

Aug 08, 2016 Version 3

Unix and Bioinformatics V.3

DOI

dx.doi.org/10.17504/protocols.io.fitbken

Benjamin Tully and Ken Youens-Clark¹

¹EARTHCUBE OCEANOGRAPHY AND GEOBIOLOGY ENVIRONMENTAL 'OMICS

ECOGEO



Elisha M Wood-Charlson

KBase

OPEN  ACCESS



DOI: dx.doi.org/10.17504/protocols.io.fitbken

Protocol Citation: Benjamin Tully and Ken Youens-Clark 2016. Unix and Bioinformatics. [protocols.io](#)

<https://dx.doi.org/10.17504/protocols.io.fitbken>

License: This is an open access protocol distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited

Protocol status: Working

Created: August 08, 2016

Last Modified: March 28, 2018

Protocol Integer ID: 3379

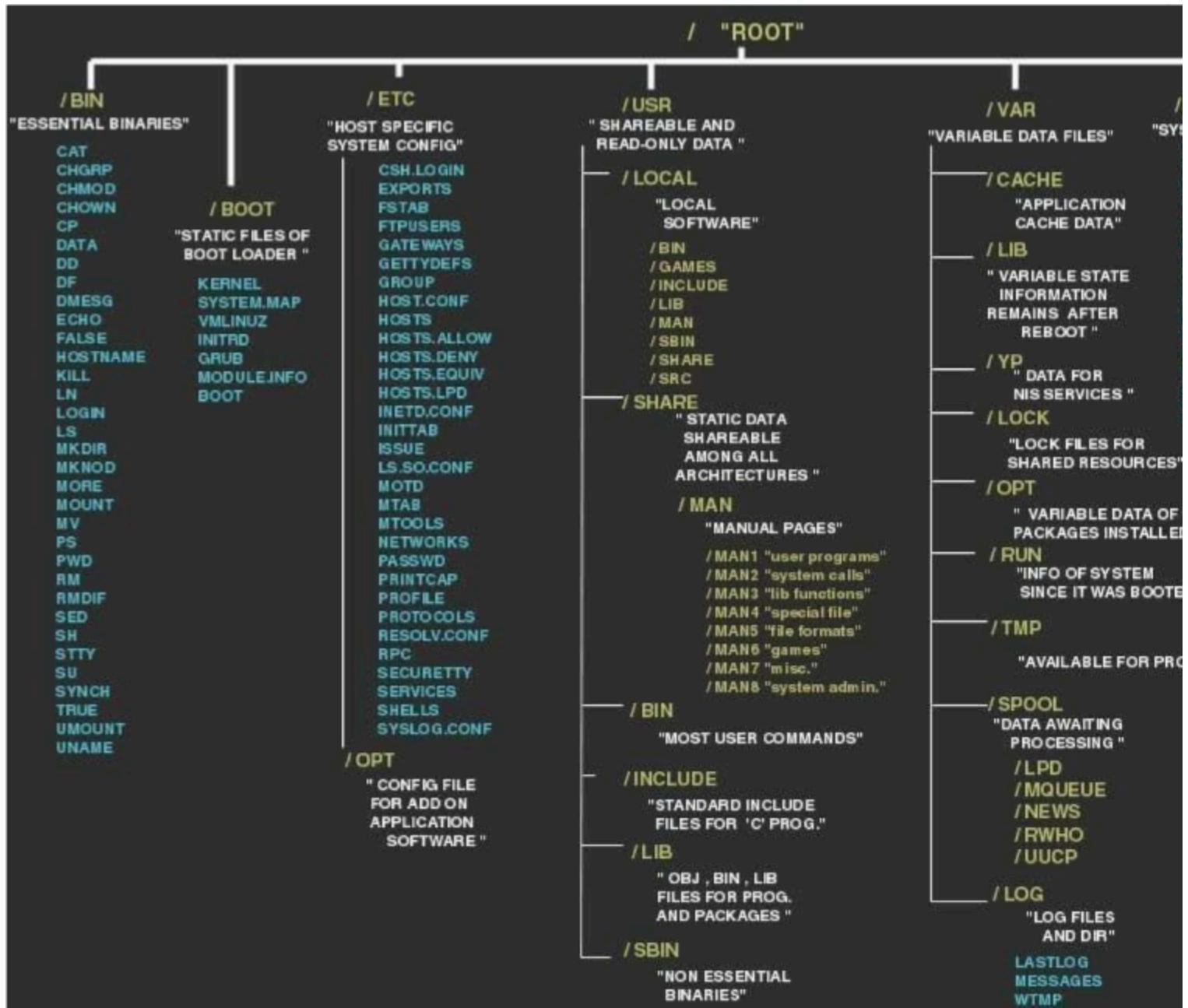
Abstract

This protocol details the use of various unix commands commonly used in bioinformatics.

Guidelines

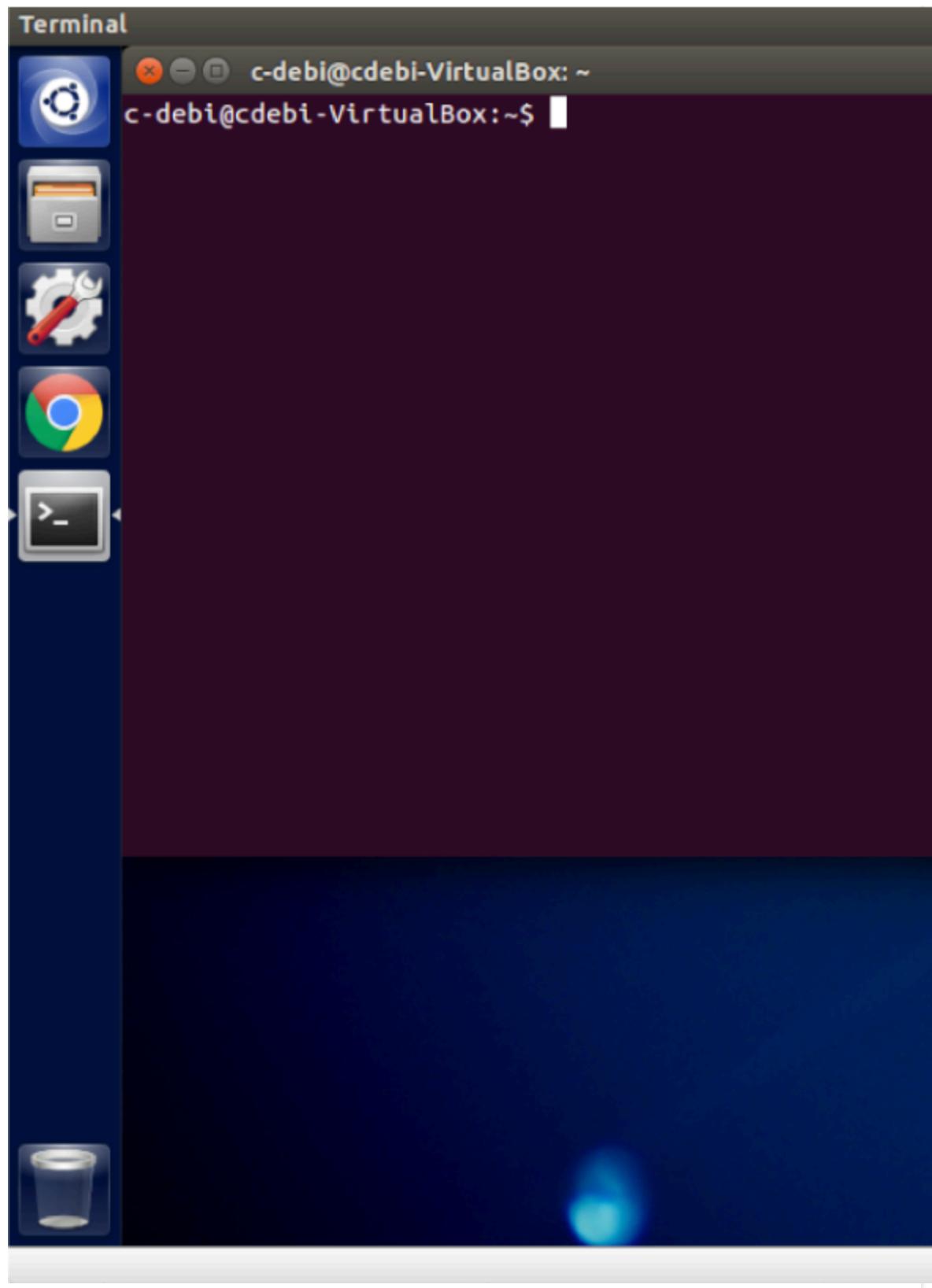
Unix Commands

	pwd	rm	grep	tail	install
ls	'>'	sed	cut		
cd	cat	nano	top		
mkdir	'<'	history	screen		
touch	' '	\$PATH	ssh		
cp	sort	less	df		
mv	uniq	head	rsync/scp		



The Start

- 1 Open terminal window



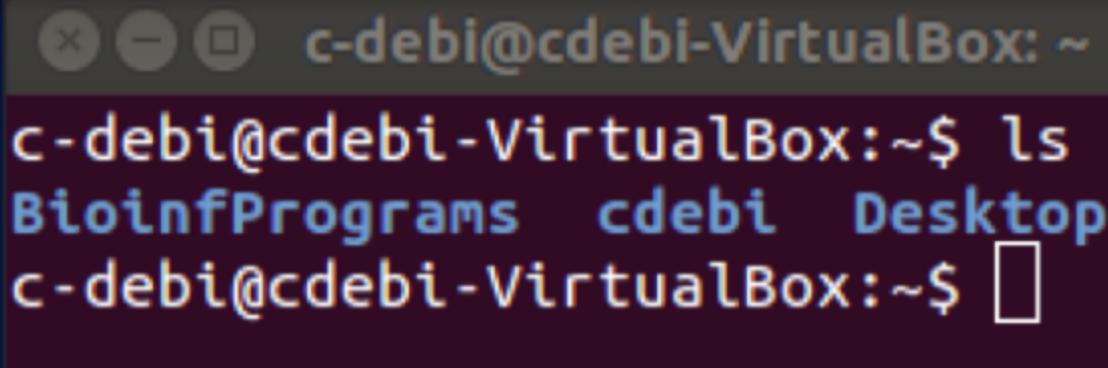
- 2 Use ls to list items in the current directory.

Command

lists items in the current directory

```
ls
```

Expected result



c-debi@cdebi-VirtualBox: ~
c-debi@cdebi-VirtualBox: ~\$ ls
BioinfPrograms cdebi Desktop
c-debi@cdebi-VirtualBox: ~\$ █

A screenshot of a terminal window titled 'c-debi@cdebi-VirtualBox: ~'. The window shows the command 'ls' being run, which lists three directories: 'BioinfPrograms', 'cdebi', and 'Desktop'. The window has standard OS X-style window controls (close, minimize, maximize) at the top.

- 3 Many commands have additional options that can be set by a '-'

Command

lists all files/directories, including hidden files '!'

lists the long format

lists the long format, but ordered by date last modified

```
ls -a
```

```
ls -l
```

```
ls -lt
```

Expected result

```
c-debi@cdebi-VirtualBox: ~
c-debi@cdebi-VirtualBox:~$ ls
BioinfPrograms  cdebi  Desktop  Downloads  ecogeo
c-debi@cdebi-VirtualBox:~$ ls -a
.                      .com.zerog.registry.xml  .install4
..                     .config                   .InstallA
.bash_history          .dbus                    .jalview_
.bash_logout           .Dendroscope.def    .java
.bashrc                Desktop                 .jswingre
BioinfPrograms        Downloads               .kde
.biojs_templates       ecogeo                  .local
.cache                 .gconf                  .mozilla
cdebi                 .gnome                 .pki
.compiz               .ICEauthority        .profile
c-debi@cdebi-VirtualBox:~$ ls -l
total 20
drwxrwxr-x 28 c-debi c-debi 4096 Jul 17 22:13 Bioinf
drwxrwxr-x  6 c-debi c-debi 4096 Dec  8  2015 cdebi
drwxr-xr-x  2 c-debi c-debi 4096 Jul  4 10:00 Desktop
drwxr-xr-x  7 c-debi c-debi 4096 Jul 17 22:14 Downloads
drwxrwxr-x 11 c-debi c-debi 4096 Jul 17 22:13 ecogeo
c-debi@cdebi-VirtualBox:~$ ls -lt
total 20
drwxr-xr-x  7 c-debi c-debi 4096 Jul 17 22:14 Downloads
drwxrwxr-x 11 c-debi c-debi 4096 Jul 17 22:13 ecogeo
drwxrwxr-x 28 c-debi c-debi 4096 Jul 17 22:13 Bioinf
drwxr-xr-x  2 c-debi c-debi 4096 Jul  4 10:00 Desktop
drwxrwxr-x  6 c-debi c-debi 4096 Dec  8  2015 cdebi
c-debi@cdebi-VirtualBox:~$ █
```

Directory System

- 4 cd - change directory

Command

```
cd ecogeo/
```

5 List the contents of the current directory.

6 Move into the directory called **unix**

7 `pwd` (present working directory) can be used to show the current directory.

Command

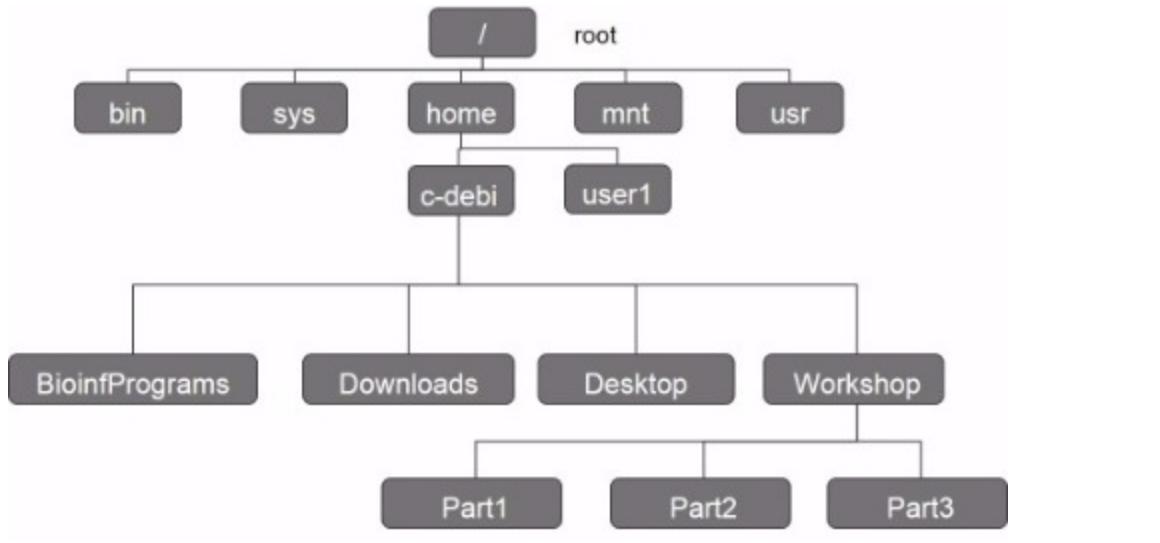
prints the path to the current directory

```
pwd
```

Expected result

```
cd /home/c-debi/ecogeo/unix
```

8 Move to the root directory.



Note

This is where everything is stored in the computer. All the commands we are running live in `/bin`.

Command

```
cd /
```

- 9 Change directory to **home**
- Change directory to **c-debi**
- Change directory to **ecogeo**
- Change directory to **unix**
- List contents
- Change directory to **data**
- Change directory to **root**

Note

Tabs can be used to auto complete names.

- 10 Change directory to **unix/data** in one step

Command

```
$ cd /home/c-debi/ecogeo/unix/data
```

- 11 cd '..' allows you to step back up through the path directory. Display present working directory path.

Command

moves back in the path directory

```
cd ..  
pwd
```

Expected result

```
/home/c-debi/ecogeo/unix
```

- 12 Step back up to the c-debi directory.
- 13 Change directory to BioinfPrograms
- 14 List contents

Expected result

```
c-debi@cdebi-VirtualBox: ~/BioinfPrograms
c-debi@cdebi-VirtualBox:~/ecogeo/unix/data$ 
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ cd .
c-debi@cdebi-VirtualBox:~/ecogeo$ cd ..
c-debi@cdebi-VirtualBox:~$ pwd
/home/c-debi
c-debi@cdebi-VirtualBox:~$ ls
BioinfPrograms  cdebi  Desktop  Downloads
c-debi@cdebi-VirtualBox:~$ cd BioinfPrograms
c-debi@cdebi-VirtualBox:~/BioinfPrograms$ l
amos-2.0.8          FastQC
anvio-2.0.2          FastTree
anvio-2.0.2.tar.gz   FigTree_v1.4.2
anvi-ubuntu-setup.sh hmmer-3.1b2-linu
AUTHORS             idba-1.1.1
bin                 include
bowtie-1.1.2         Jalview
building.html        jalview.jar
cutadapt            Jalview.lax
dendroscope         lax.jar
Dendroscope_unix_3_5_7.sh lib
diamond              LICENSE
EMIRGE               megahit
ESOM                 MetaRNA_to_FastQ
examples             mothur
c-debi@cdebi-VirtualBox:~/BioinfPrograms$
```

15 Change directory to unix/

- 16 Make a directory named "storage".

Command

```
mkdir storage
```

- 17 List contents of directory.

- 18 Move into the storage directory.

Manipulating files

- 19 The 'touch' command allows you to create a blank file of the input name.

Command

creates a blank file of the input name

```
touch temp.txt
```

- 20 The 'cp' command allows you to copy a file and can be used to move a copy of a file to a directory.

Command

```
$ cp
```

- 21 The 'mv' or move command "destroys" the original and places the content elsewhere.

Command

```
$ mv
```

- 22 Using copy:

Command

```
$ cp temp.txt newtemp.txt  
$ cp temp.txt ../
```

- 23 Change directory up a level.

- 24 List contents.

Expected result

```
c-debi@cdebi-VirtualBox: ~/ecogeo/unix
c-debi@cdebi-VirtualBox:~/ecogeo/unix/storage
c-debi@cdebi-VirtualBox:~/ecogeo/unix/storage
c-debi@cdebi-VirtualBox:~/ecogeo/unix/storage
temp.txt
c-debi@cdebi-VirtualBox:~/ecogeo/unix/storage
c-debi@cdebi-VirtualBox:~/ecogeo/unix/storage
newtemp.txt temp.txt
c-debi@cdebi-VirtualBox:~/ecogeo/unix/storage
c-debi@cdebi-VirtualBox:~/ecogeo/unix/storage
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ ls
data storage temp.txt
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ █
```

25 Change directory to storage.

26 Utilize move command:

Command

```
$ mv newtemp.txt oldtemp.txt
$ mv oldtemp.txt /home/c-debi/ecogeo/unix/data
```

27 Change directory to data, list content.

28 List current working directory.

Command

```
/home/c-debi/ecogeo/unix/data
```

29 The 'rm' remove command deleted a file PERMANENTLY

Command

```
rm oldtemp.txt
```

30 Change directory to **storage**.

31 Remove **temp.txt**

32 Change directory to **unix**

33 Remove storage directory:

Command

```
$ rm -r storage
```

Expected result

```
c-debi@cdebi-VirtualBox:~/ecogeo/un  
c-debi@cdebi-VirtualBox:~/ecogeo/un  
temp.txt  
c-debi@cdebi-VirtualBox:~/ecogeo/un  
c-debi@cdebi-VirtualBox:~/ecogeo/un  
c-debi@cdebi-VirtualBox:~/ecogeo/un  
data storage temp.txt  
c-debi@cdebi-VirtualBox:~/ecogeo/un  
c-debi@cdebi-VirtualBox:~/ecogeo/un  
data temp.txt  
c-debi@cdebi-VirtualBox:~/ecogeo/un
```

- 34 Create a directory called **bestdirectoryever**
Change directory to **bestdirectoryever**
Create a file called **glam.txt**
Change **glam.txt** to **formerglam.txt**
Remove **formerglam.txt**
Change directory to **unix**
Remove **bestdirectoryever**

Expected result

```
c-debi@cdebi-VirtualBox: ~/ecogeo/unix
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ mkdir bestdir
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ cd bestdir
c-debi@cdebi-VirtualBox:~/ecogeo/unix/bestdirectoryever
c-debi@cdebi-VirtualBox:~/ecogeo/unix/bestdirectoryever/glam.txt
c-debi@cdebi-VirtualBox:~/ecogeo/unix/bestdirectoryever/formerglam.txt
c-debi@cdebi-VirtualBox:~/ecogeo/unix/bestdirectoryever
c-debi@cdebi-VirtualBox:~/ecogeo/unix/bestdirectoryever
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ ls
bestdirectoryever  data
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ rm -r bestdir
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ ls
data
c-debi@cdebi-VirtualBox:~/ecogeo/unix$ █
```

35 Change directory to data.

36 List contents.

37 Remove oldtemp.txt

38 group12_contigs.fasta
group20_contigs.fasta
group24_contigs.fasta

FASTA files - specific format

> Header line, contains ID and information about...

ATGATAGCTAGCAGCAGCTA[...] 80bp and then a newline.

Looking at the contents of a file

- 39 'head' will allow you to view the first 10 lines of a file.

Command

default displays the first 10 lines

```
$ head [filename]
```

- 40 'tail' allows you to view the last 10 lines of a file.

Command

default displays last 10 lines

```
$ tail [filename]
```

- 41 'less' allows you to scroll through a file using arrow keys or spacebar = advanced page | b = reverse page | q = quit

Command

```
$ less [filename]
```

- 42 Use head to display the first 10 lines of **group12_contigs.fasta**
Display the first 5 lines of **group12_contigs.fasta**
Display the last 10 lines of **group12_contigs.fasta**
Display the last 5 lines of **group12_contigs.fasta**

43 grep - file pattern searcher

Command

```
$ grep
```

44 wc - count the number of words, lines, characters

45 Use grep on group12_contigs.fasta

46 How many? Combine grep and wc?

Use the "|" (pipe) symbol

47 Repeat but add the option -l to wc

48 Use the same technique to determine the number of sequences in
group20_contigs.fasta

49 What about the number of matches to "47" in **group12_contigs.fasta**?
Or "_47"?

Note

```
grep '>' group12_contigs.fasta | grep 47
```

50 Redirecting output to file:

51 Look at the contents of **group12_ids**

52 cat - has multiple functions:

Command

With a single input - prints file contents

```
$ cat group12_ids_with_47
```

- 53 With '>' cat has the same function as cp

Command

```
$ cat group12_ids_with_47 > temp1_ids  
$ cp group12_ids_with_47 temp2_ids
```

- 54 Double check to make sure **temp1_ids = temp2_ids**

- 55 Concatenate files with cat - most important function:

- 56 Check contents of duplicate_ids using less or cat

- 57 Grab all of the contigs IDs from **group20_contigs.fasta** that contain the number "51"

Command

```
$ grep 51 group20_contigs.fasta
```

- 58 Concatenate the new IDs to the duplicate_ids file in a file called **multiple_ids**

- 59 uniq - can be used to remove duplicates or identify lines with 1 occurrence or multiple occurrences

Command

```
$ uniq
```

- 60 sort - sort lines in a file alphanumerically

Command

```
$ sort
```

- 61 Compare **multiple_ids** before and after uniq

Command

```
$ uniq multiple_ids
```

- 62 Why was there no change?
uniq has a weakness, can only identify duplicates in adjacent lines

- 63 Clear all present files with temp in title

Command

'*' - acts as a wildcard, so any file that starts with temp would be identified and removed, no matter the suffix

```
$ rm temp*
```

- 64 How do **temp1_ids** & **temp2_ids** compare?

Command

```
$ sort multiple_ids | uniq -d > temp1_ids  
$ sort multiple_ids | uniq -u > temp2_ids
```

- 65 Identify duplicates:

- 66 Identify unique entries:

- 67 **temp1_ids = group12_ids_with_47 &**
temp2_ids = group20_ids_with_51

- 68 Remove all present files with temp in title

- 69 sed - modify files a file based on the issued commands

Command

```
$ sed
```

- 70 Want a list of sequence IDs without the '>?'

Note

`sed 's/C/c/'`
between the single quotes, substitute the occurrence of upper case C to lower case c

Command

```
$ sed 's/C/c/' clean_ids  
$ sed 's/_/./' clean_ids  
$ sed 's/>//' clean_ids > newclean_ids
```

- 71 seqmagick

Wrapper designed to utilize built in Biopython modules to manipulate and change FASTA files

Requires Biopython

<http://fhcrc.github.io/seqmagick/>

- 72 Discuss:

convert - produce a modified new file

mogrify - change the input file

info - present information of files in a directory

Additionally: backtrans-align, extract-ids, quality-filter, and primer-trim

Command

```
$ seqmagick
```

- 73 Execute seqmagick convert:

Command

```
$ seqmagick convert --include-from-file newclean_ids  
group12_contigs.fasta newgroup12_contigs.fasta
```

- 74 How many sequences are in **newgroup12_contigs.fasta**? Using grep '>':

Command

```
$ seqmagick extract-ids newgroup12_contigs.fasta | wc  
$ seqmagick info *fasta
```

Expected result

```
c-debi@cdebi-VirtualBox:~/Workshop/Part1_Unix/data$ seqmagick info *fasta  
name      alignment  min_len  max_len  avg_len  num_seqs  
group12_contigs.fasta  FALSE      5136    116409   22974.30    132  
group20_contigs.fasta  FALSE      5029    22601     7624.38    203  
group24_contigs.fasta  FALSE      5024    81329   12115.70    139  
newgroup12_contigs.fasta FALSE      5587    30751   16768.14      7  
c-debi@cdebi-VirtualBox:~/Workshop/Part1_Unix/data$ 
```

- 75 Store the information generated by 'seqmagick info' in a new file
fasta_info

Command

cut - pulling out columns from a table file
-d allows for the assignment of the type of delimiter between fields, if not TAB
-f delineates which fields to preserve, starting at 1

```
$ cut  
$ cut -f 2 fasta_info  
$ cut -f 2,4 fasta_info  
$ cut -f 2-4 fasta_info
```

Some additional tools

- 76 history - prints a sequential list of all commands in the current session

echo \$PATH - lists the directories for which the OS is checking for commands and data

- 77 nano - in window text editor

Command

Additional text can be entered like any text editor
To close out - Ctrl+X, hit 'Y', then ENTER
Create a new file - nano and then enter file name after Ctrl+X

```
$ nano fasta_info
```

- 78 Simple bash scripts: Text file with a list of commands that can be executed as a batch.
Look at the contents of **simplebashscript**
- 79 chmod - change file modes

Note

```
chmod 755 simplebashscript
```

Command

```
$ chmod 775 simplebashscript
```

- 80 Plain text file → executable text file.

Command

```
$ ./simplebashscript
```