

Cloud Computing and Unix: An Introduction

Dr. Sophie Shaw

University of Aberdeen, UK

s.shaw@abdn.ac.uk



Aberdeen

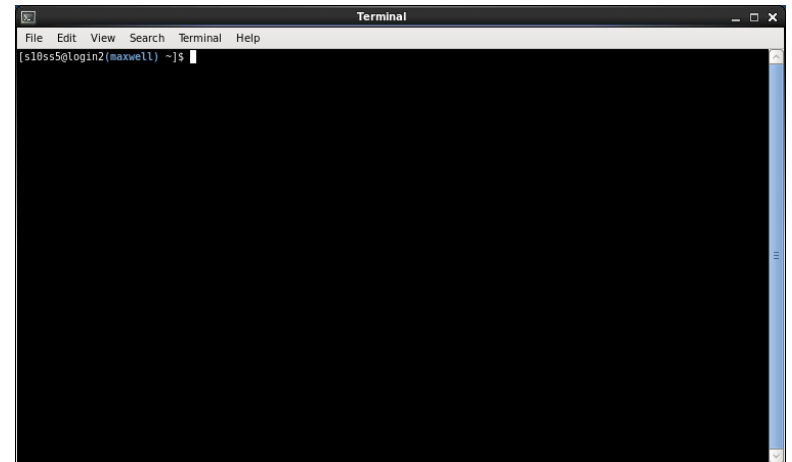


London

Exeter

What We're Going To Do

- Why Unix?
- Cloud Computing
- Connecting to AWS
- Introduction to Unix Commands



Etiquette

- PowerPoint interspersed with Challenges
- Ask me questions
- Ask demonstrators
- Work together
- Cheat!



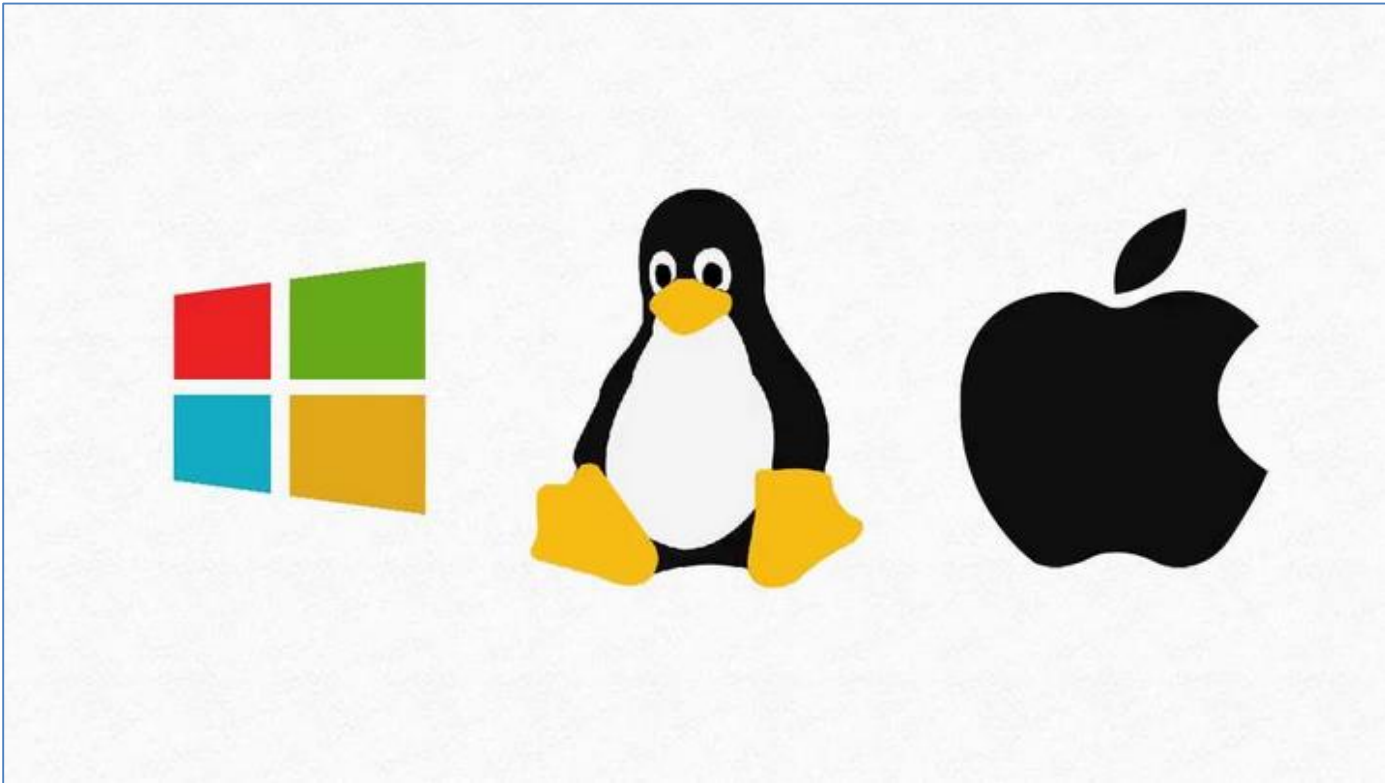
File Commands	System Info
ls - directory listing	date - show the current date and time
ls -al - formatted listing with hidden files	cal - show this month's calendar
cd dir - change directory to <i>dir</i>	uptime - show current uptime
cd - change to home	w - display who is online
pwd - show current directory	whoami - who you are logged in as
mkdir dir - create a directory <i>dir</i>	finger user - display information about <i>user</i>
rm file - delete <i>file</i>	uname -a - show kernel information
rm -r dir - delete directory <i>dir</i>	cat /proc/cpuinfo - cpu information
rm -f file - force remove <i>file</i>	cat /proc/meminfo - memory information
rm -rf dir - force remove directory <i>dir</i> *	man command - show the manual for <i>command</i>
cp file1 file2 - copy <i>file1</i> to <i>file2</i>	df - show disk usage
cp -r dir1 dir2 - copy <i>dir1</i> to <i>dir2</i> ; create <i>dir2</i> if it doesn't exist	du - show directory space usage
mv file1 file2 - rename or move <i>file1</i> to <i>file2</i> if <i>file2</i> is an existing directory, moves <i>file1</i> into directory <i>file2</i>	free - show memory and swap usage
ln -s file link - create symbolic link <i>link</i> to <i>file</i>	whereis app - show possible locations of <i>app</i>
touch file - create or update <i>file</i>	which app - show which <i>app</i> will be run by default
cat > file - places standard input into <i>file</i>	
more file - output the contents of <i>file</i>	Compression
head file - output the first 10 lines of <i>file</i>	tar cf file.tar files - create a tar named <i>file.tar</i> containing <i>files</i>
tail file - output the last 10 lines of <i>file</i>	tar xf file.tar - extract the files from <i>file.tar</i>
tail -f file - output the contents of <i>file</i> as it grows, starting with the last 10 lines	tar czf file.tar.gz files - create a tar with Gzip compression
Process Management	tar xzf file.tar.gz - extract a tar using Gzip
ps - display your currently active processes	tar cjf file.tar.bz2 - create a tar with Bzip2 compression
top - display all running processes	tar xjf file.tar.bz2 - extract a tar using Bzip2
kill pid - kill process id <i>pid</i>	gzip file - compresses <i>file</i> and renames it to <i>file.gz</i>
killall proc - kill all processes named <i>proc</i> *	gzip -d file.gz - decompresses <i>file.gz</i> back to <i>file</i>
bg - lists stopped or background jobs; resume a stopped job in the background	Network
fg - brings the most recent job to foreground	ping host - ping <i>host</i> and output results
fg n - brings job <i>n</i> to the foreground	whois domain - get whois information for <i>domain</i>
File Permissions	dig domain - get DNS information for <i>domain</i>
chmod octal file - change the permissions of <i>file</i> to <i>octal</i> , which can be found separately for user, group, and world by adding:	dig -x host - reverse lookup <i>host</i>
<ul style="list-style-type: none"> • 4 - read (r) • 2 - write (w) • 1 - execute (x) 	wget file - download <i>file</i>
Examples:	wget -c file - continue a stopped download
chmod 777 - read, write, execute for all	Installation
chmod 755 - rwx for owner, rx for group and world	Install from source:
For more options, see man chmod .	./configure
SSH	make
ssh user@host - connect to <i>host</i> as <i>user</i>	make install
ssh -p port user@host - connect to <i>host</i> on port <i>port</i> as <i>user</i>	dpkg -i pkg.deb - install a package (Debian)
ssh-copy-id user@host - add your key to <i>host</i> for <i>user</i> to enable a keyed or passwordless login	rpm -Uvh pkg.rpm - install a package (RPM)
Searching	Shortcuts
grep pattern files - search for <i>pattern</i> in <i>files</i>	Ctrl+C - halts the current command
grep -r pattern dir - search recursively for <i>pattern</i> in <i>dir</i>	Ctrl+Z - stops the current command, resume with fg in the foreground or bg in the background
command grep pattern - search for <i>pattern</i> in the output of <i>command</i>	Ctrl+D - log out of current session, similar to exit
locate file - find all instances of <i>file</i>	Ctrl+W - erases one word in the current line
	Ctrl+U - erases the whole line
	Ctrl+R - type to bring up a recent command
	!! - repeats the last command
	exit - log out of current session

* use with extreme caution.



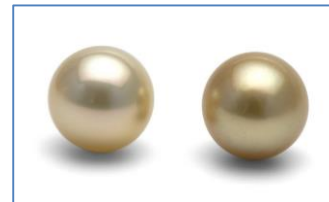
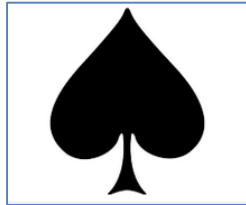
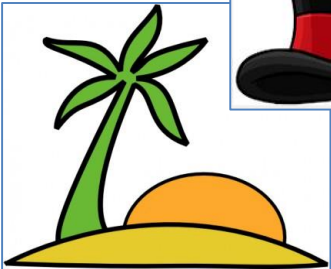
What is Unix?

- Operating System



Why Unix?

- Bioinformatics software designed to run on Unix platforms.
- Large amounts of data.
- Much faster than your Windows PC.



How Can We Use Unix?

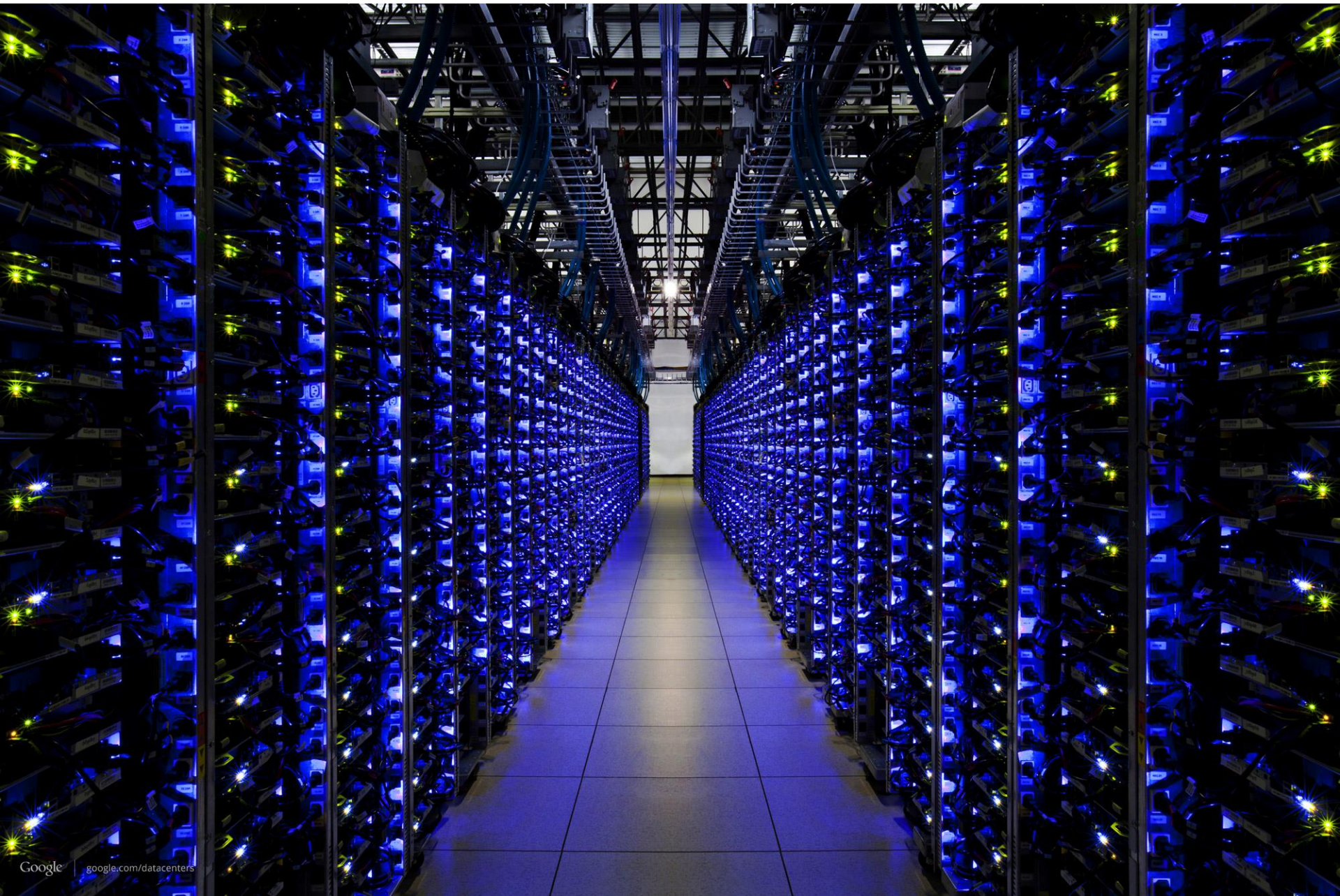
- Linux computers or servers.
- Compute clusters.
- The cloud.
 - What we're going to use this week

Download more graphics at www.psdgraphics.com



So What is Cloud Computing?



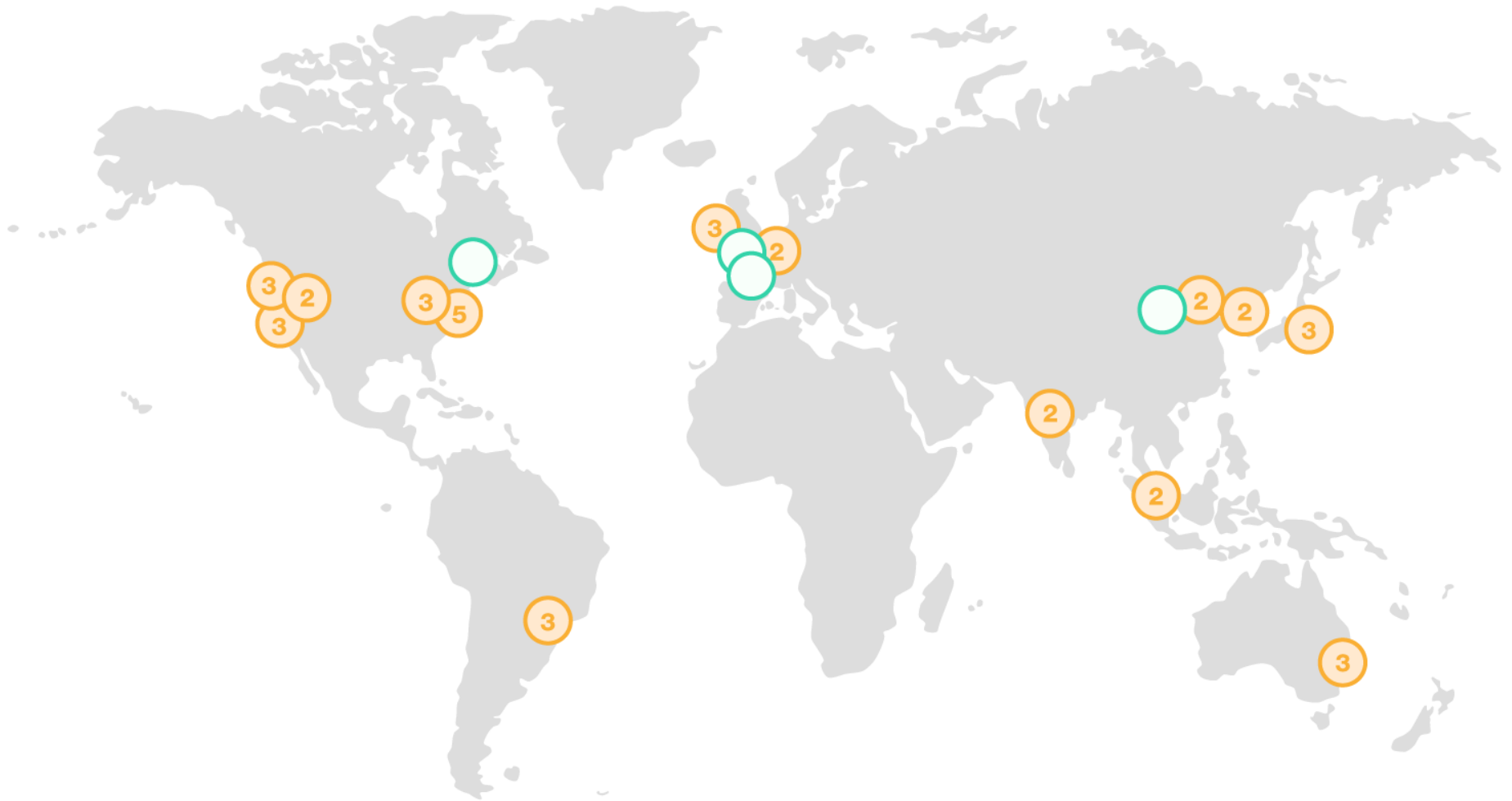


Cloud Computing Solutions



Google Compute Engine

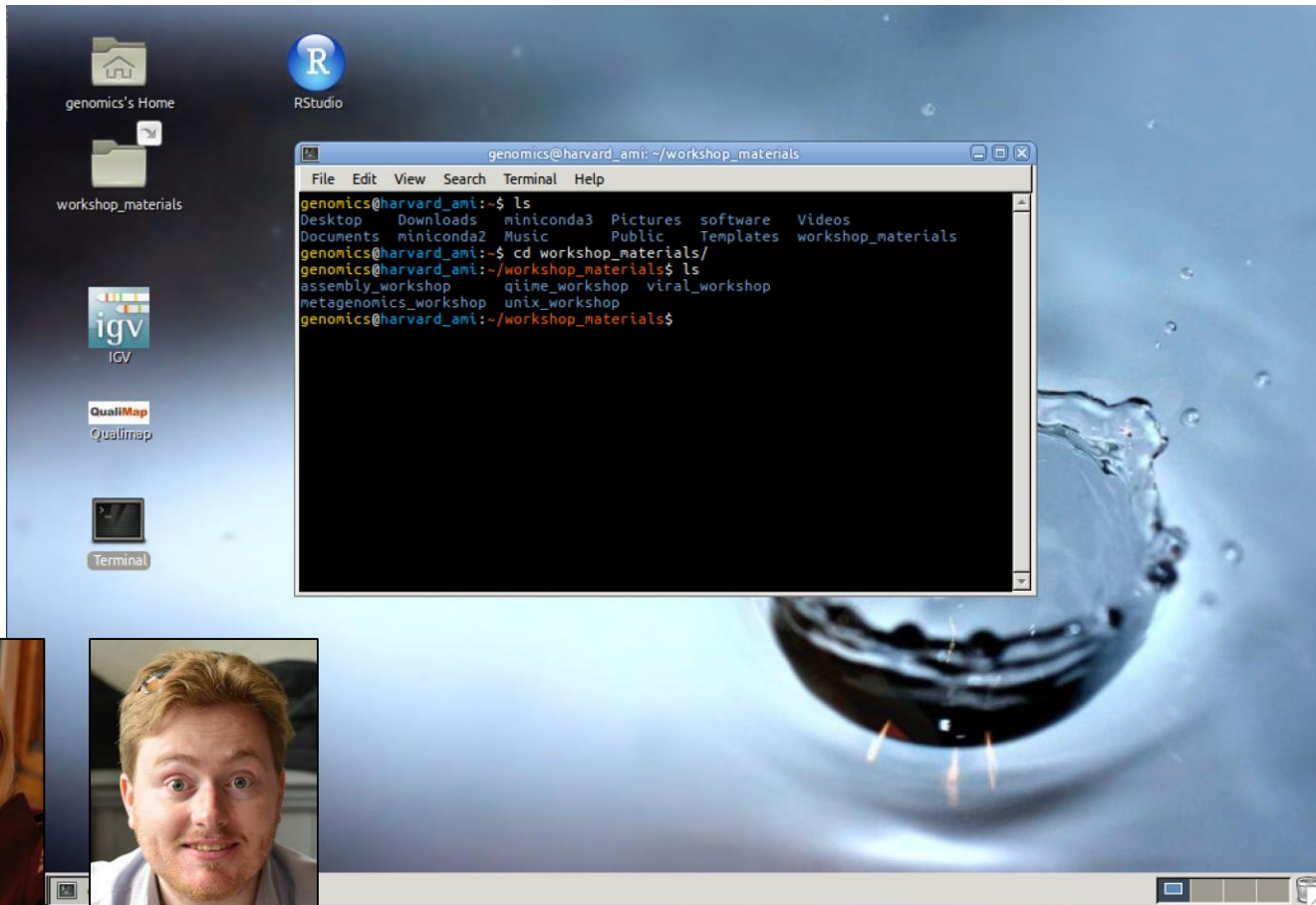




AWS “Availability Zones” and Data Centres

How it Works

AMI (“Amazon Machine Image”)
Base computer with all data and software



Terminology

- Creating an instance – *buying a brand new computer with software already installed.*
- Starting an instance – *turning that computer on.*
- Stopping an instance – *turning that computer off.*
- Terminating an instance – *setting that computer on fire and throwing it out of the window.*

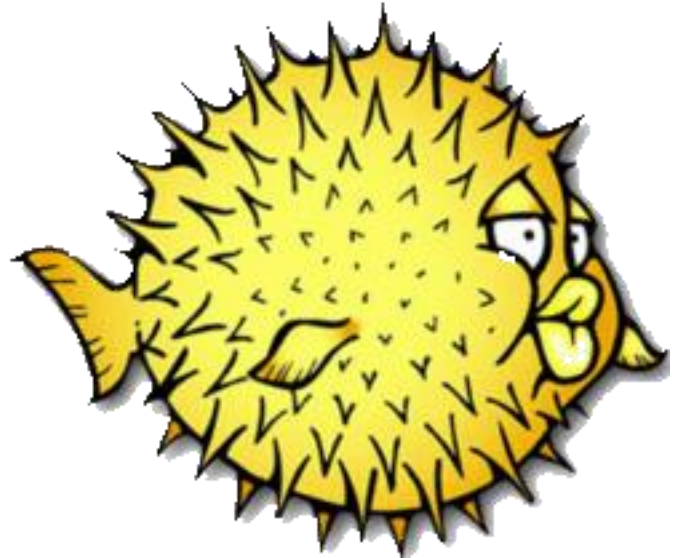
The Rules

- Only create one instance each.
- Stop your instance at the end of each day (unless you have software running).
- Name your instance (with YOUR name! No Bruce Waynes please)
- Only start or stop your own instance.
- Only terminate your own instance.

Connecting to Your Instance



Remote Desktop
Software
e.g. X2Go

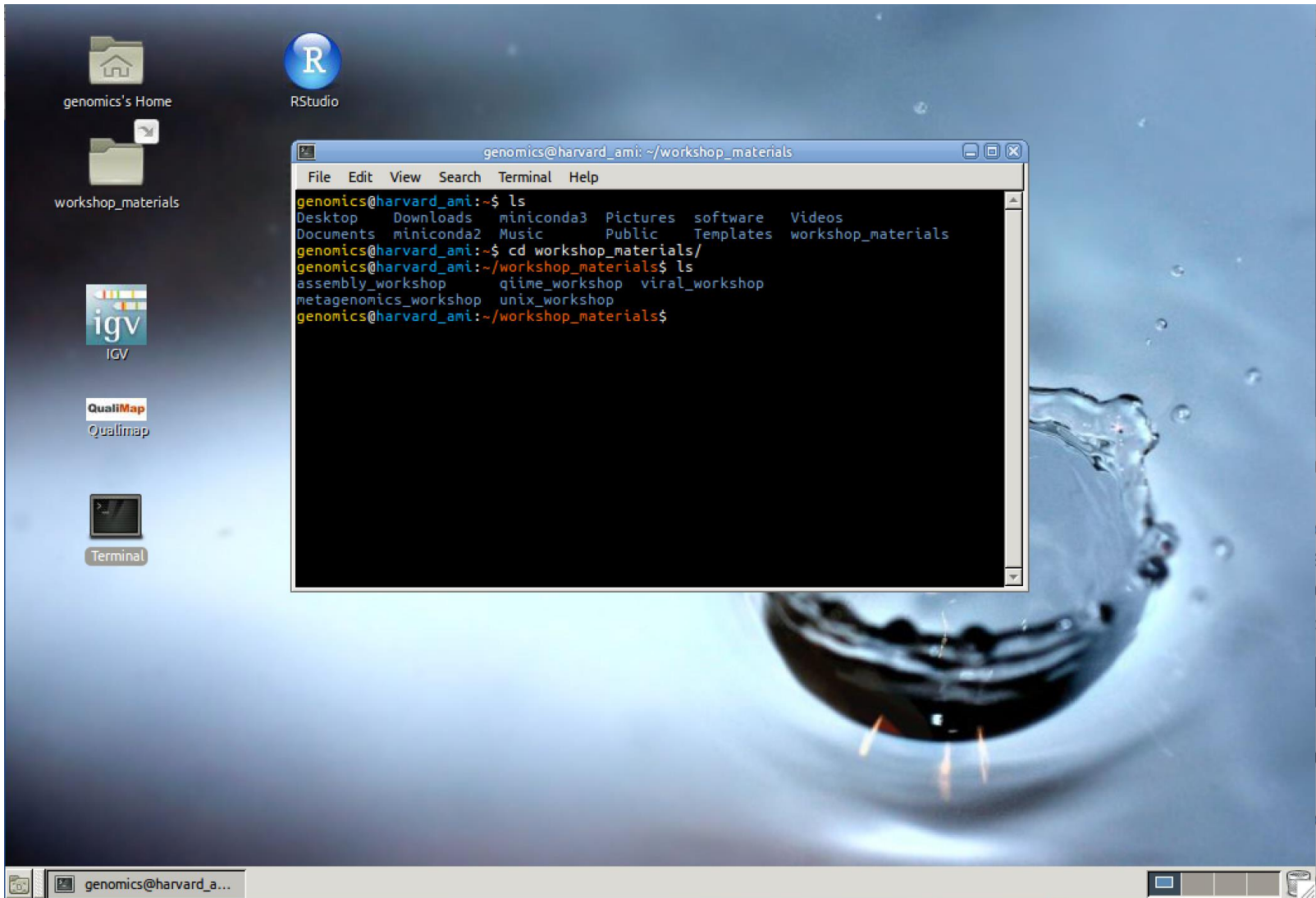


Secure Shell –
“SSH”
e.g. SSH or PuTTY

Now What?!

- You're each going to create, start and connect to your own instance.

INSERT LIVE DEMO

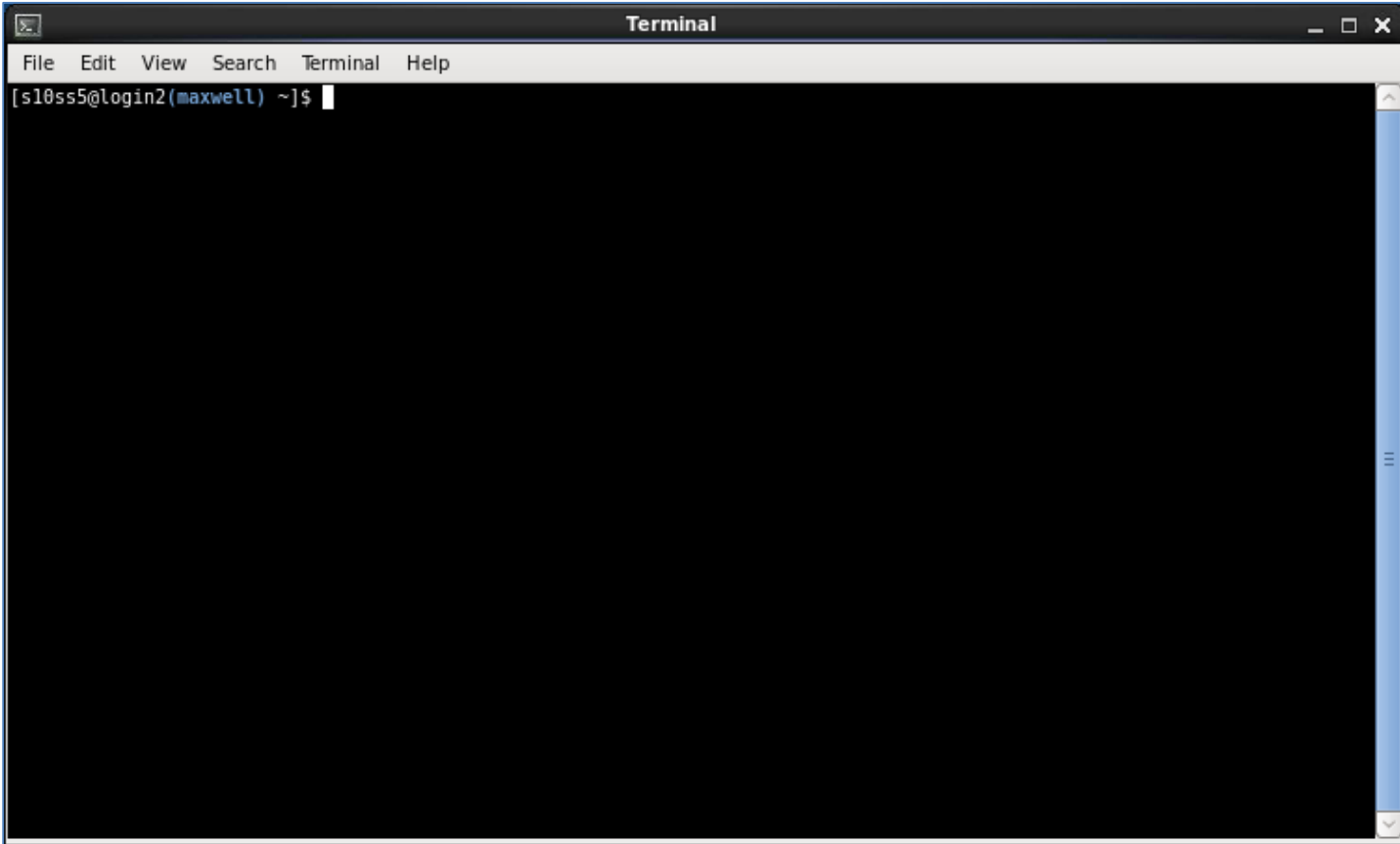


You're now connected to your instance and you're ready to learn some Unix!

Any Questions So Far?



The Terminal Window



The Command Line, The Shell, The Prompt

Where you see this “\$” followed by text, I want you to type the text on your command line

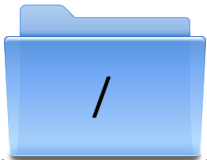
Location is Important

First Task – Where am I?

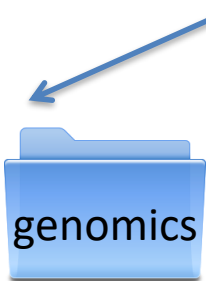
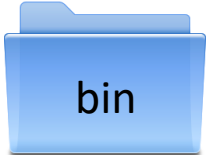
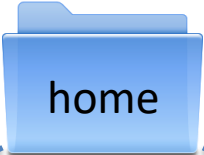
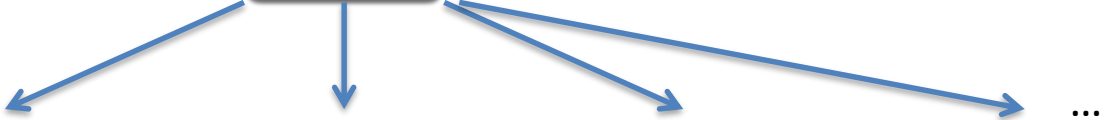
```
$ pwd
```

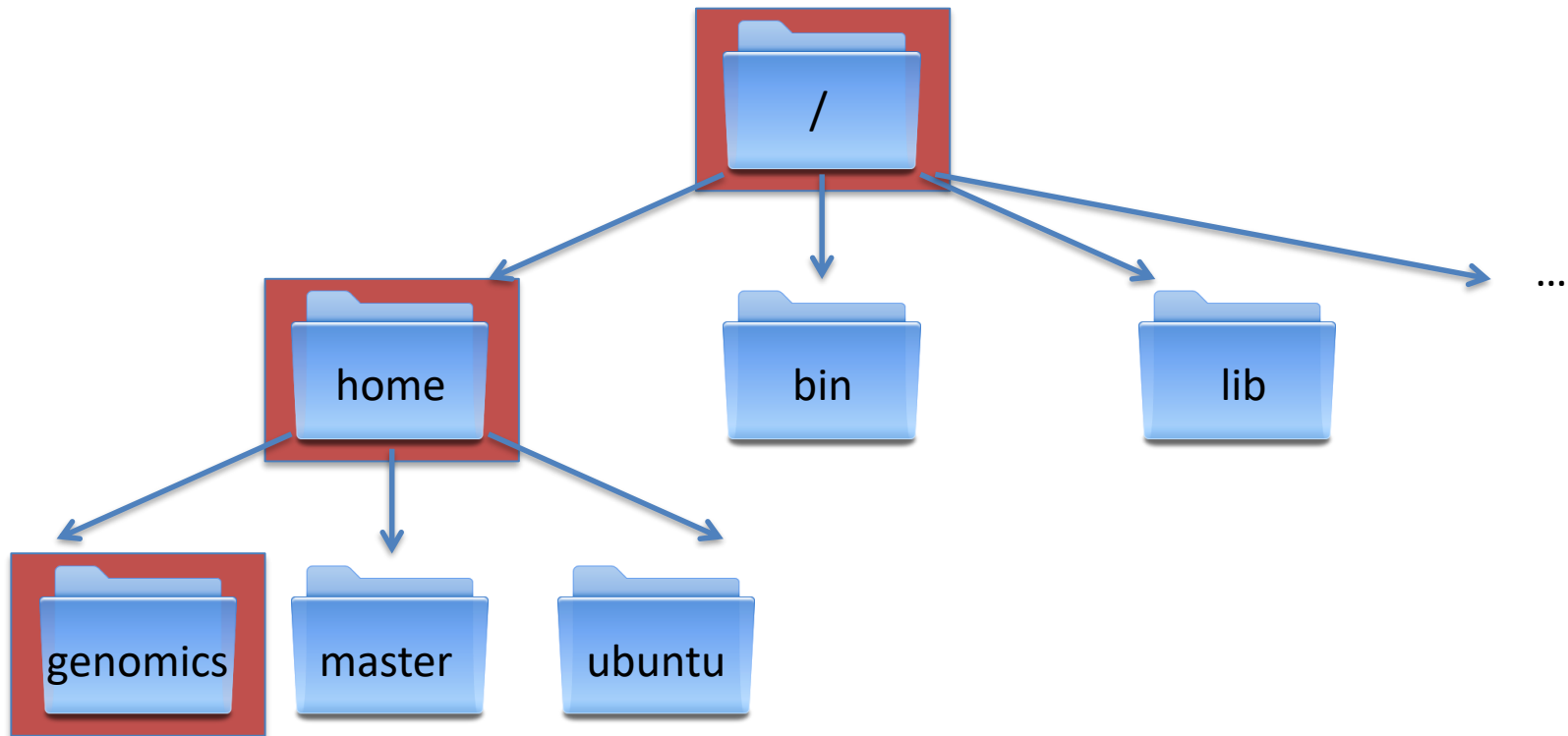
```
genomics@harvard_ami:~$ pwd  
/home/genomics  
genomics@harvard_ami:~$
```

This is your “present working directory”

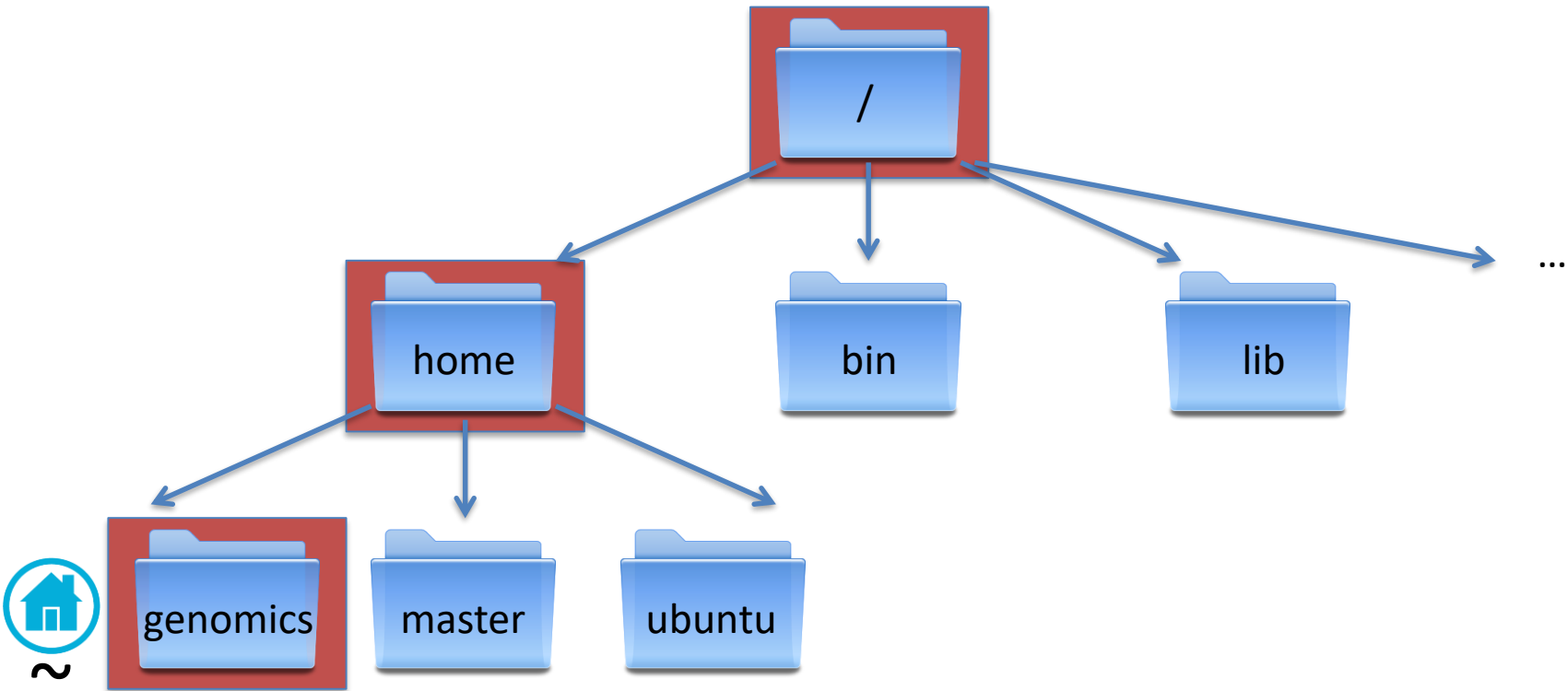


= ROOT





```
genomics@harvard_ami:~$ pwd  
/home/genomics  
genomics@harvard_ami:~$
```



```
genomics@harvard_ami:~$ pwd
/home/genomics
genomics@harvard_ami:~$
```

This location is also known as your Home Directory

Tilde is shorthand for Home:

~

Now let's create some directories and files

Make a directory

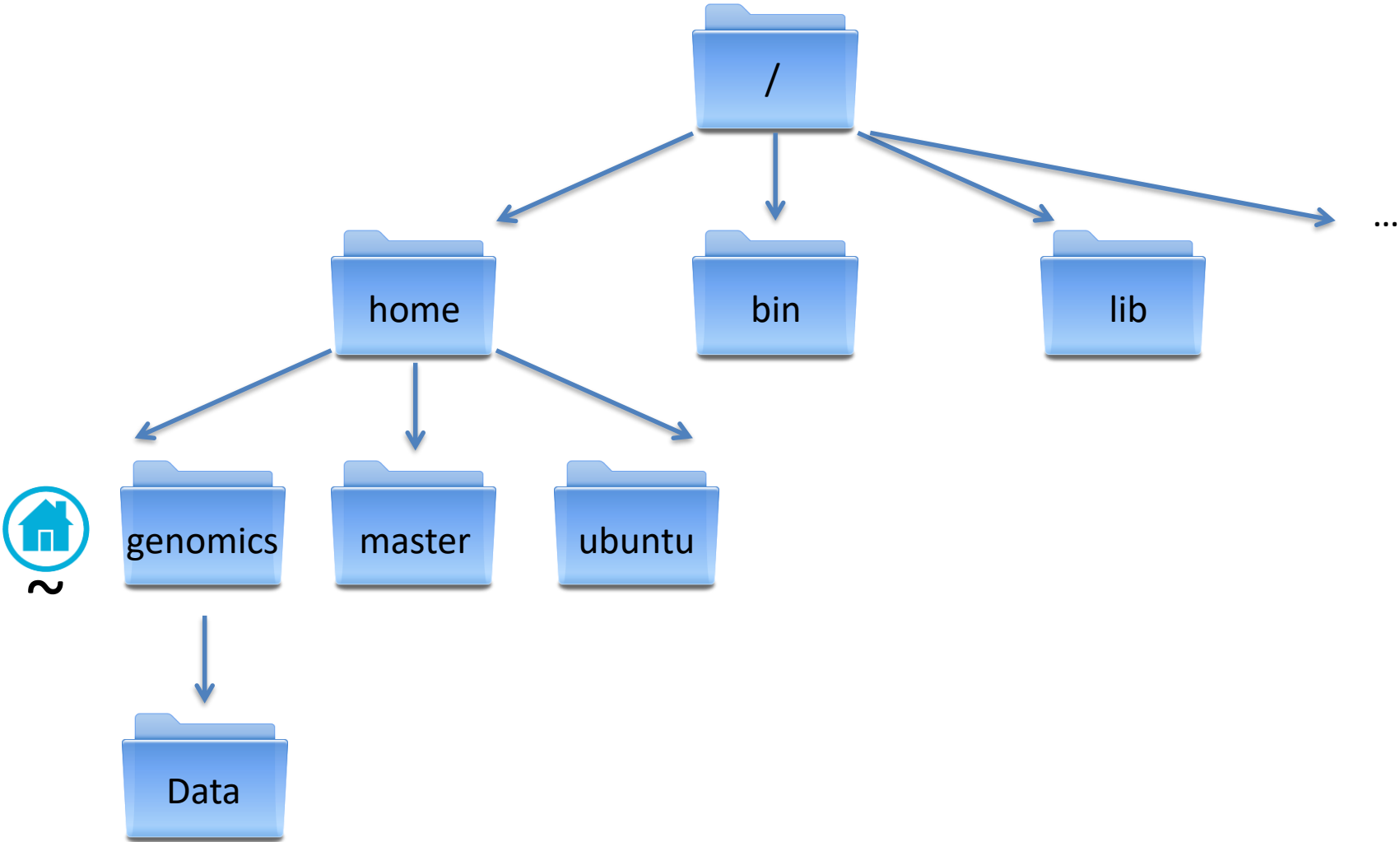
```
$ mkdir Data
```

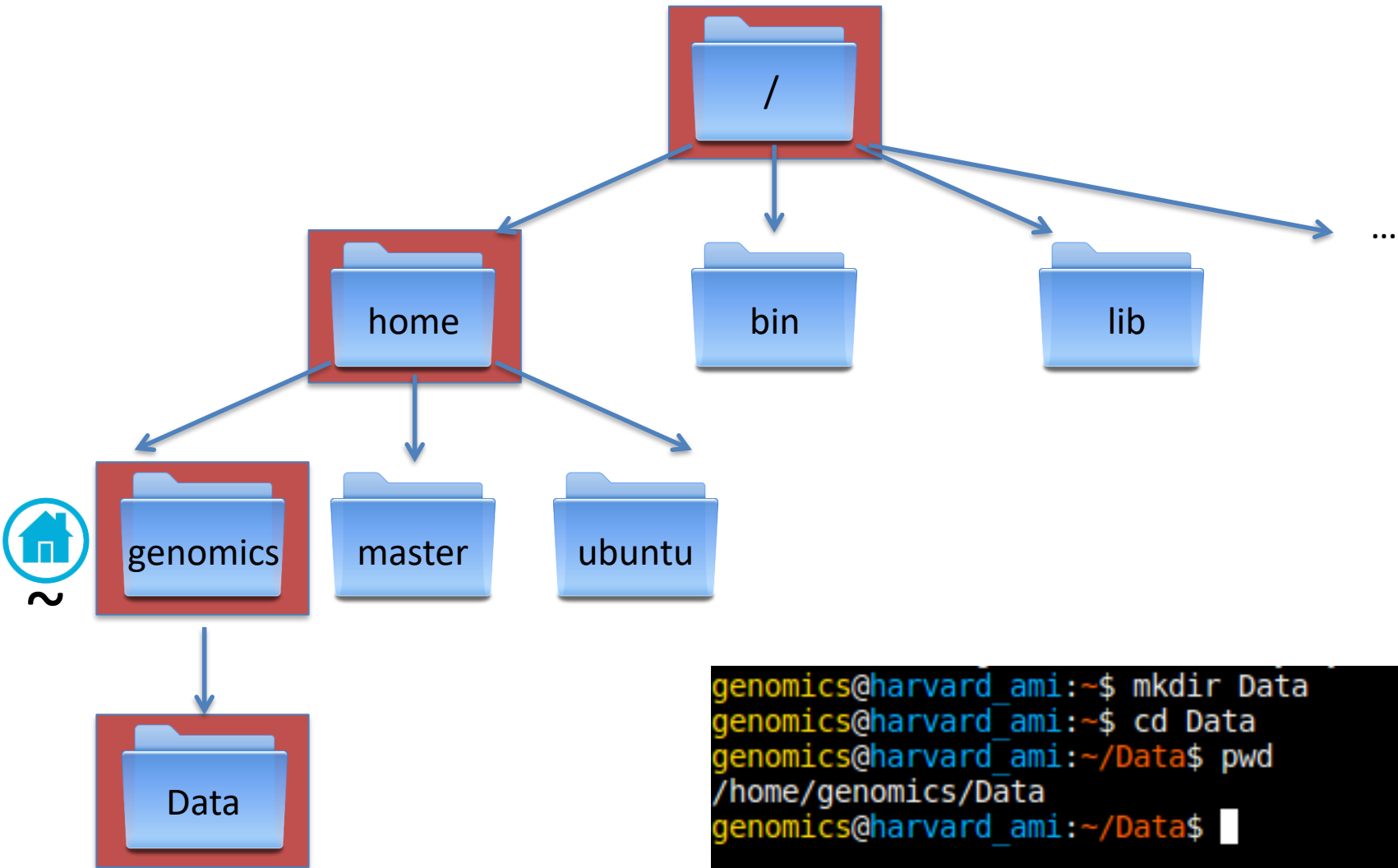
Change into this directory

```
$ cd Data
```

Now what is your present working directory?

NOTE! Directory names (and file names for the matter) can not contain spaces. Underscores are often used instead if you want to separate words.





```
genomics@harvard_ami:~$ mkdir Data
genomics@harvard_ami:~$ cd Data
genomics@harvard_ami:~/Data$ pwd
/home/genomics/Data
genomics@harvard_ami:~/Data$
```

Now let's create some directories and files

Make an empty file

```
$ touch rags
```

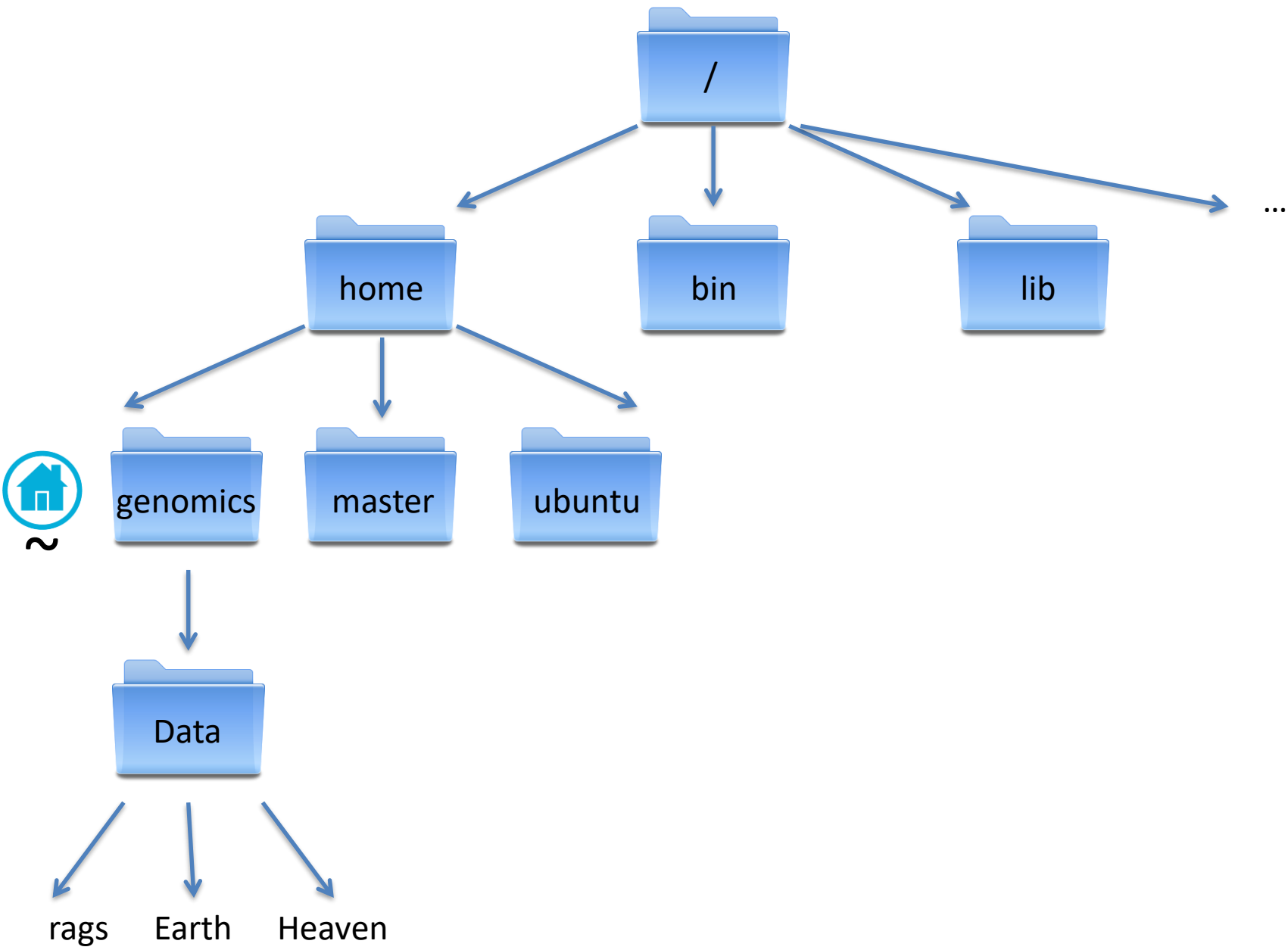
And another two

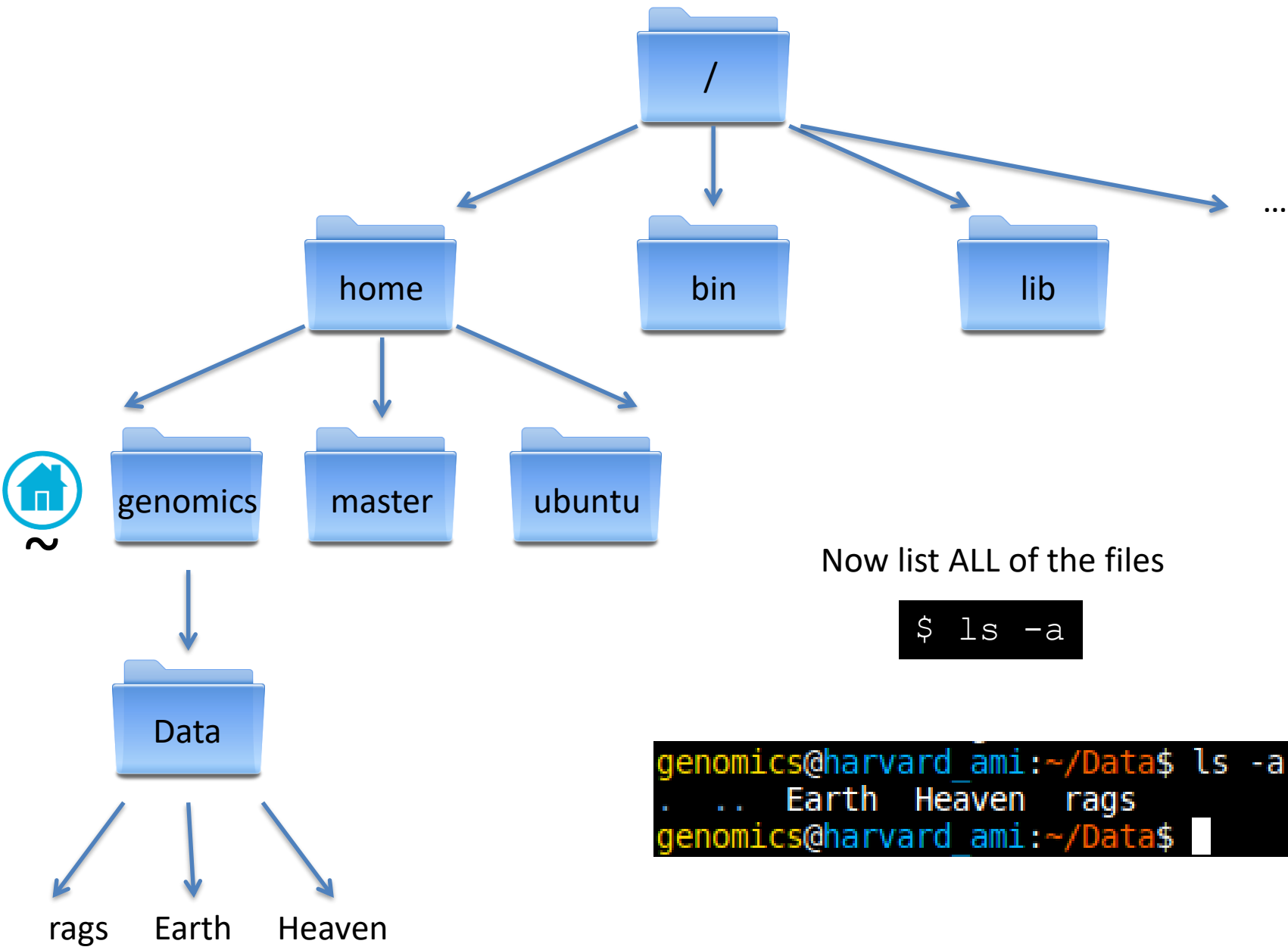
```
$ touch Earth Heaven
```

Now let's list the contents of the current directory (Data)

```
$ ls
```

```
genomics@harvard_ami:~/Data$ touch rags
genomics@harvard_ami:~/Data$ touch Earth Heaven
genomics@harvard_ami:~/Data$ ls
Earth Heaven rags
genomics@harvard_ami:~/Data$
```

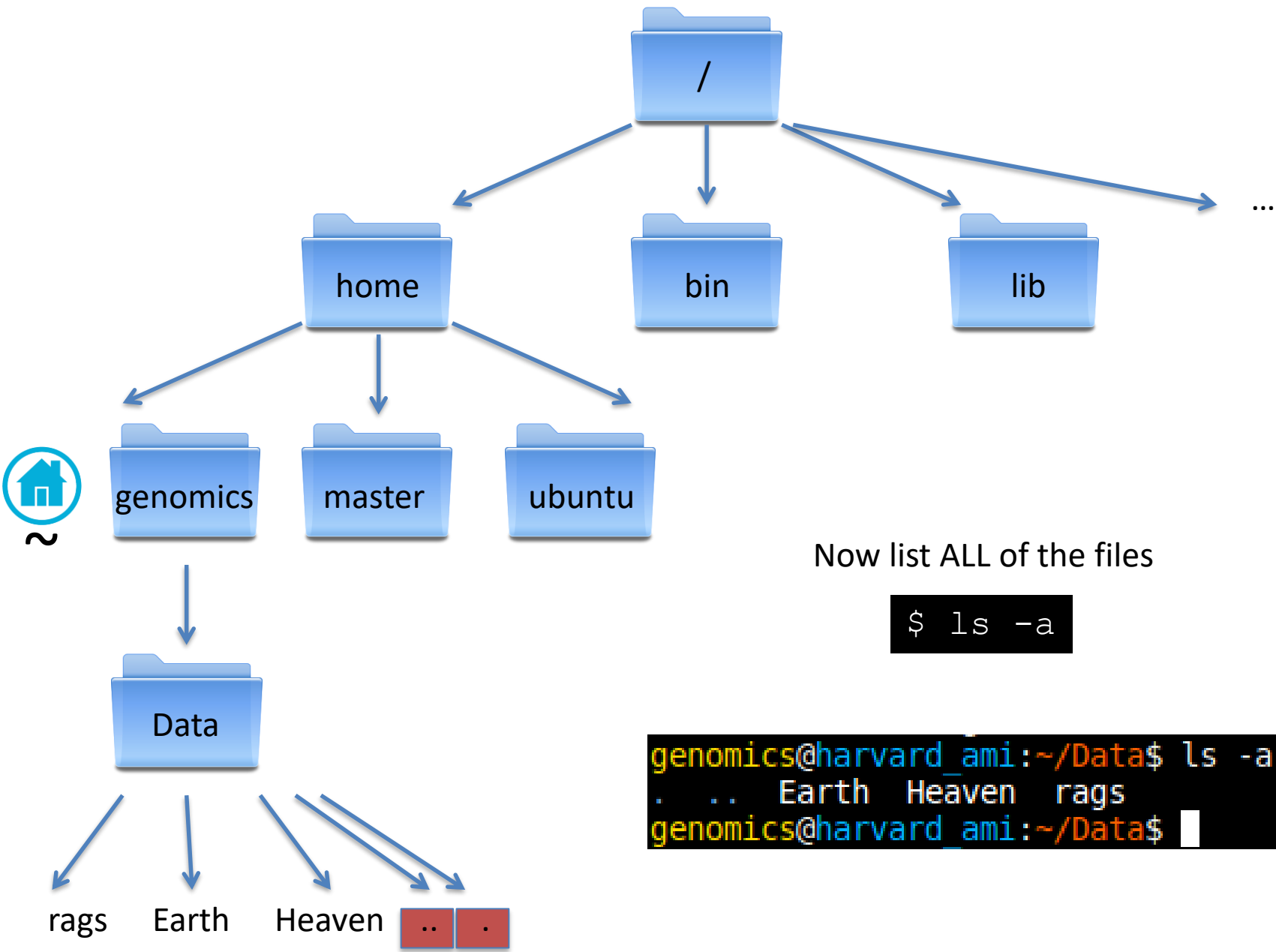




Now list ALL of the files

```
$ ls -a
```

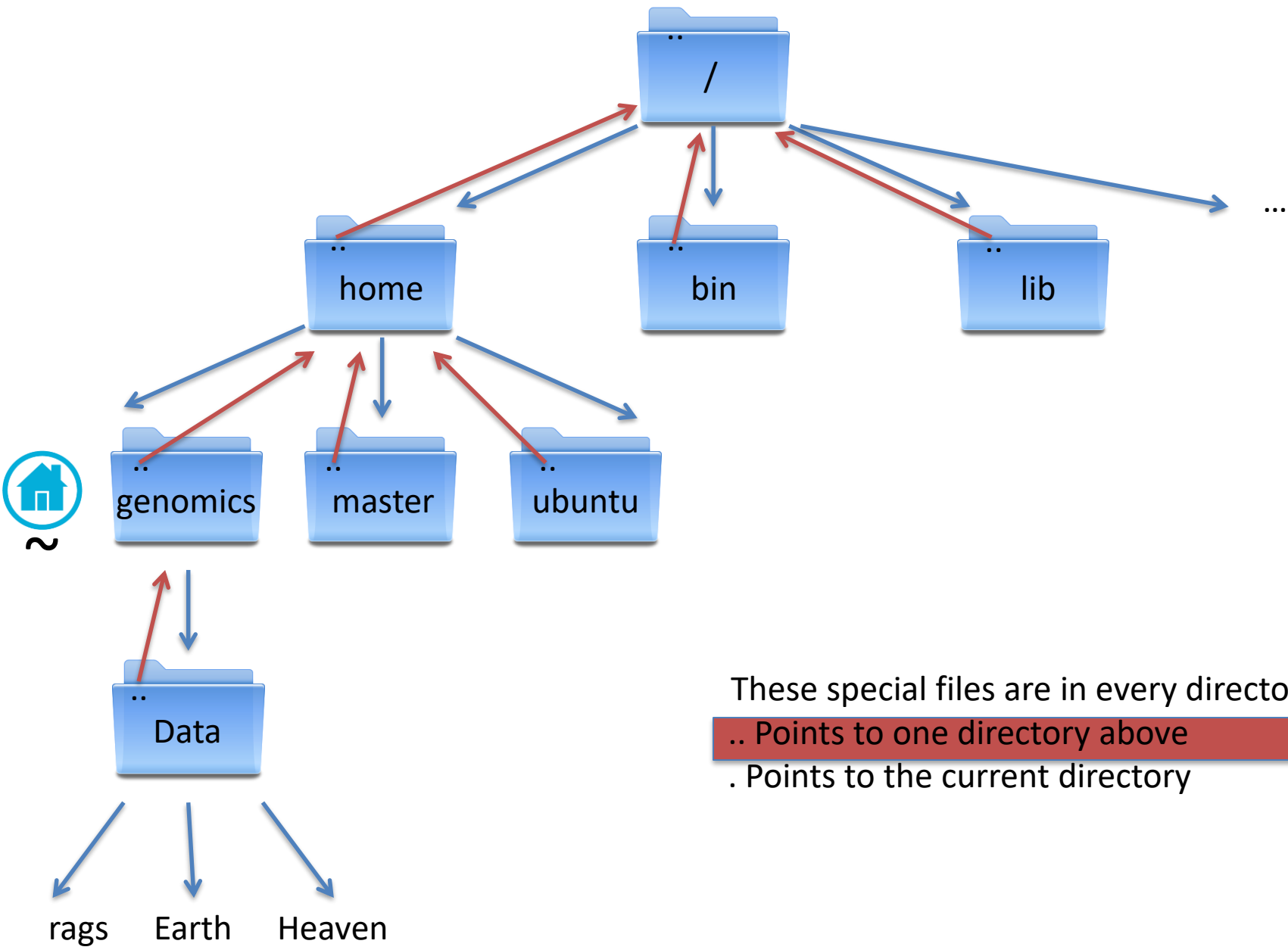
```
genomics@harvard_ami:~/Data$ ls -a
.  ..  Earth  Heaven  rags
genomics@harvard_ami:~/Data$
```

Now list ALL of the files

```
$ ls -a
```

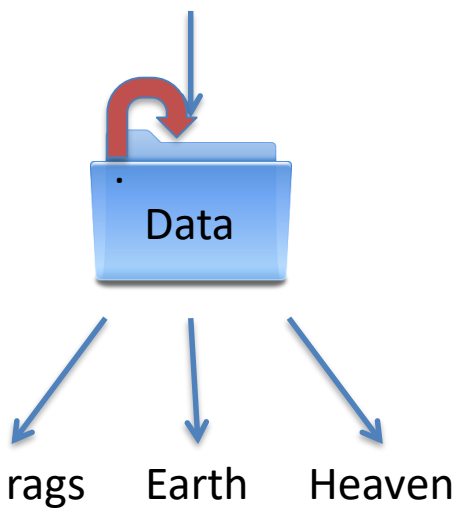
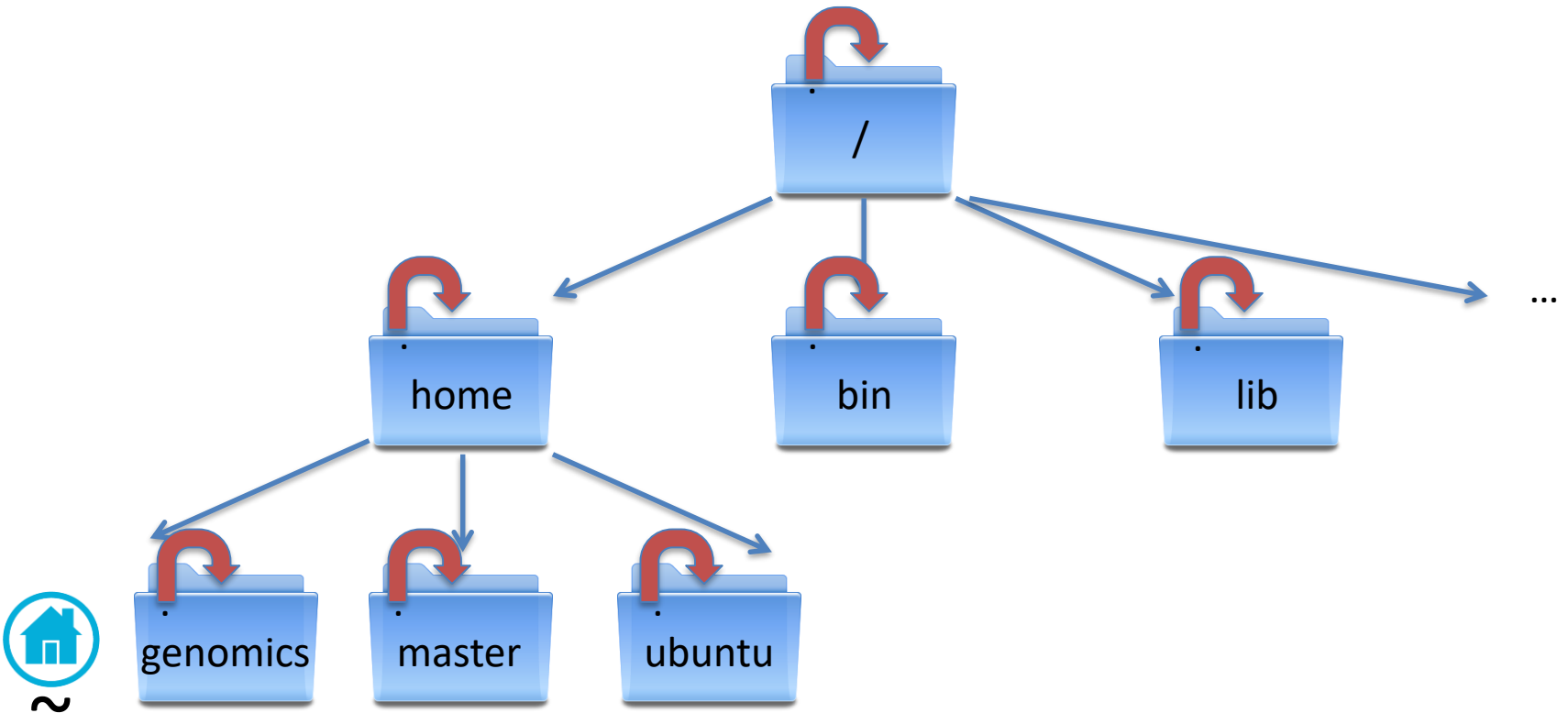
```
genomics@charvard_ami:~/Data$ ls -a
.  ..  Earth  Heaven  rags
genomics@charvard_ami:~/Data$
```



These special files are in every directory

.. Points to one directory above

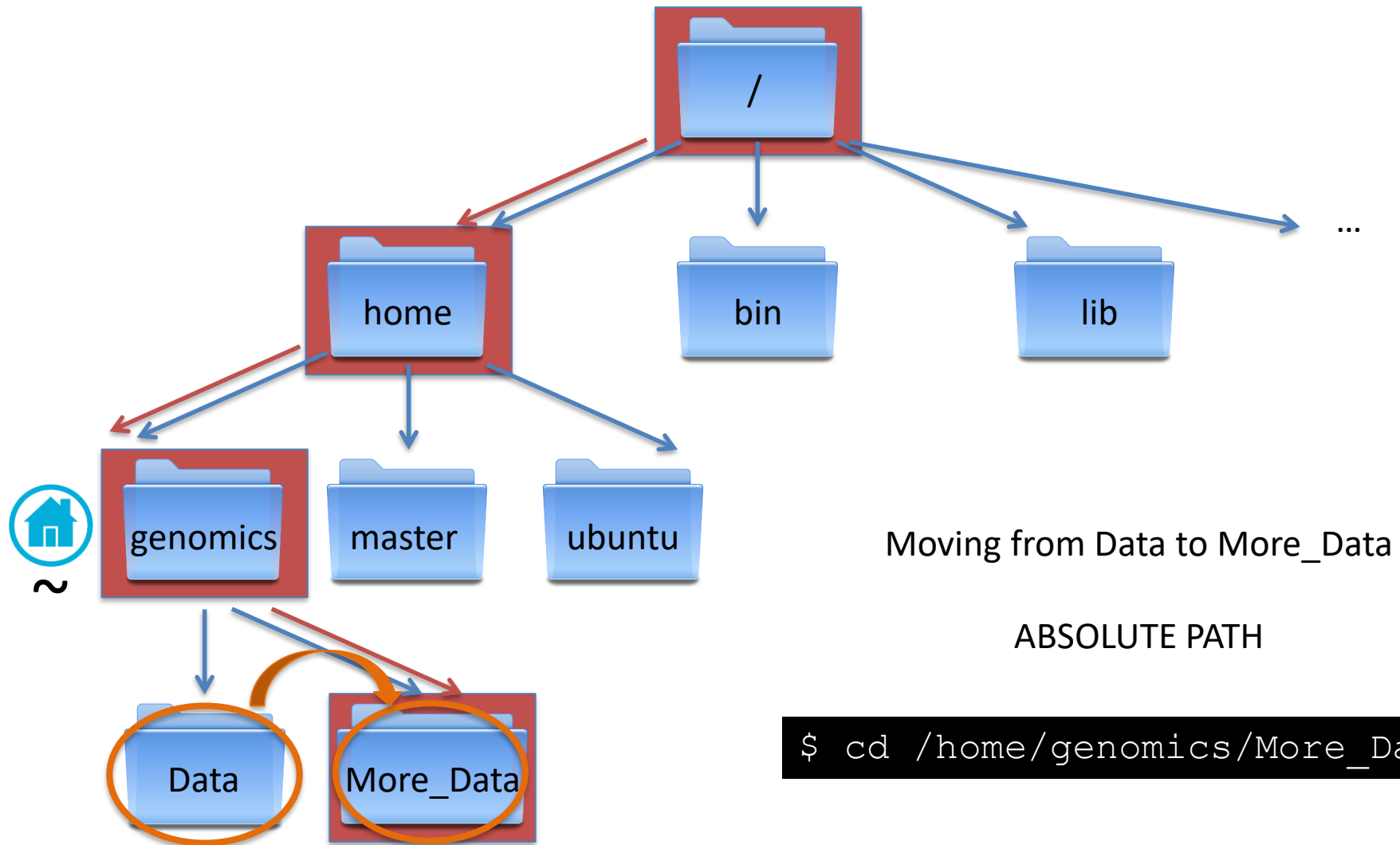
. Points to the current directory



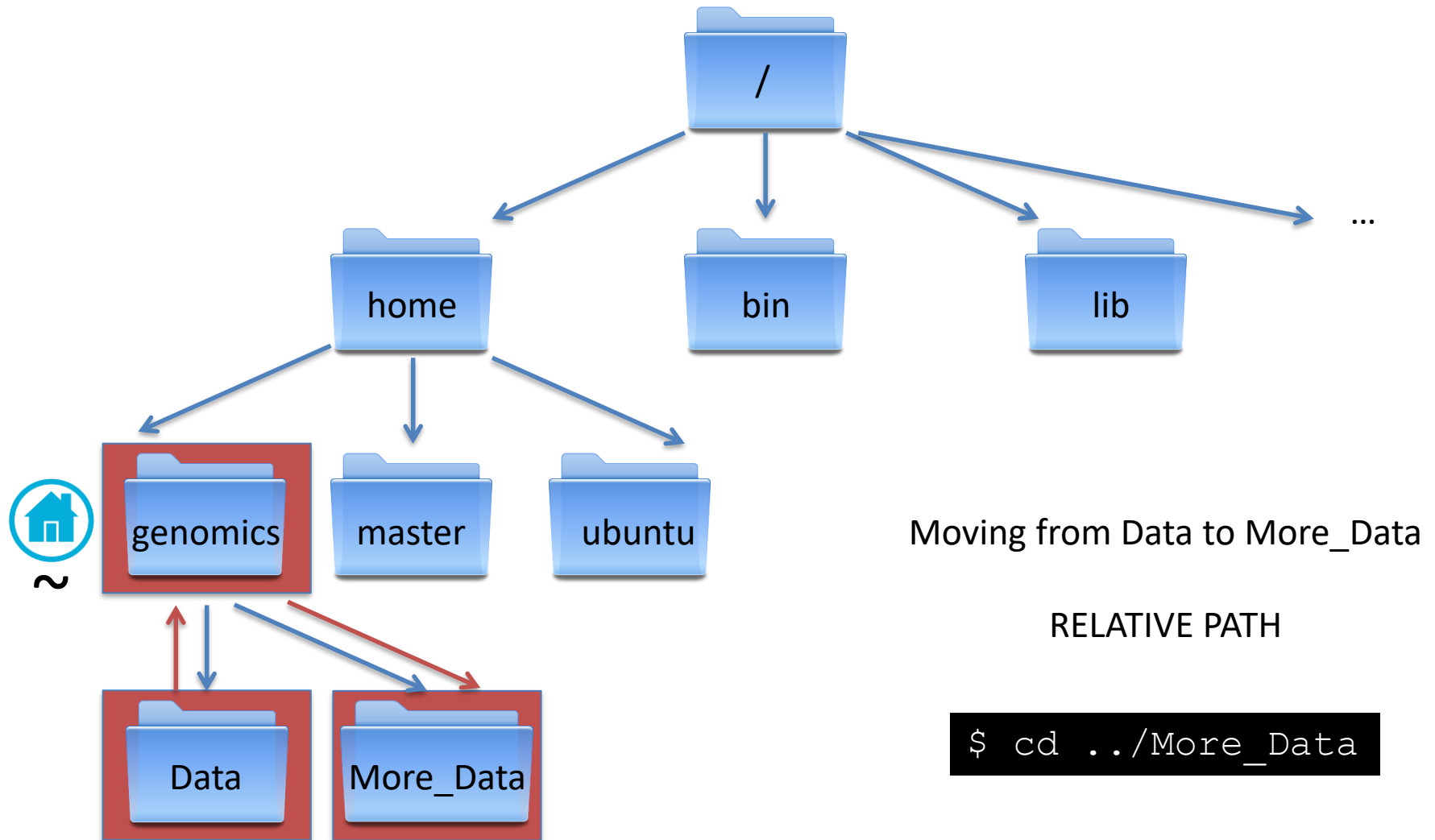
These special files are in every directory
.. Points to one directory above
. Points to the current directory

ABSOLUTE AND RELATIVE PATHS: GETTING FROM ONE PLACE TO ANOTHER

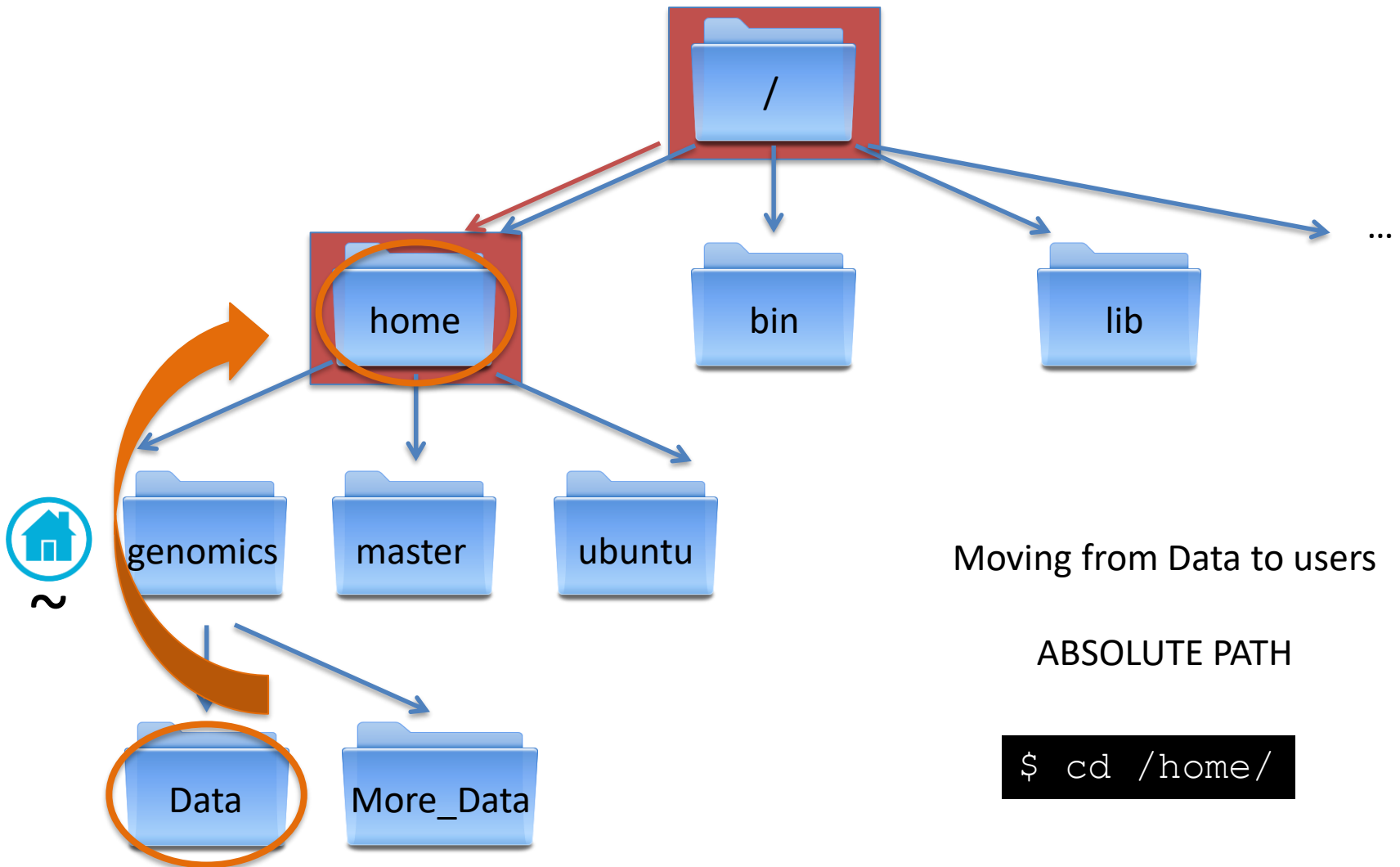
Absolute and Relative Paths



Absolute and Relative Paths



Absolute and Relative Paths



Let's put this to practice

Where am I right now? (Should be the Data directory)

```
$ pwd
```

Change to the directory above

```
$ cd ..
```

Let's list the contents of the Data directory

```
$ ls ./Data
```

Let's put this to practice

Where am I right now? (Should be the Data directory)

```
$ pwd
```

Change to the directory above

```
$ cd ..
```

Let's list the contents of the Data directory

```
$ ls ./Data
```

CHALLENGE 1!

1. Move into the Data directory and list the contents of your home directory
2. In Data, make a new directory and move into this location
3. From this new directory, move into your home directory IN ONE COMMAND and check your location

If You're Typing, You're Doing Something Wrong!

Tab complete is a nice trick to save you typing paths

For this examples we are going to list everything in directory /var/run

Start by typing:

```
$ ls /
```

Followed by tab twice quickly

```
genomics@harvard ami:~$ ls /
bin/          lib/          root/         usr/
boot/         lib64/        run/          var/
dev/          lost+found/   sbin/         vmlinuz
etc/          media/        snap/         vmlinuz.old
home/         mnt/          srv/
initrd.img    opt/          sys/
initrd.img.old  proc/        tmp/
genomics@harvard ami:~$ ls /
```

This shows the contents of the root directory

If You're Typing, You're Doing Something Wrong!

Now type:

```
$ ls /v
```

Followed by tab once. The path to the /var/ directory has filled in.

```
$ ls /var/
```

Now type:

```
$ ls /var/r
```

Followed by tab once. The path to the /var/run/ directory has filled in.

```
$ ls /var/run
```

Tab complete will fill in paths, save you time in typing
and prevent typos!

If You're Typing, You're Doing Something Wrong!

Two more tricks for less typing!

* Represents a special character

For example:

```
$ ls /home/genomics/*.txt
```

Will list everything in my home directory ending .txt

The up arrow can be used to re-run commands

Press your up arrow and see

If you want all of these commands listed, simply type

```
$ history
```

Any Questions So Far?



Binary programs

These are all programs installed on the Unix machine.

They can be found in /bin

```
$ ls /bin
```

```
genomics@harvard_ami:~$ ls /bin/
bash          chmod          hciconfig     mv             pwd            systemd-tty-ask-password-agent
btrfs         chown          hostname     nano           rbash          tailf
btrfs-calc-size  chvt           ip            nc             readlink       tar
btrfsck        cp             journalctl   nc.openbsd    red            tempfile
btrfs-convert  cpio           kbd mode     netcat         rm             touch
btrfs-debug-tree  dash          kill          netstat        rmdir          true
btrfs-find-root  date          kmod          networkctl    rnano          udevadm
btrfs-image      dd            less          nisdomainname run-parts      ulockmgr_server
btrfs-map-logical  df            lessecho     ntfs-3g        sed            umount
btrfs-select-super  dir           lessfile     ntfs-3g.probe setfacl        uname
btrfs-show-super  dmesg         lesskey      ntfs-3g.secaudit setfont        uncompress
btrfstune        dnsdomainname  lesspipe     ntfs-3g.usermap  setupcon      unicode_start
btrfs-zero-log    domainname     ln            ntfsclust     sh             vdir
bunzip2          dumpkeys      loadkeys     ntfscluster   sh.distrib    vmmouse_detect
busybox          echo           login         ntfsncmp       sleep          wdctl
bzip2            ed             loginctl     ntfsfallocate  ss             which
bzip             egrep          lowntfs-3g   ntfsfix        static-sh     whiptail
bzdiff           false          ls            ntfsinfo       stty           ypsdomainname
bzegrep          fgconsole     lsblk        ntfsls         su             zcat
bzexe            fgrep          lsmode       ntfsmove       sync           zcmp
bzfgrep          findmnt       mkdir         ntfstuncate    systemctl     zdiff
bzgrep           fsck.btrfs    mkfs.btrfs   ntfswipe       systemd        zegrep
bzip2            fuser         mknod        open            systemd-ask-password  zfgrep
bzip2recover     fusermount    mktemp       openvt         systemd-escape  zforce
bzless           getfacl       more          pidof          systemd-hwdb    zgrep
bzmore           grep           mount        ping            systemd-inhibit  zless
cat              gunzip        mountpoint   ping6           systemd-machine-id-setup  zmore
chacl            gzexe         mt            plymouth       systemd-notify  znew
chgrp            gzip          mt-gnu       ps              systemd-tmpfiles
```

These include pwd, mkdir, ls ...

Every binary program has a manual

To view the manual page, type man followed by the name of the program

```
$ man <PROGRAMME>
```

Open the manual page for ls

Scroll through (enter) and find the options for: long listing format, human-readable sizes and sort by modification time

Exit the manual page (type q) and give these ls options a go in your Data directory



Every binary program has a manual

To view the manual page, type man followed by the name of the program

```
$ man <PROGRAMME>
```

Open the manual page for ls

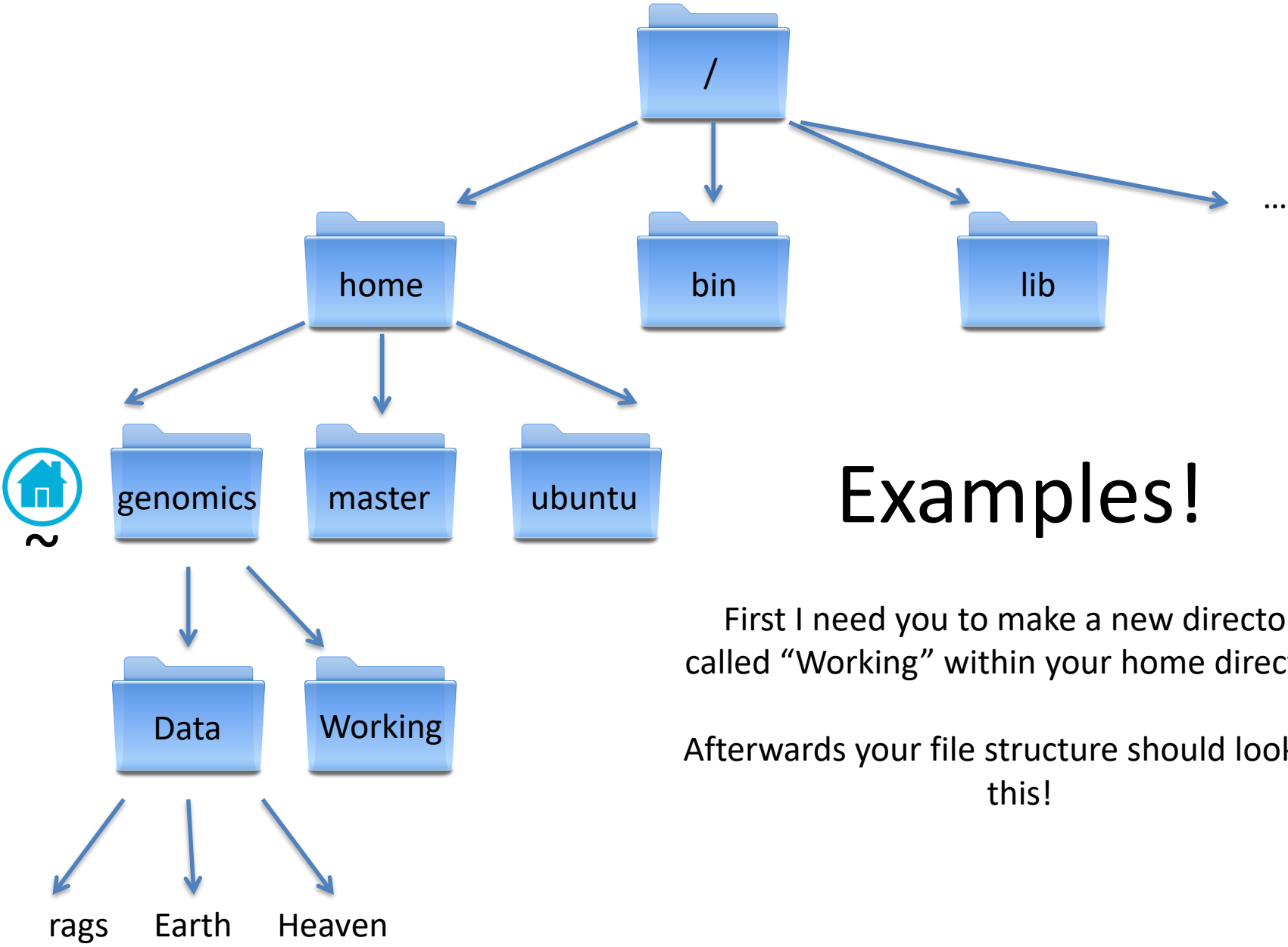
```
$ man ls
```

Scroll through (enter) and find the options for: long listing format (-l), human-readable Sizes (-h) and sort by modification time (-t)

Exit the manual page (type q) and give these ls options a go in your Data directory

```
$ ls -l OR $ ls -h OR $ ls -t
```

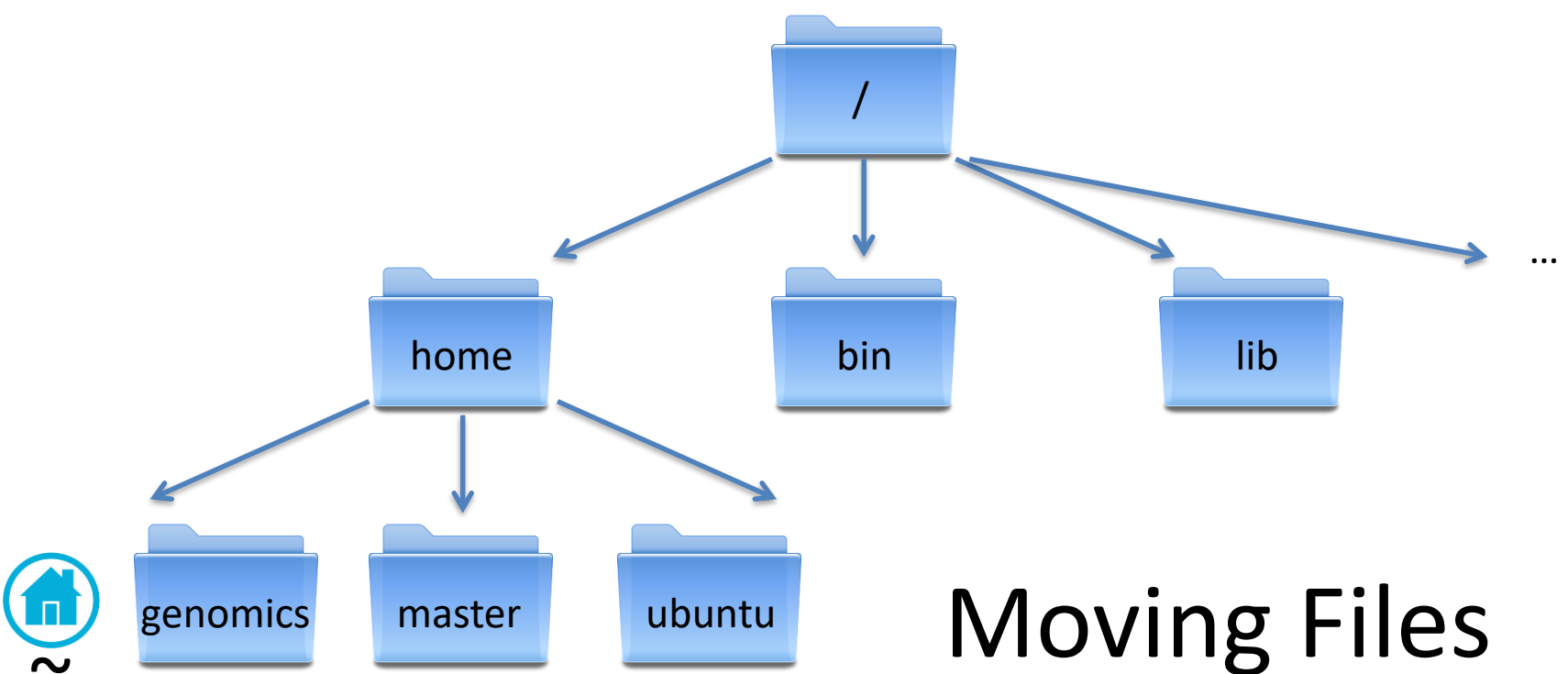




Examples!

First I need you to make a new directory called “Working” within your home directory.

Afterwards your file structure should look like this!



Moving Files

Lets move Heaven and Earth from Data to Working

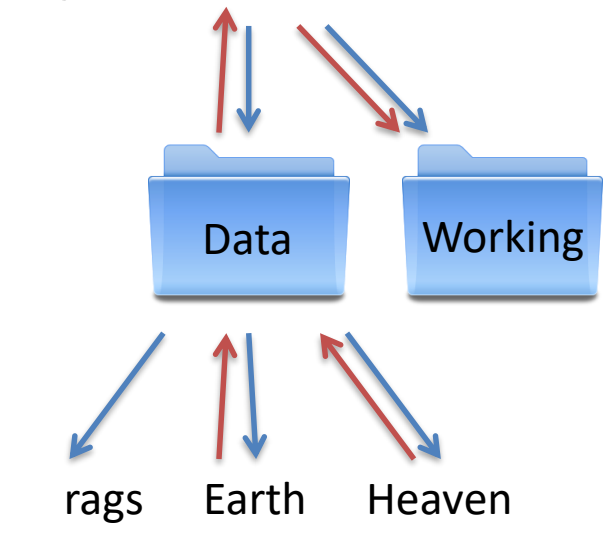
```

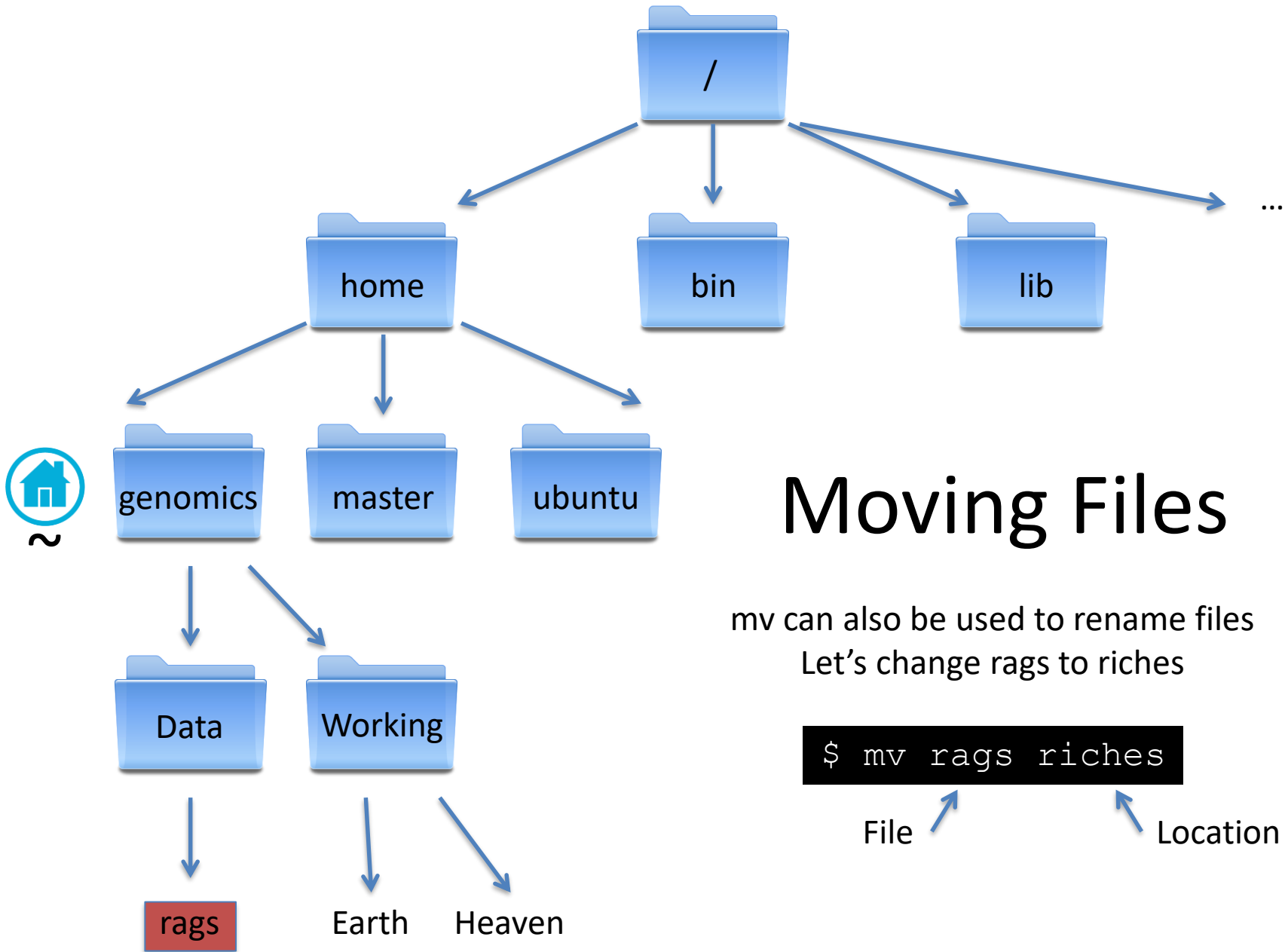
$ cd ~/Data
$ mv Earth ../Working/
  
```

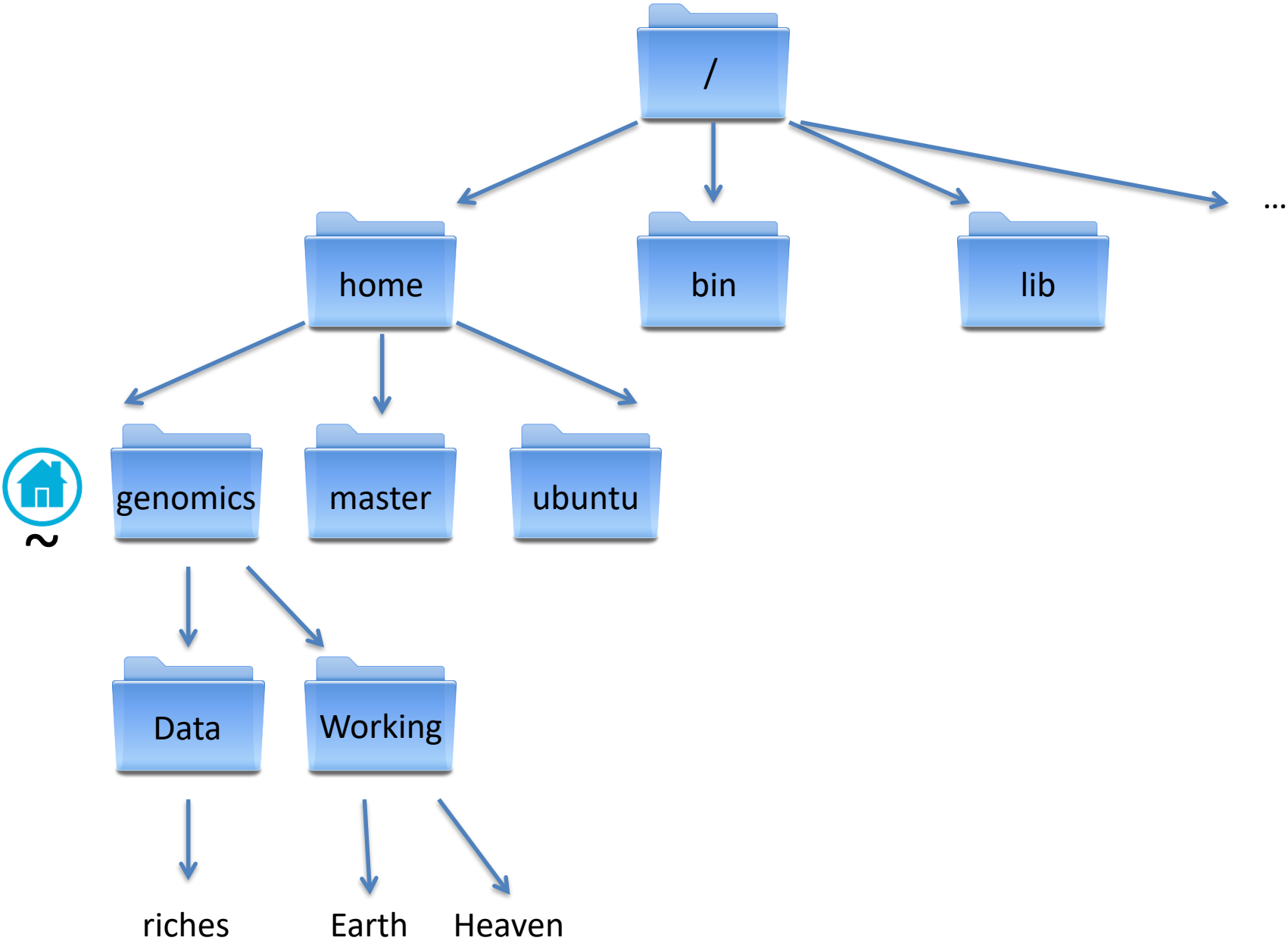
File Location

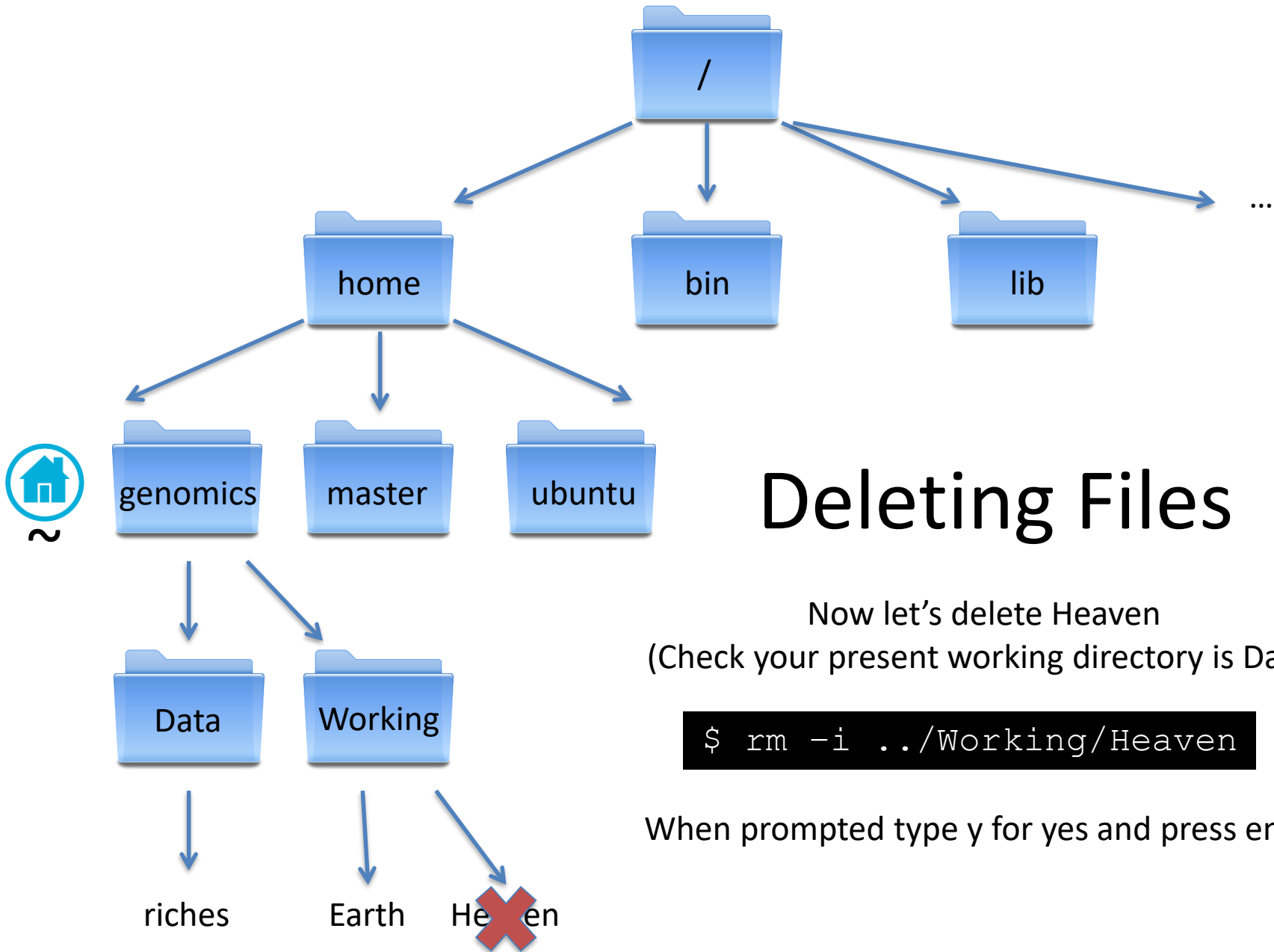
Now move Heaven too

REMEMBER TAB COMPLETE!







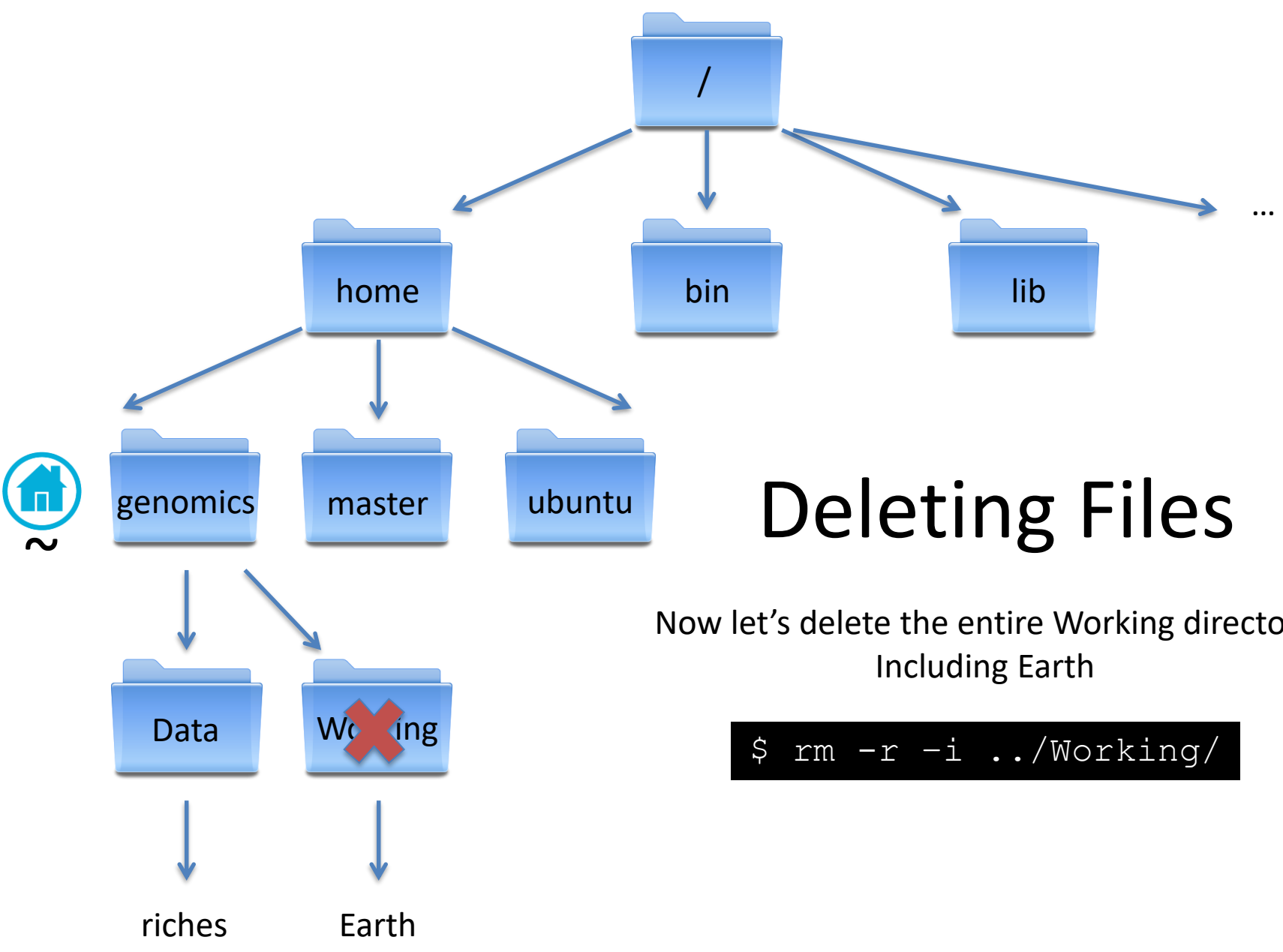


Deleting Files

Now let's delete Heaven
(Check your present working directory is Data)

```
$ rm -i ../Working/Heaven
```

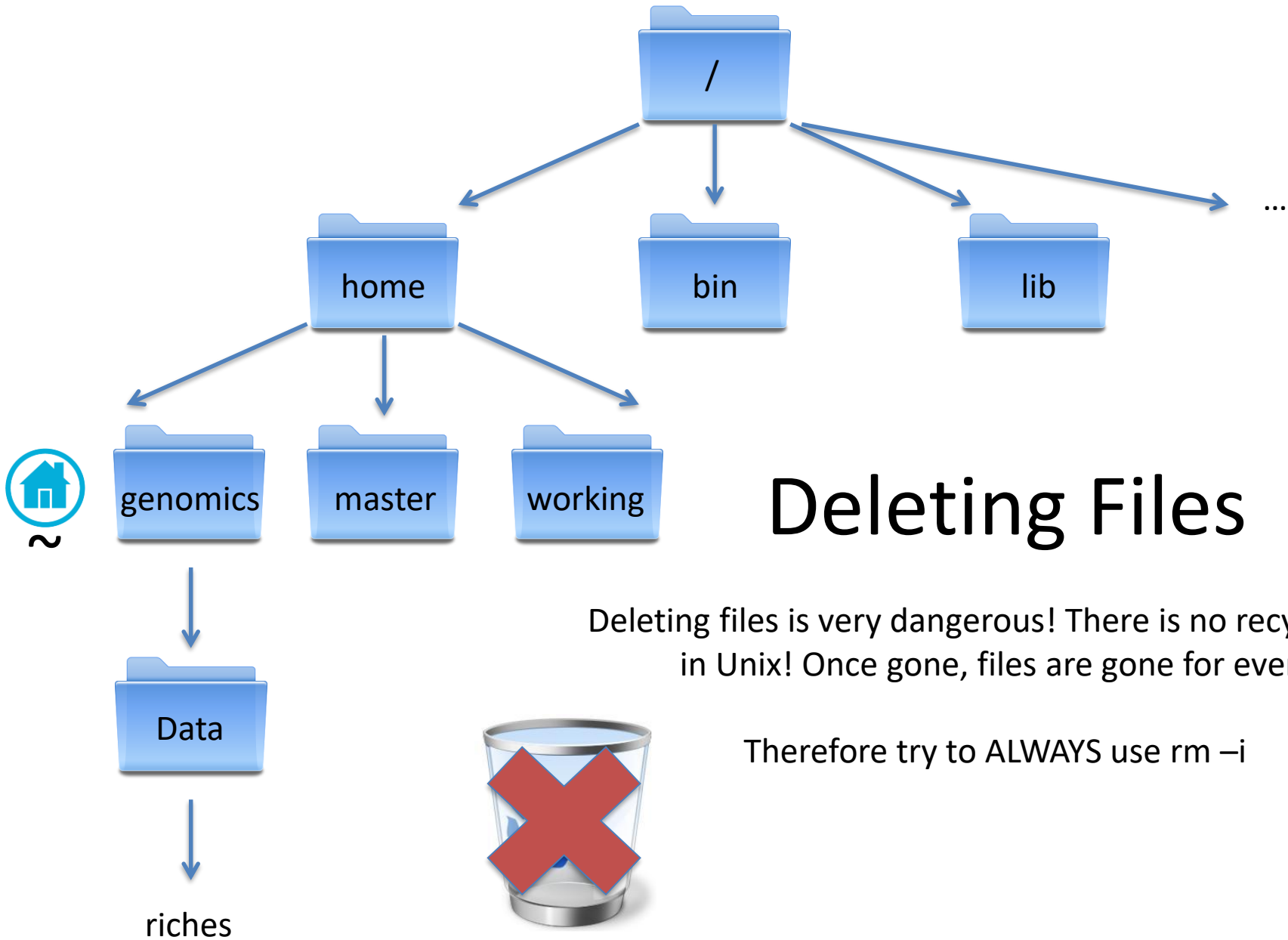
When prompted type y for yes and press enter



Deleting Files

Now let's delete the entire Working directory
Including Earth

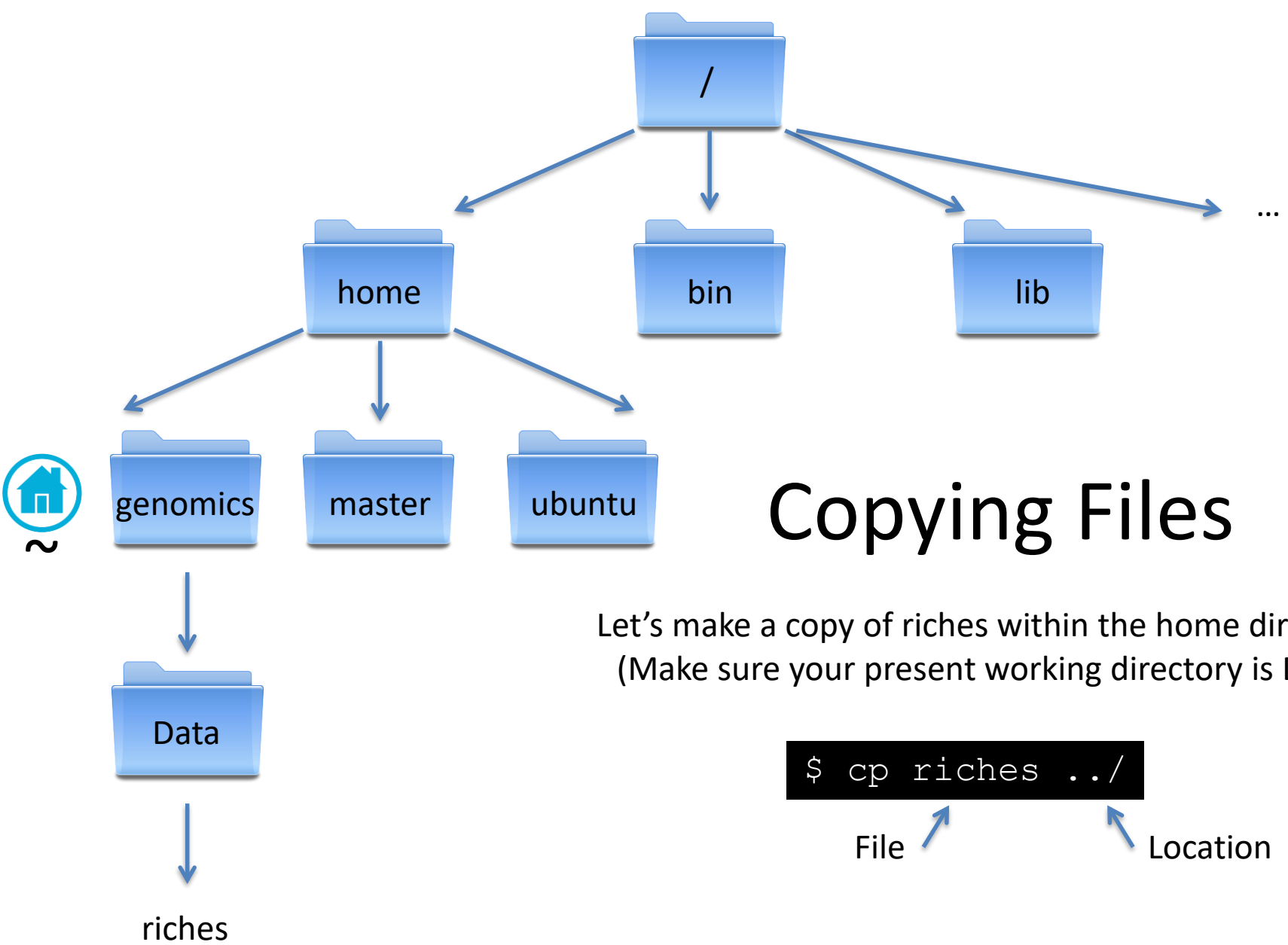
```
$ rm -r -i ../Working/
```



Deleting Files

Deleting files is very dangerous! There is no recycle bin in Unix! Once gone, files are gone for ever!

Therefore try to ALWAYS use `rm -i`

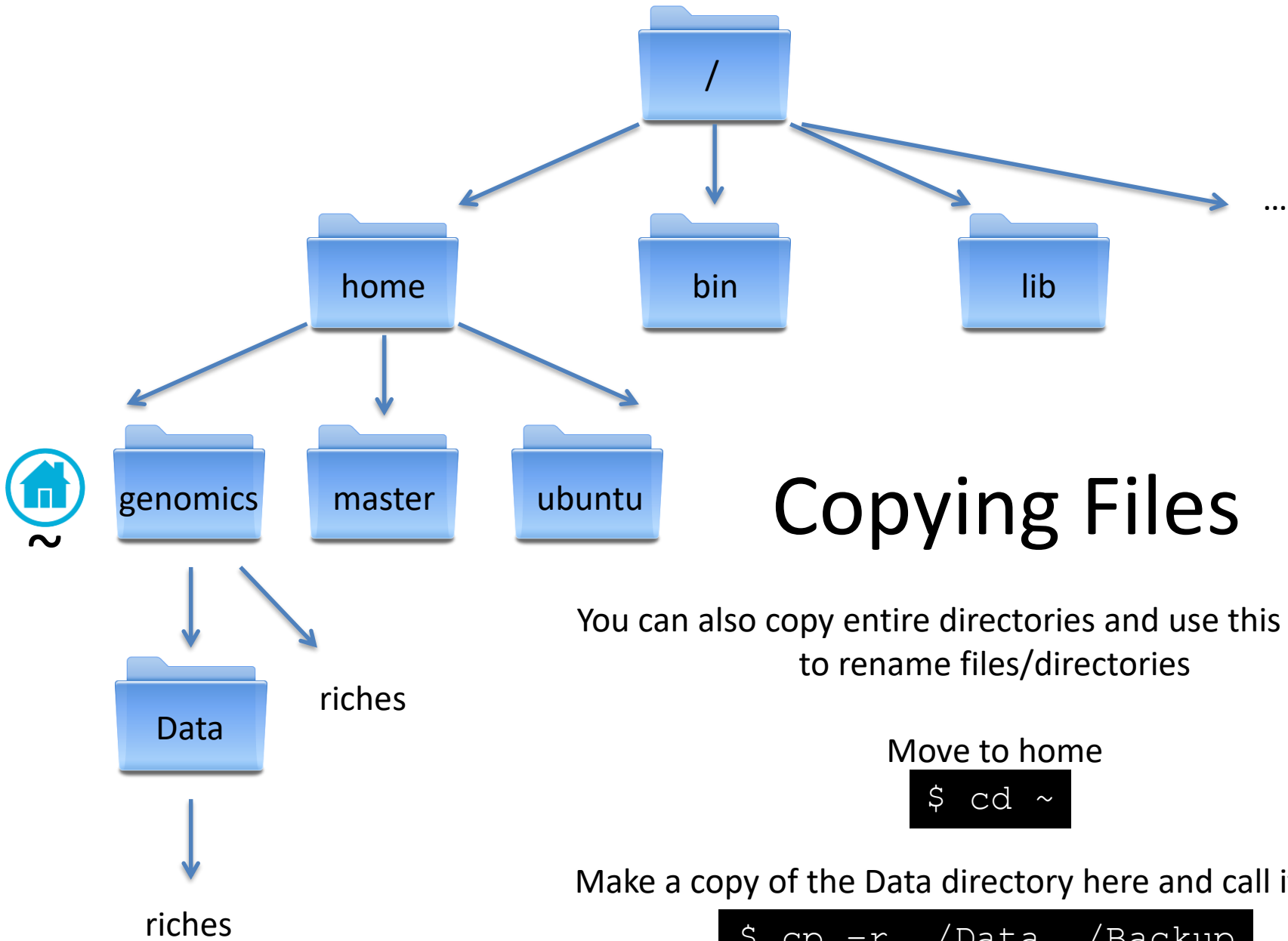


Copying Files

Let's make a copy of riches within the home directory
(Make sure your present working directory is Data)

```
$ cp riches ../
```

File Location



Copying Files

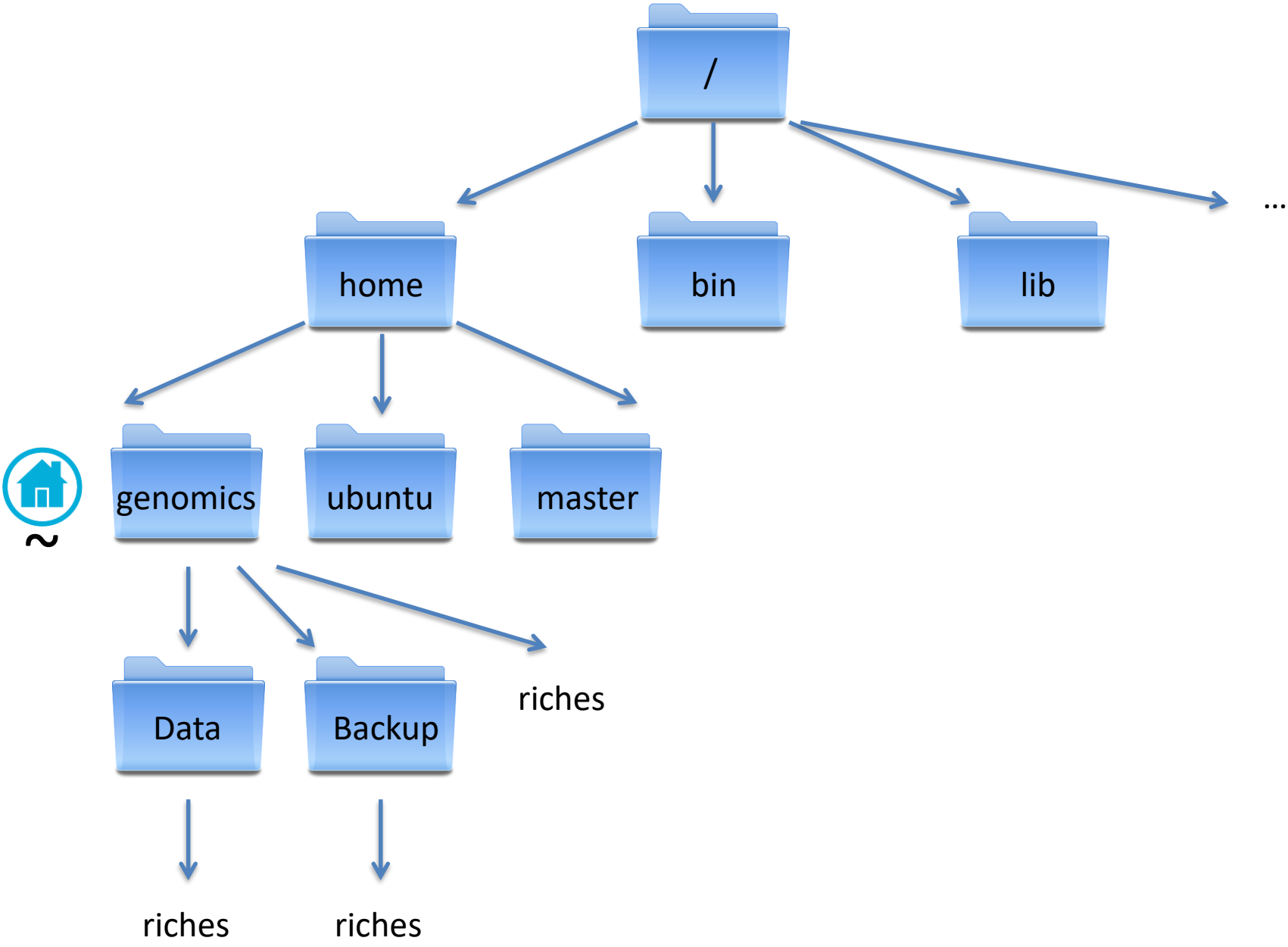
You can also copy entire directories and use this function to rename files/directories

Move to home

```
$ cd ~
```

Make a copy of the Data directory here and call it Backup

```
$ cp -r ./Data ./Backup
```



Typical File Sizes



One Sequencing Sample on the Illumina NextSeq
3,000,000 reads = 1 Gb

But typically you will sequence more than one sample
You may have different patients, different locations, replicates etc...

The size of the sequencing data file can easily become 100s of Gb

(or even bigger depending on the sequencer used)

Archived/Compressed Files

Commonly, people will compress large files so that they are easier to store or share

Here's an example:

sequences.tar.gz

.tar – means that it is a tape archive

.gz – means that it is gzipped

These can be used alone or in combination

To uncompress

A Tar Archive

```
$ tar -xvf <filename>
```

(x = extract, v = verbose, f = all files)

A Gzipped file

```
$ gunzip <filename>
```

A Gzipped Tar archive

```
$ tar -xzvf <filename>
```

Any Questions So Far?



Challenge 2!

1. Change to the `unix_workshop` directory at the following path:

```
$ cd ~/workshop_materials/unix_workshop
```

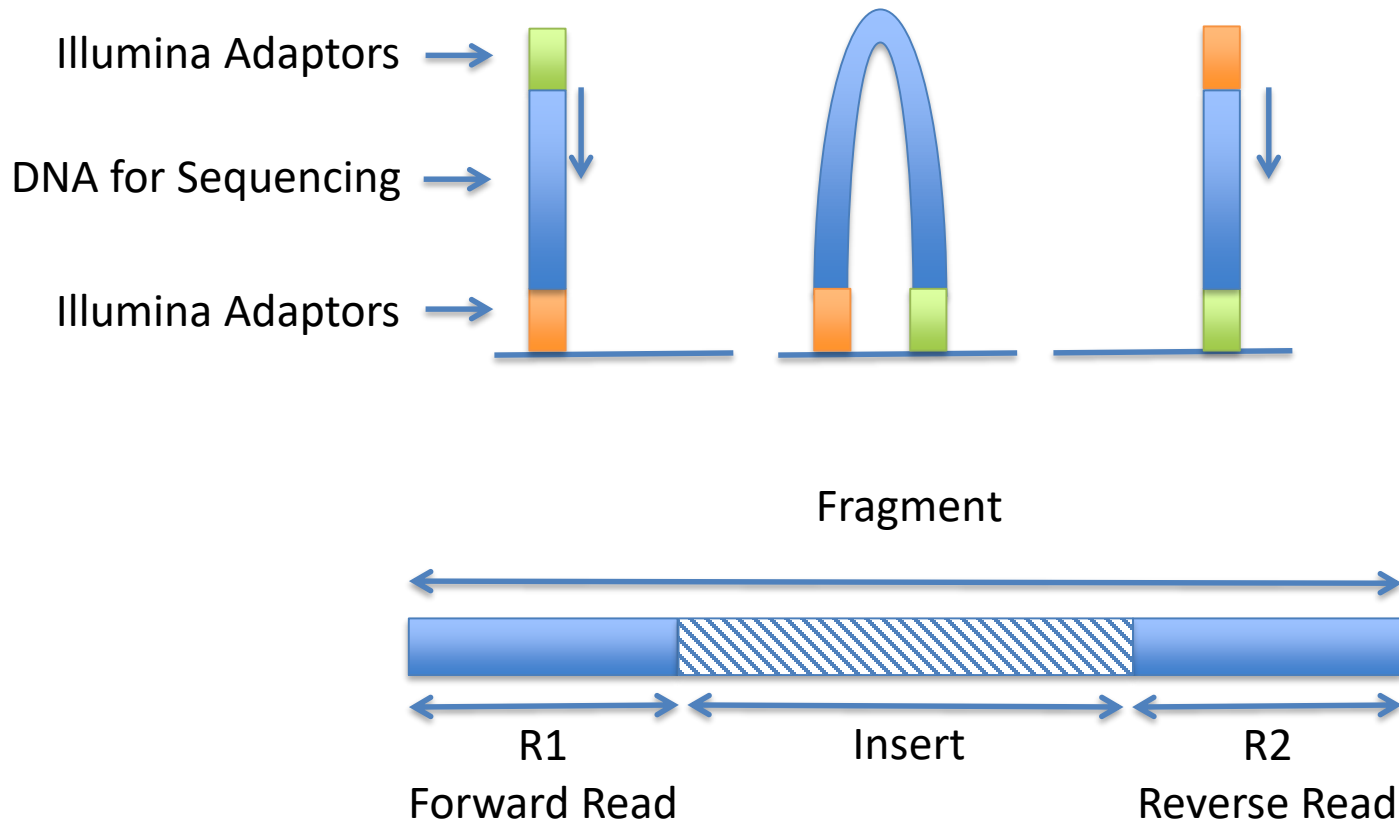
You should find a compressed directory:

```
Sequences.tar
```

2. Make a copy of this file in a Backup directory
3. Un archive the directory
4. Unzip the read files
4. Rename the unarchived files – `sequence_1.fq` and `sequence_2.fq`
5. Delete the original `.tar` file

```
tar          gunzip
cp           mv
rm -i       mkdir
cd
```

Paired Reads



An example: 300 bp paired end reads with a 700 bp fragment size
R1 = 300 bp, R2 = 300 bp, Insert = 100bp

Looking at File Contents

head	tail	more	less	cat
Shows the top lines of a file	Shows the bottom lines of a file	Shows the file one full screen at a time	Shows the file one full screen at a time	Shows an entire file all at once
-n specifies the number of lines (default 10)	-n specifies the number of lines (default 10)	Enter to scroll one line Space to scroll a page q to quit	Enter to scroll one line Space to scroll a page q to quit / to search	Ctrl + C to stop

Use these command line programs to look at the sequence files

Let's put this to use

```
genomics@harvard_ami:~/workshop_materials/unix_workshop$ cd Sequences
genomics@harvard_ami:~/workshop_materials/unix_workshop/Sequences$ head sequence_1.fq
@E.-371320/1
GCTGGTCAACCAGGATAAAACCACCACTGACCCGATGGCGGTTGTTGACTGGATCAACATGTTTGCCTGGCAGTGAACGAAGAGAACGCTGCTGG
CGGTCGCGTGGTACTGCGCCGACTAACGGTGCCTGCGGGATTATCCCGGCAGTCTGCGTACTACGACAAGTTTATCCGCGAAGTGAACGCTAA
CTCACTGGCTCGTTACCTGCTGGTAGCCAGCGCCATTGGTACTCTTTATAAGATGAAC
+
????,BBBBDD<BD?<FGFFGFCFFIIHIIHIGIIDHGIIDIIIIIGIFHHHIIHHHIIIII-HHDIIHIIHIIIFIIHHHFIGIHHH:HGI=G
HHHHFGIBEGHHHHBFH=HHHFEEFGGEGHGDEBFG?FIFEG:8EBEEDG:GEEBGGGGGGG(GEGGGGE?FECGFFCGFDFGFGEFEE??GG?
GCGE*FG?:6E/FGCEEHC:FGF-:G?6GCGGAAGGG6G)EGEC:GGFE'G;GC?G8
```

Fastq File Format:

Header

Sequence

Second Header (often +)

Phred Quality Score

Lot's of analysis software like paired reads to be in the same order

Use head to check that the top three headers are in the same order in
sequence_1.fq and sequence_2.fq

Sequencing Stats

How many reads?

Count the number of lines

```
$ wc -l sequence_1.fq
```

742640 lines THEREFORE 185660 reads

This is the letter l

Are there the same number of reverse reads?

How about just counting the header lines?

Each header line starts @E

```
$ grep -c "@E" sequence_1.fq
```

219153

BUT the numbers from the two programs don't match?!

```
$ grep "@E" sequence_1.fq
```

How about with this

```
$ grep -c "^@E" sequence_1.fq
```

185660

^ matches this pattern at the start of the line – this is an example of a regular expression

Any Questions So Far?



**AND NOW FOR A BRIEF SEGWAY
INTO SCRIPTS...**

Shell Scripts

Imagine you have a complicated command to run. Take this as an example:

```
ref_map.pl -o ./stacks_gsnap/ -T 4 -O ./popmap -B middleton2_  
radtags -b 1 -s ./aligned_gsnap/s13_an_01.bam -s ./aligned_  
gsnap/s13_an_02.bam -s ./aligned_gsnap/s13_an_03.bam -s  
./aligned_gsnap/s13_an_04.bam -s ./aligned_gsnap/s13_an_05.bam  
-s ./aligned_gsnap/s13_an_06.bam -s ./aligned_gsnap/s13_an_07  
.bam -s ./aligned_gsnap/s13_an_08.bam -s ./aligned_gsnap/s13_  
fw_01.bam -s ./aligned_gsnap/s13_fw_02.bam -s ./aligned_gsnap/  
s13_fw_03.bam -s ./aligned_gsnap/s13_fw_04.bam -s ./aligned_  
gsnap/s13_fw_05.bam -s ./aligned_gsnap/s13_fw_06.bam -s  
./aligned_gsnap/s13_fw_07.bam -s ./aligned_gsnap/s13_fw_08.bam
```

But what if you make a mistake?

Or want to run this command 10 times?

You have to type it out every time ☹️

Shell Scripts

Instead we can put this command inside a script.

Then it can easily be edited and ran multiple times

To understand shell scripts, we're going to look at a few topics:

- Shell scripting languages
 - Text editors
- How to write a script
- How to run a script

Scripts make a great record of what you've done, when and with what.

You should also aim to keep a computational biology lab book.

What is a Shell Script?

A computer program designed to be run by the Unix shell, the command line interpreter.

There are various types of shell scripts. These are scripting languages.

Today we are going to look at bash

First, let's run a simple bash command:

```
$ echo Hello World
```

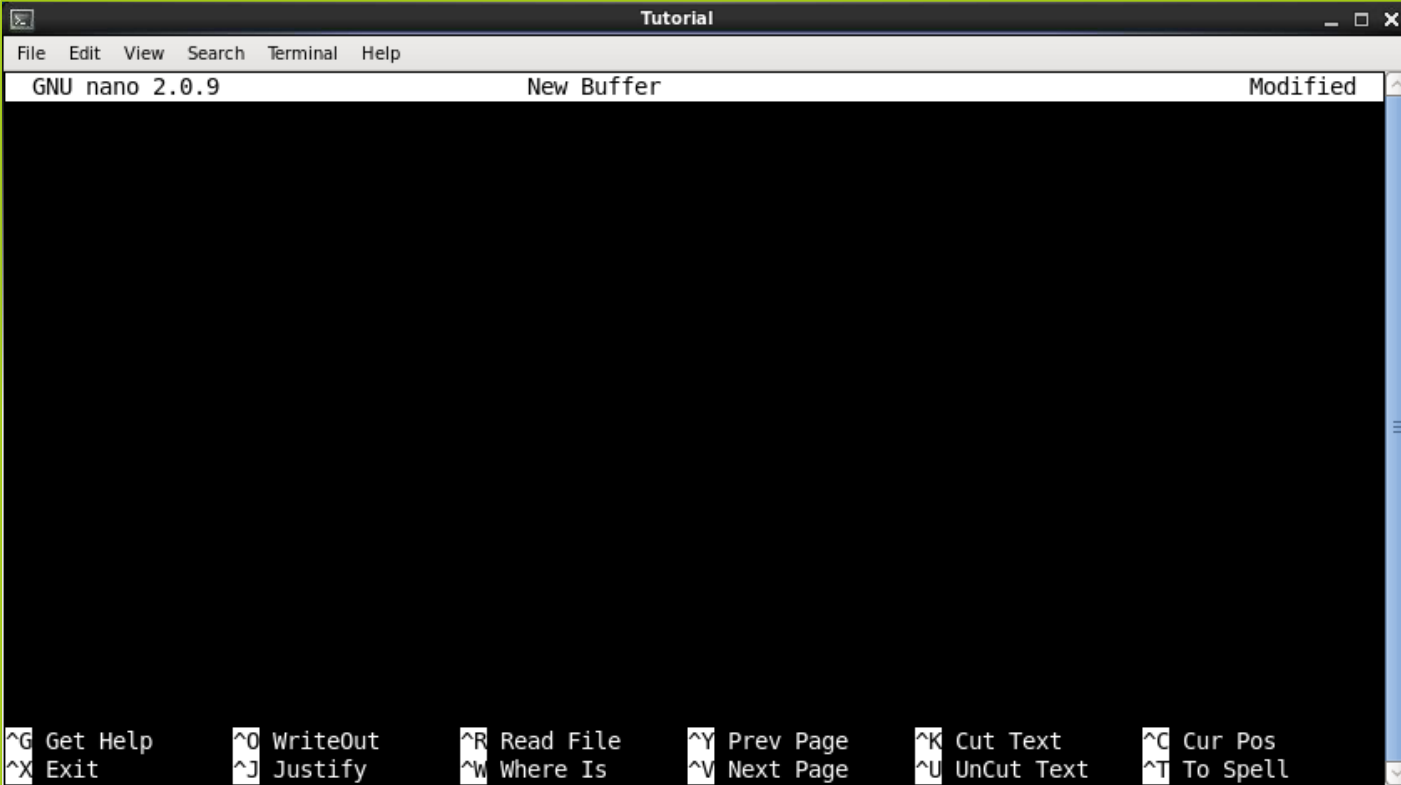
```
$ echo Hello World  
Hello World
```

Try using echo with a different phrase

Your First Script

Let's start by opening nano

```
$ nano
```



The screenshot shows the GNU nano 2.0.9 terminal editor interface. The window title is "Tutorial". The menu bar includes "File", "Edit", "View", "Search", "Terminal", and "Help". The status bar at the top shows "GNU nano 2.0.9", "New Buffer", and "Modified". The main editing area is black. The bottom status bar displays various keyboard shortcuts: `^G` Get Help, `^O` WriteOut, `^R` Read File, `^Y` Prev Page, `^K` Cut Text, `^C` Cur Pos; `^X` Exit, `^J` Justify, `^W` Where Is, `^V` Next Page, `^U` UnCut Text, `^T` To Spell.

Key Nano Commands

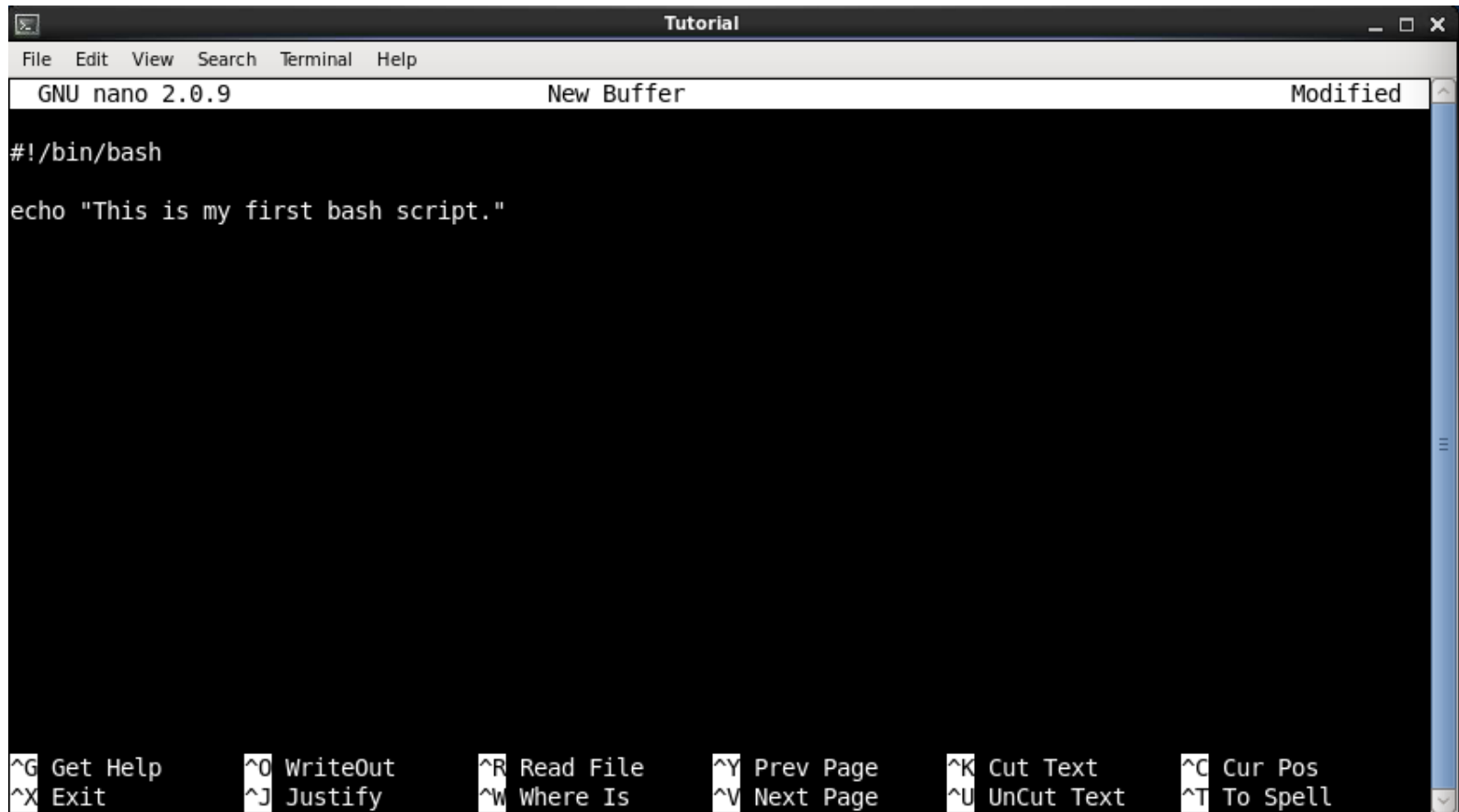
Ctrl + O – This saves the file. You will be asked for a file name.
Type the name and press enter.



Ctrl + X – This exits nano. If the file is unsaved, you will be asked at this point if you'd like to save it.



Your First Script

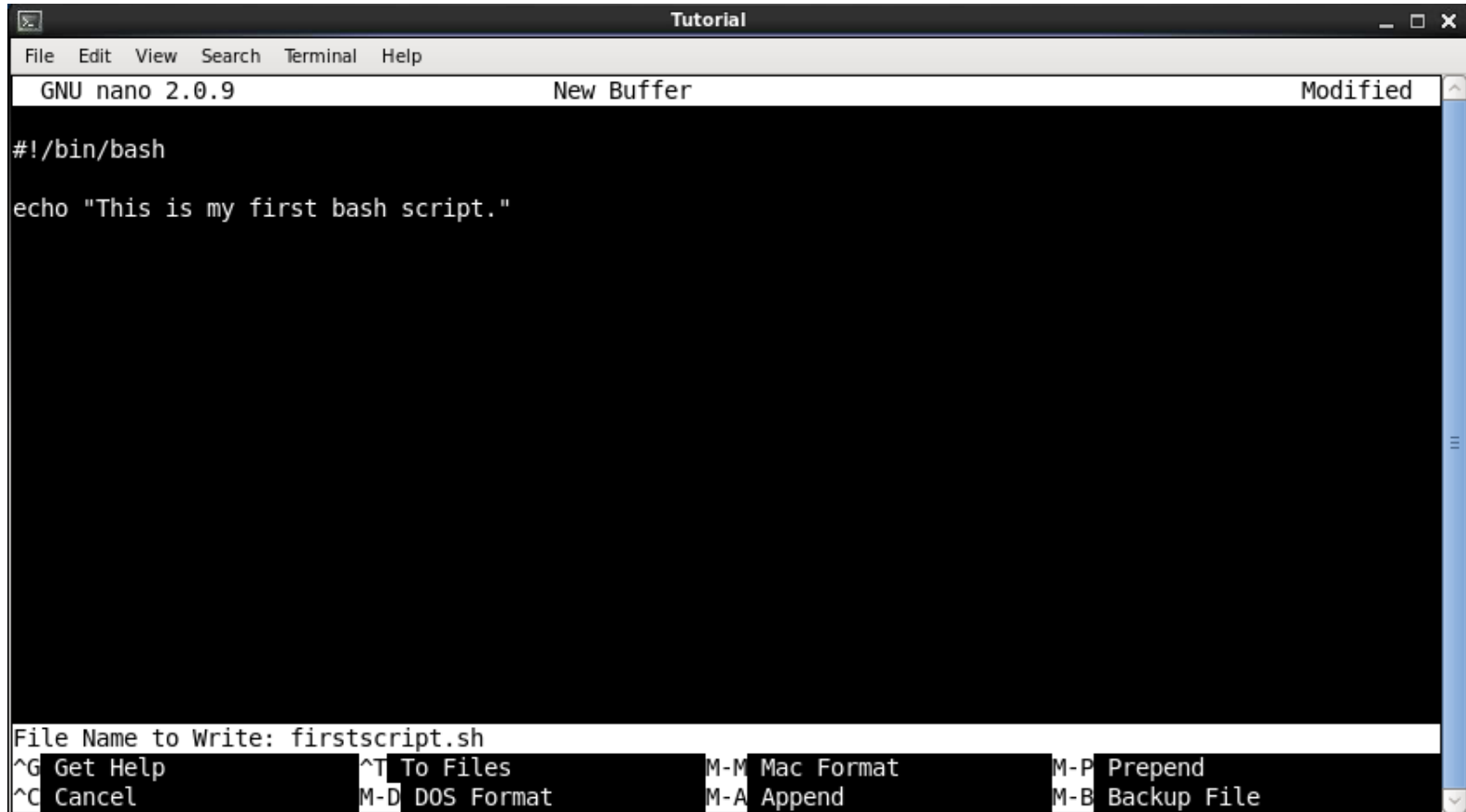


```
GNU nano 2.0.9          New Buffer          Modified
#!/bin/bash
echo "This is my first bash script."

^G Get Help      ^O WriteOut     ^R Read File    ^Y Prev Page    ^K Cut Text     ^C Cur Pos
^X Exit          ^J Justify      ^W Where Is     ^V Next Page    ^U UnCut Text   ^T To Spell
```

`#!/bin/bash` tells the computer that this script is in the language bash. It always needs to go at the top of any bash script.

Your First Script



The image shows a terminal window titled "Tutorial" with a menu bar containing "File", "Edit", "View", "Search", "Terminal", and "Help". The editor is GNU nano 2.0.9, editing a file named "New Buffer". The content of the file is a simple bash script:

```
#!/bin/bash
echo "This is my first bash script."
```

At the bottom of the window, the "File Name to Write: firstscript.sh" is displayed. Below this, a list of keyboard shortcuts is shown:

^G	Get Help	^T	To Files	M-M	Mac Format	M-P	Prepend
^C	Cancel	M-D	DOS Format	M-A	Append	M-B	Backup File

Then use Ctrl + O to save and give the file the name firstscript.sh.

Then use Ctrl + X to exit.

Now Run Your Script

Simply Type:

```
$ bash firstscript.sh
```

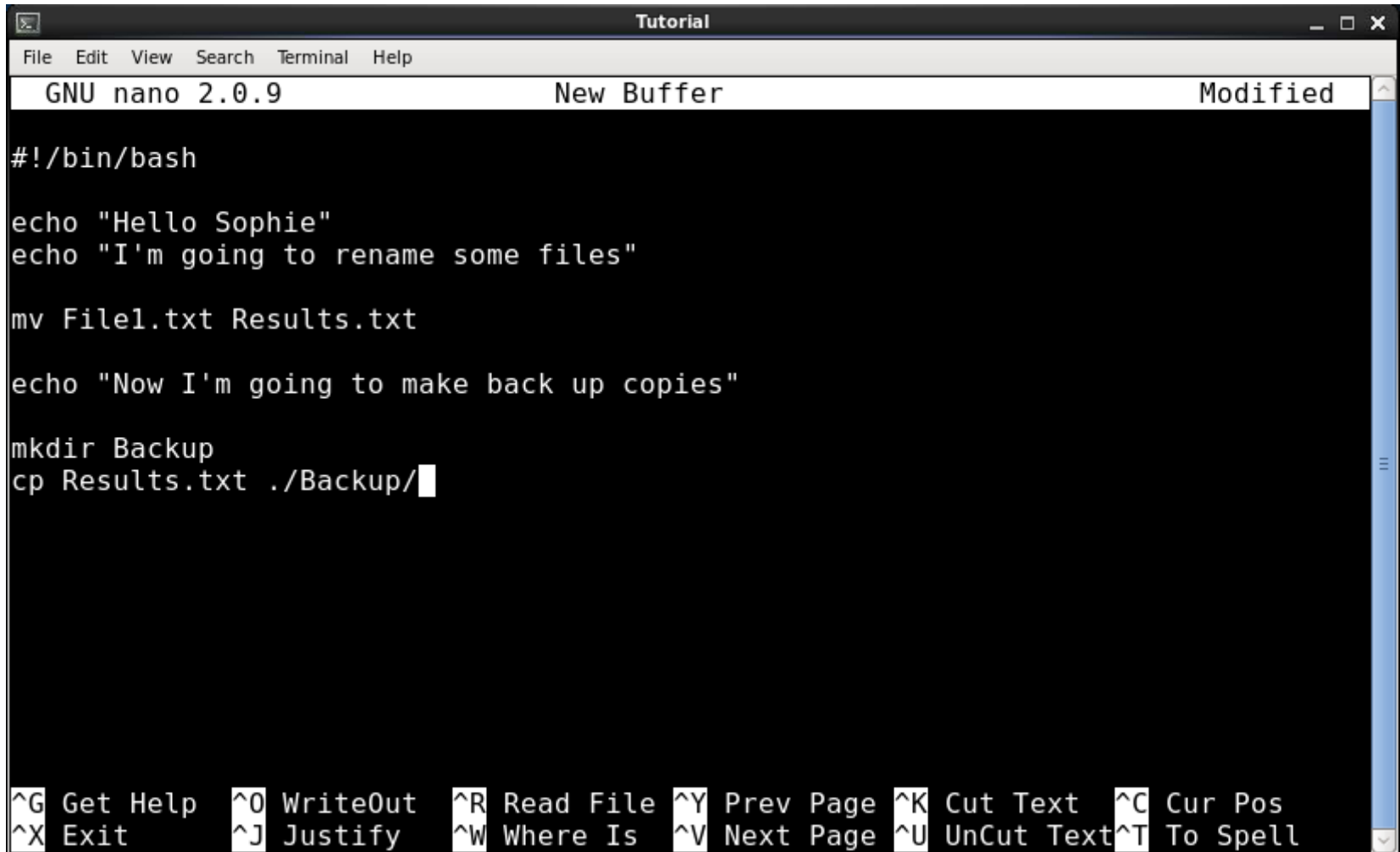
Reopen the same script:

```
$ nano firstscript.sh
```

Change the phrase, save it and run the script again

Bash Scripts

Bash scripts can be used to run binary programs like cd, mv, cp etc...



The image shows a screenshot of a terminal window titled "Tutorial" running the GNU nano 2.0.9 text editor. The editor is editing a new buffer. The script content is as follows:

```
#!/bin/bash

echo "Hello Sophie"
echo "I'm going to rename some files"

mv File1.txt Results.txt

echo "Now I'm going to make back up copies"

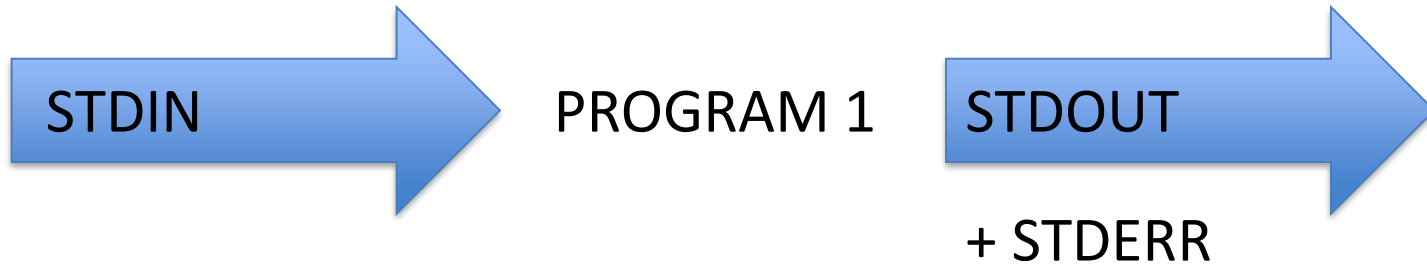
mkdir Backup
cp Results.txt ./Backup/
```

At the bottom of the window, there is a help menu with the following shortcuts:

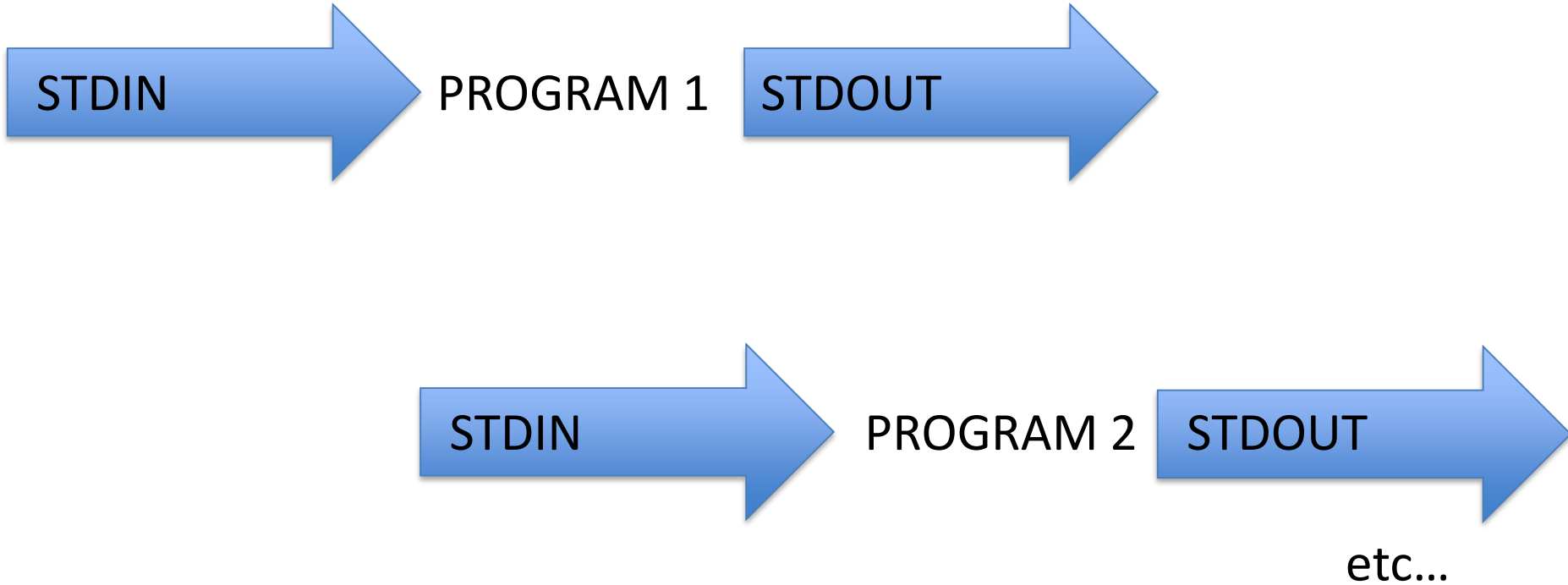
^G Get Help	^O WriteOut	^R Read File	^Y Prev Page	^K Cut Text	^C Cur Pos
^X Exit	^J Justify	^W Where Is	^V Next Page	^U UnCut Text	^T To Spell

PIPELINES (TIME DEPENDING)

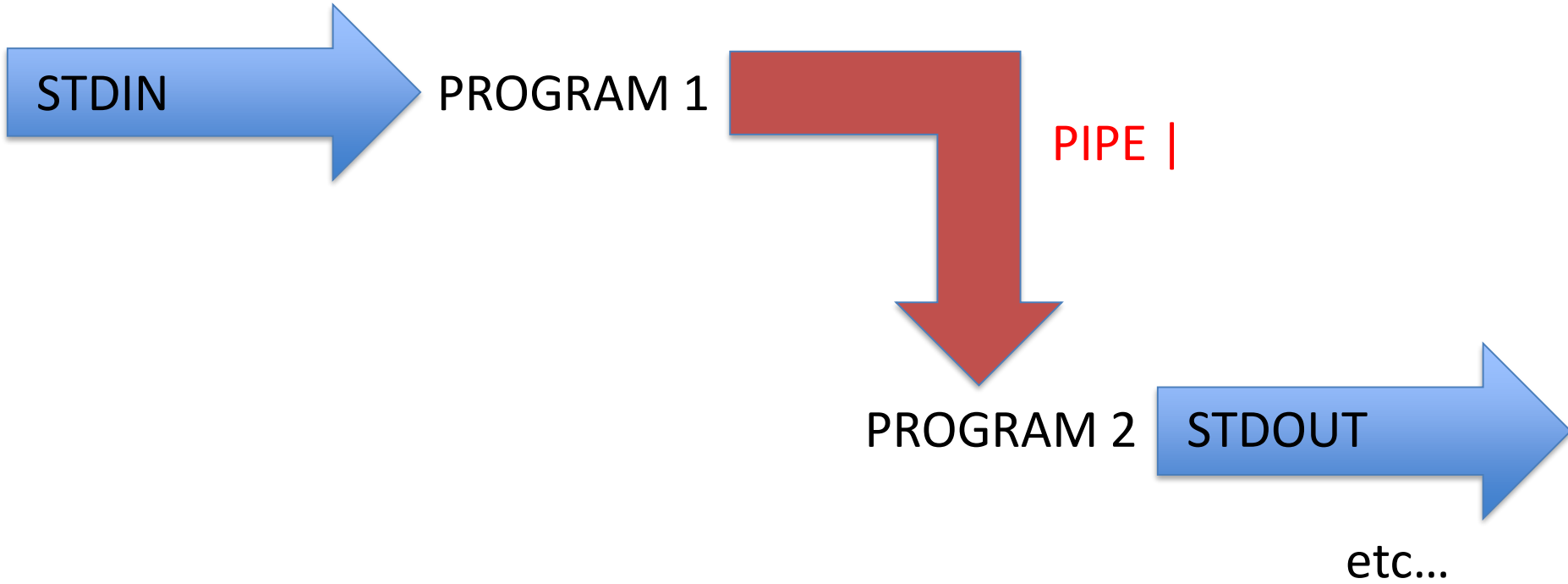
Pipelines



Pipelines



Pipelines



Let's put this to practice: Building Pipelines

Count the number of files and folders in your home directory

Let's build the first part of the pipeline, listing the files:

Number 1

```
$ ls -l /home/genomics/
```



PIPE this into `wc -l` to count the number of lines:
(i.e. the number of files and folders)

```
$ ls -l /home/genomics/ | wc -l
```



Letter l

Let's put this to practice: Building Pipelines

How many base pairs in first sequence?

Firstly let's get the top two lines of the sequence file:

```
$ head -n 2 sequence_1.fq
```

Now let's PIPE this into tail to get just the sequence line

```
$ head -n 2 sequence_1.fq | tail -n 1
```

Finally PIPE this into word count of characters to count the base pairs

```
$ head -n 2 sequence_1.fq | tail -n 1 | wc -c
```

Is the first reverse read the same length?

Some More Examples

Within the Unix Workshop directory you should find a file called scientists.txt

```
$ cd ~/workshop_materials/unix_workshop/
```

Take a look at the contents

```
$ more scientists.txt
```

First	Last	DOB
Charles	Darwin	12 February 1809
Marie	Curie	07 November 1867
Stephen	Hawking	08 January 1942
Rosalind	Franklin	25 July 1920
Isaac	Newton	04 January 1643
Richard	Dawkins	26 March 1941

Some More Examples

```
$ cat scientists.txt | cut -f 1
```

First
Charles
Marie
Stephen
Rosalind
Isaac
Richard

Now take a look at the original file

```
$ more scientists.txt
```

Some More Examples

```
$ cat scientists.txt | cut -f 1,3
```

First	DOB
Charles	12 February 1809
Marie	07 November 1867
Stephen	08 January 1942
Rosalind	25 July 1920
Issac	04 January 1643
Richard	26 March 1941

Some More Examples

```
$ cat scientists.txt | cut -f 1 | sort
```

Charles
First
Isaac
Marie
Richard
Rosalind
Stephen

What if you wanted to keep the sorted list?

```
$ cat scientists.txt | cut -f 1 | sort > newfile.txt
```



Redirects

CHALLENGE 3!

Looking at the *Saccharomyces cerevisiae* gff3 file

GFF = general feature format

This is a file which lists all of the genome features, their coordinates, and info about them (genes, tRNAs, exons etc...)

```
##gff-version 3
##sequence-region I 1 230218
##sequence-region II 1 813184
##sequence-region III 1 316620
##sequence-region IV 1 1531933
##sequence-region IX 1 439888
##sequence-region Mito 1 85779
##sequence-region V 1 576874
##sequence-region VI 1 270161
##sequence-region VII 1 1090940
##sequence-region VIII 1 562643
##sequence-region X 1 745751
##sequence-region XI 1 666816
##sequence-region XII 1 1078177
##sequence-region XIII 1 924431
##sequence-region XIV 1 784333
##sequence-region XV 1 1091291
##sequence-region XVI 1 948066
#genome-build SGD R64-1-1
#genome-version R64-1-1
#genome-date 2011-09
#genome-build-accession GCA_000146045.2
#genbuild-last-updated 2011-12
I SGD chromosome 1 230218 . . . ID=chromosome:I;Alias=BK006935.2
###
I ensembl gene 335 649 . + . ID=gene:YAL069W;biotype=protein_coding;description=Dubious open reading frame%3B unlikely to encode a functional protein%2C based on available experimental and comparative sequence data [Source:SGD%3BAcc:S000002143];gene_id=YAL069W;logic_name=sgd
I ensembl transcript 335 649 . + . ID=transcript:YAL069W;Parent=gene:YAL069W;Name=YAL069W;biotype=protein_coding;transcript_id=YAL069W
I ensembl exon 335 649 . + . Parent=transcript:YAL069W;Name=YAL069W.1;constitutive=1;ensembl_end_phase=0;ensembl_phase=0;exon_id=YAL069W.1;rank=1
I ensembl CDS 335 649 . + 0 ID=CDS:YAL069W;Parent=transcript:YAL069W;protein_id=YAL069W
###
I ensembl gene 538 792 . + . ID=gene:YAL068W-A;biotype=protein_coding;description=Dubious open reading frame%3B unlikely to encode a functional protein%2C based on available experimental and comparative sequence data%3B identified by gene-trapping%2C microarray-based expression analysis%2C and genome-wide homology searching [Source:SGD%3BAcc:S000028594];gene_id=YAL068W-A;logic_name=sgd
I ensembl transcript 538 792 . + . ID=transcript:YAL068W-A;Parent=gene:YAL068W-A;Name=YAL068W-A;biotype=protein_coding;transcript_id=YAL068W-A
I ensembl exon 538 792 . + . Parent=transcript:YAL068W-A;Name=YAL068W-A.1;constitutive=1;ensembl_end_phase=0;ensembl_phase=0;exon_id=YAL068W-A.1;rank=1
I ensembl CDS 538 792 . + 0 ID=CDS:YAL068W-A;Parent=transcript:YAL068W-A;protein_id=YAL068W-A
###
I ensembl gene 1807 2169 . - . ID=gene:YAL068C;Name=PAU8;biotype=protein_coding;description=Protein of unknown function%3B member of the seripauperin multigene family encoded mainly in subtelomeric regions [Source:SGD%3BAcc:S000002142];gene_id=YAL068C;logic_name=sgd
I ensembl transcript 1807 2169 . - . ID=transcript:YAL068C;Parent=gene:YAL068C;Name=PAU8;biotype=protein_coding;transcript_id=YAL068C
I ensembl exon 1807 2169 . - . Parent=transcript:YAL068C;Name=YAL068C.1;constitutive=1;ensembl_end_phase=0;ensembl_phase=0;exon_id=YAL068C.1;rank=1
I ensembl CDS 1807 2169 . - 0 ID=CDS:YAL068C;Parent=transcript:YAL068C;protein_id=YAL068C
###
--More-- (0%)
```

CHALLENGE 3!

```
##gff-version 3
##sequence-region I 1 230218
##sequence-region II 1 813184
##sequence-region III 1 316620
##sequence-region IV 1 1531933
##sequence-region IX 1 439888
##sequence-region Mito 1 85779
##sequence-region V 1 576874
##sequence-region VI 1 270161
##sequence-region VII 1 1090940
##sequence-region VIII 1 562643
##sequence-region X 1 745751
##sequence-region XI 1 666816
##sequence-region XII 1 1078177
##sequence-region XIII 1 924431
##sequence-region XIV 1 784333
##sequence-region XV 1 1091291
##sequence-region XVI 1 948066
##genome-build SGD R64-1-1
##genome-version R64-1-1
##genome-date 2011-09
##genome-build-accession GCA_000146045.2
##genebuild-last-updated 2011-12
I SGD chromosome 1 230218 . . ID=chromosome:I;Alias=BK006935.2
###
I ensembl gene 335 649 . + . ID=gene:YAL069W;biotype=protein_coding;description=Dubious open reading frame%3B unlikely to encode a functional protein%2C based on available experimen
tal and comparative sequence data [Source:SGD%3BAcc:S000002143];gene_id=YAL069W;logic_name=sgd
I ensembl transcript 335 649 . + . ID=transcript:YAL069W;Parent=gene:YAL069W;Name=YAL069W;biotype=protein_coding;transcript_id=YAL069W
I ensembl exon 335 649 . + . Parent=transcript:YAL069W;Name=YAL069W.1;constitutive=1;ensembl_end_phase=0;ensembl_phase=0;exon_id=YAL069W.1;rank=1
I ensembl CDS 335 649 . + 0 ID=CDS:YAL069W;Parent=transcript:YAL069W;protein_id=YAL069W
###
I ensembl gene 538 792 . + . ID=gene:YAL068W-A;biotype=protein_coding;description=Dubious open reading frame%3B unlikely to encode a functional protein%2C based on available experimen
tal and comparative sequence data%3B identified by gene-trapping%2C microarray-based expression analysis%2C and genome-wide homology searching [Source:SGD%3BAcc:S000028594];gene_id=YAL068W-A;logic_name=sgd
I ensembl transcript 538 792 . + . ID=transcript:YAL068W-A;Parent=gene:YAL068W-A;Name=YAL068W-A;biotype=protein_coding;transcript_id=YAL068W-A
I ensembl exon 538 792 . + . Parent=transcript:YAL068W-A;Name=YAL068W-A.1;constitutive=1;ensembl_end_phase=0;ensembl_phase=0;exon_id=YAL068W-A.1;rank=1
I ensembl CDS 538 792 . + 0 ID=CDS:YAL068W-A;Parent=transcript:YAL068W-A;protein_id=YAL068W-A
###
I ensembl gene 1807 2169 . - . ID=gene:YAL068C;Name=PAU8;biotype=protein_coding;description=Protein of unknown function%3B member of the seripauperin multigene family encoded mainly in
subtelomeric regions [Source:SGD%3BAcc:S000002142];gene_id=YAL068C;logic_name=sgd
I ensembl transcript 1807 2169 . - . ID=transcript:YAL068C;Parent=gene:YAL068C;Name=PAU8;biotype=protein_coding;transcript_id=YAL068C
I ensembl exon 1807 2169 . - . Parent=transcript:YAL068C;Name=YAL068C.1;constitutive=1;ensembl_end_phase=0;ensembl_phase=0;exon_id=YAL068C.1;rank=1
I ensembl CDS 1807 2169 . - 0 ID=CDS:YAL068C;Parent=transcript:YAL068C;protein_id=YAL068C
###
--More--(8X)
```

Lines that start # are comments – just run information

Column 1 = Chromosome

Column 3 = Feature/Type e.g. gene, chromosome, exon

Column 4 = Start Location

Column 5 = Stop Location

CHALLENGE 3!

1. In the Unix workshop directory you should find a gff3 file.

```
$ cd ~/Unix_Workshop/Challenge5
```

```
Saccharomyces_cerevisiae.R64-1-1.85.gff3.gz
```

2. Unzip the file.

3. How many feature entries are there?

4. List and count all the different types of features

5. Which chromosome is the longest?

Hints!

- Use head to work with 10 lines whilst testing what your pipe does!
 - This is a tab delimited file with a column layout.
 - Google “gff3 format” to find out what each of the columns are.
 - Remember that cat opens an entire file at once.
- There are a number of info lines at the start which begin with a hash. Look into grep with invert matches to skip these.
- Cut can be used to isolate certain columns. You’ll want the field option.
 - The programs sort and uniq may be helpful.
 - Sort must be used before uniq.
 - Uniq has a counting option.
 - Sort uses the key option to sort by a column.

more	gunzip	head	uniq
cp	mv	grep	wc
rm -i	mkdir	cut	
cd	cat	sort	man

Any Questions So Far?

