-o output_file]
-s comma delimited list of ids to extract
-i instead of -s, you can provide a file
containing a list of ids
```

### Notes

... run on each line of the input file.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

## I.E. Working with Blast

- Formatting the  
`/Users/mahdi/work/blasts/blastdb/kbay.fas`  
file
- Blasting `test.fas` against `kbay.fas` in tabular format
- Extracting hit ids with `evalues > 1e-50`
- Extracting the ids from the original fasta file using `fastacmd`

### Notes

... run on each line of the input file.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

## Exercise 6

- Use the blast report:  
`/Users/mahdi/work/blasts/contigs_viruses_all.csv`
- Extract all the ids whose hit description contains the keyword `Rhizosolenia`
  - Make sure you have no redundancies
- Extract the corresponding sequences using `fastacmd`
  - The original sequences are in  
`/Users/mahdi/work/blasts/blastdb/kbay.fas`
  - Remember you have to `formatdb` before running `fastacmd`
- Once you have your input, blast it (`blastn`) against the already formatted db:  
`/Users/mahdi/work/blasts/blastdb/marine/marine.fas`
  - Generate a table-formatted output (-m 8)

### Notes

... run on each line of the input file.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

... calculates how big the sum of the seq. elements will be.

... 100K+ entries, some interesting statistics are shown on the screen.

## Feedback Time

- Any thoughts?
- Suggestions for next workshops?
- Subjects suggestions?

### Notes

---

---

---

---

### Notes

---

---

---

---

### Notes

---

---

---

---

# Introduction to Command-Line for Bioinformatics Cheat Sheet

## Regular Expressions

### Anchors

X<sup>^</sup> Match X at the beginning of a line  
X\$ Match X at the end of a regex

### Quantifiers

<b>Explicit</b>	Extracts fasta sequences from a blast-formatted DB
X{N}	Match X exactly N times
X{N,M}	Match X at least N times
No more than M times	
X{N, }	Match X at least N times

### Implicit

X*	Match X, 0 or more times
X+	Match X, 1 or more times
X++	Match X, 0 or 1 times

### Alternation

regex_A   regex_B	Match regex_A or regex_B
Character Classes [XYZ]	Match any <u>one</u> letter between [] . (X or Y or Z)
[A – M]	Match any <u>one</u> letter in range A to M
[0 – 9]	Match any digit
.	Match any character

Usually, ^ and \$ are respectively first and last elements of a

To get help, try man command  
If manual pages are not available for command, try:  
command --help or command -h or

### Miscellaneous symbols

> Redirection symbol

| Piping symbol

## Elementary Commands

cd [dir-path] Change directory  
mkdir name Create new directory  
rmdir <dir> Remove empty directory  
cp file1 file2 Create a copy of file 1  
cp file1 file2 Copy files 1 through n to dir  
..file1 dir Copy directory1 to directory2  
cp -r dir1 dir2 Copy files 1 through n  
rm file1 Recursively remove directory  
Notes

- .. / references one level up the hierarchy
- Paths can be *relative* or *absolute*.
  - Absolute paths starts at the root /
  - Relative paths start at the *current working directory*

• rmdir requires dir to be empty

## File Summary Commands

Count Characteristics of a file

wc -l Count lines

-m Count character

Count words

-w Sort a file

sort -n Use numerical order

-r Use descending order

## Blastall Suite

formatdb	Formats a database for blastall or fastacmd
formatdb -i <input file> -p <type> -o T	for protein, F for nucleotides

fastacmd	Extracts fasta sequences from a blast-formatted DB
fastacmd -d <database> -s <ids> [-o output_file]	Comma delimited list of ids to extract
-s	Instead of -s, you can provide a file containing a list of ids
-i	

### blastall

blastall -p <program> -i <input> -d <database>	an input file on a database
-o output [other options]	
Program name (blastn, blastx, tblastx, ...)	
Program name (Fasta input)	
blast formatted DB	
The output file	
Changes the output type. Use -m 8 for table output	
Set the e-value threshold (usually 1e-05 for nucleotide blast.)	
(F T) Filter sequences of low complexity (T by default)	
Changes the word size, for small input sequence	

## Blast Notes

grep [OPTIONS] PATTERN [File\_1, ..., File\_n]

-c Count the number of occurrences of the pattern

-i	Ignore case of the pattern
-n	Print the line number containing hit
-v	Print lines that DO NOT contain pattern
-AN	Print N lines after the pattern
-BN	Print N lines before the pattern

sed	Perform editing operations on a file
sed [options] s/PATTERN/STRING/[g] file	Give a regex as a pattern (-r under linux)

-E	Backup a file and edit it directly
-i postfix	

awk	Parse and process pattern delimited files
awk [options] 'program text' input_file	Contains the logic for processing your input file
program text	

\$	Ex. column_a"SEPARATOR" column_n,...
\$1,\$2,\$3,...	Prefix column of the input
awk '{ if (<condition>) {execute this block} else{execute this other block} } <condition>	Column 1, 2, 3, etc ...

Implements a conditional test. The format is:	
(COL >   <   =   != Value)	Another print statement

block to execute	
------------------	--

sort	
-n	
-r	

history

pwd - print working directory, shows what directory you are at

cd .. goes up a directory