Lab introduction

Gemma Martínez Redondo & Karin Steffen

Many slides taken from Dr. Sophie Shaw



Daily workshop material, slides & exercises http://evomics.org/2024-workshop-on-phylogenomics-cesky-krumlov/

\leftarrow \rightarrow C \bigcirc $\stackrel{\sim}{\sim}$ evomics.org		☆	Q Search			${igardown}$	ABP 📘	ш	U.	f 1	»	=
EVOLUTION AND GEN Intensive and immersive training opportunitie	IOMICS es	WOR	RKSHOPS	LEARNING	PEOPLE	AF	PPLY	INFO	RMATIC	N		
	and the second second	Work	kshop on Gen	omics							-	
State of the second sec		Workshop on Population and Speciation Genomics										
	2024 Workshop on Phylogenomics, Cesky Krumlov	Work	kshop on Phyl	ogenomics								
	2019 Workshop on Phylogenomics, Cesky Krumlov	Work	kshop on Mol	ecular Evolution							and the second s	
	2017 Workshop on Phylogenomics, Cesky Krumlov	Harv	ard University	y Workshops								
STATE NORTH	MERICAN PERSONNE	Work	kshop on Micr	robiome and Tran	scriptome Ar	nalysis,	Durban,	South A	frica	The second		
Immersive training oppor	tunities in modern computational	Adva	anced Topics							'S		
A A A A A A A A A A A A A A A A A A A		WS &	ANNOUN	CEMENTS	- D	A	1		Ormania D		t	

Daily workshop material, slides & exercises

EVOLUTION AND GENOMICS Intensive and immersive training opportunities	WORKSHOPS	LEARNING	PEOPLE	APPLY	INFORMATIO	N
Important information:						
Group photo!						
Online booking of your free time activities! (also available at the Info-Centre)						
FIND YOUR WAY AROUND ČESKÝ KRUMLOV						
Český Krumlov town information						
• Map of the 2024 changes in Český Krumlov due to the bridge reparation – Check before arrival to choose the bus stop closest to your accomodation!!						
 Map of the main locations in the schedule (Save to your Google Maps!) 						
 Map of some of the main Český Krumlov restaurants (Save to your Google Maps!) 						
 Map of some of the main Český Krumlov bars (Save to your Google Maps!) 						
 Map of Český Krumlov service points: shops, pharmacy, ATM, hospital (Save to your Google Maps!) 						
 Map of the historic City Centre of Český Krumlov 						
FACULTY AND GENOMICS						
 Click here to see our Faculty arrival and departure dates 						
 Faculty Bios and Teaching Assistant Bios 						
Faculty Lunches: sign-up sheet						\uparrow
Instance IPs						

-

Daily workshop material, slides & exercises

EVOL Intensive	UTION and immersive to	AND G	ENOMICS unities		WORKSHOPS	LEARNING	PEOPLE	APPLY	INFORMATIO	N
	Faculty Luncl	hes: sign-up	sheet							
	Instance IPs									
•	Cheatsheet – copying files from/to the instances									
COMPET	ITIONS									
• CODE OF • Week 1 :	BINGO! CONDUCT Code of condu	ict, contact p y, 2024	eople Karin and Jacob							
DATE	DAY	TIME	PRESENTER	TOPIC				LO	CATION	
Jan 21	Sunday	18 - 22	Everyone	Reception				Hot	tel Zlaty Andel	
Jan 22	Monday	09 – 12	Anna Karnkowska	Introduction & Orientation, O	City Information			Тоу	vn Theatre	
	Monday	14 - 17	Workshop Team	Lab introduction				Но	use of Prelate	
	Monday	19 – 22	Everyone	Scientific speed networking				Kru	mlov mill	
Jan 23	Tuesday	09 – 12	Rosa Fernández	Introduction to Phylogenomi	ics			Тоу	vn Theatre	
	Tuesdav	14 - 17	Workshop Team	Alignment and Alignment Tri	imming			Но	use of Prelate	

Good workshop practice

- Work together
- Ask lots of questions
- Take breaks
- Use cheat sheets
- Have lots of fun



Cheat sheets

Unix/Linux Command Reference



File Commands	System Info
1s – directory listing	date - show the current date and time
1s -al - formatted listing with hidden files	cal – show this month's calendar
cd dir - change directory to dir	uptime - show current uptime
cd – change to home	\mathbf{w} – display who is online
pwd - show current directory	whoami – who you are logged in as
mkdir <i>dir</i> – create a directory <i>dir</i>	finger user - display information about user
rm <i>file</i> – delete <i>file</i>	uname -a - show kernel information
rm -r dir – delete directory dir	cat /proc/cpuinfo - cpu information
rm -f file – force remove file	cat /proc/meminfo - memory information
rm -rf <i>dir</i> – force remove directory <i>dir</i> *	man <i>command</i> – show the manual for <i>command</i>
cp file1 file2 - copy file1 to file2	df - show disk usage
cp -r dir1 dir2 - copy dir1 to dir2; create dir2 if it	du – show directory space usage
doesn't exist	free - show memory and swap usage
mv file1 file2 – rename or move <i>file1</i> to <i>file2</i>	where is app - show possible locations of app
if file 2 is an existing directory, moves file 1 into	which app - show which app will be run by default
directory file 2	Comprossion
ln -s file link - create symbolic link link to file	Compression
touch file - create or update file	tar cf file.tar files – create a tar named
cat > file – places standard input into <i>file</i>	file.tar containing files
more file – output the contents of file	tar xi file.tar – extract the files from file.tar
head file – output the first 10 lines of file	Grin compression
tail file – output the last 10 lines of file	top wif file top gr extract a top using Grip
tail -f file - output the contents of file as it	tar xzi iiie.tar.gz - extract a tar using Ozip
grows, starting with the last 10 lines	compression
Process Management	tar wif file tar bz2 - extract a tar using Bzin?
ps – display your currently active processes	gzip file - compresses file and renames it to
top – display all running processes	file 97
kill pid - kill process id pid	gzip -d file.gz - decompresses file.gz back to
killall proc - kill all processes named proc *	file
bg – lists stopped or background jobs; resume a	J
stopped job in the background	Network
fg – brings the most recent job to foreground	ping host – ping host and output results
fg \boldsymbol{n} - brings job \boldsymbol{n} to the foreground	whois domain – get whois information for domain
File Permissions	dig domain – get DNS information for domain
chmod octal file - change the permissions of file	dig -x host - reverse lookup host
to <i>octal</i> , which can be found separately for user,	wget file - download file
group, and world by adding:	wget -c file - continue a stopped download
• $4 - read(r)$	Installation
• 2 – write (w)	Install from source:
• 1 - execute (x)	/configure
Examples:	make
chmod 777 - read, write, execute for all	make install
chmod 755 – rwx for owner, rx for group and world	dpkg -i pkg. deb - install a package (Debian)
For more options, see man chmod.	rom -Uvh pkg. rom - install a package (RPM)
SSH	-rr- pgr
ssh user@host = connect to hast as user	Shortcuts
ssh -p port user@host - connect to host on port	Ctrl+C – halts the current command
port as user	Ctrl+Z - stops the current command, resume with
ssh-copy-id user@host - add your key to host for	fg in the foreground or bg in the background
user to enable a keyed or passwordless login	Ctrl+D - log out of current session, similar to exit
Searching	Ctrl+W – erases one word in the current line
	Ctrl+U – erases the whole line
grep pattern files - search for pattern in files	Ctrl+R – type to bring up a recent command
grep -r pactern air - search recursively for	!!! - repeats the last command

exit - log out of current session * use with extreme caution.

grep pat grep -r p pattern in dir command | grep pattern - search for pattern in the output of command

locate file - find all instances of file

Cheat sheets





Bard

Stack overflow, ChatGTP, Bard (Google).

Unix/Linux Command Reference



File Commands	System Info
1s – directory listing	date - show the current date and time
1s -al - formatted listing with hidden files	cal – show this month's calendar
cd dir - change directory to dir	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 dir	finger user – display information about user
rm file – delete file	uname -a - show kernel information
rm - r dir – delete directory dir	cat /proc/cpuinfo - cpu information
rm -f file – force remove <i>file</i>	cat /proc/meminfo - memory information
rm - rf dir – force remove directory dir *	man command – show the manual for command
cp file1 file2 - copy file1 to file2	df – show disk usage
cp - r dirl dir2 - copy dirl to dir2; create dir2 if it	du – snow directory space usage
doesn't exist	free – show memory and swap usage
mv filel file2 – rename or move file1 to file2	where is app - snow possible locations of app
directory file?	which app - show which app will be run by default
ln -2 filo link create symbolic link link to file	Compression
touch file - create or update file	tar cf file.tar files - create a tar named
$cot \sim file$ places standard input into file	file.tar containing files
more file output the contents of file	tar xf file.tar - extract the files from file.tar
head file – output the first 10 lines of file	tar czf file.tar.gz files - create a tar with
tail file - output the last 10 lines of file	Gzip compression
tail -f file - output the contents of file as it	tar xzf file.tar.gz - extract a tar using Gzip
grows, starting with the last 10 lines	tar cjf file.tar.bz2 - create a tar with Bzip2
Brooce Management	compression
	tar xjf file.tar.bz2 - extract a tar using Bzip2
ps – display your currently active processes	gzip file - compresses file and renames it to
top – display all running processes	file.gz
kill pid - kill processes named proc *	gzip -d file.gz - decompresses file.gz back to
$\mathbf{b}_{\mathbf{r}}$ = lists stopped or background jobs: resume a	file
stopped job in the background	Network
$f \sigma = hrings$ the most recent job to foreground	ping host - ping host and output results
fg n - brings the most recent job to foreground	who is domain - get who is information for domain
File Dermissions	dig domain – get DNS information for domain
File Permissions	dig -x host - reverse lookup host
chmod octal file - change the permissions of file	wget file-download file
to octal, which can be found separately for user,	wget -c file - continue a stopped download
group, and world by adding:	
• $4 - 1$ cad (1) • $2 - $ write (w)	Installation
• $1 - execute(\mathbf{x})$	Install from source:
Examples:	./configure
chmod 777 - read, write, execute for all	make
chmod 755 – rwx for owner, rx for group and world	make install
For more options, see man chmod.	dpkg -1 pkg.deb - install a package (Debian)
20U	rpm - Ovn <i>pkg.rpm</i> – install a package (KPM)
33 П	Shortcuts
ssn user(host - connect to host as user	Ctrl+C – halts the current command
ssn -p port usergnost - connect to nost on port	Ctrl+Z - stops the current command, resume with
port as user	fg in the foreground or bg in the background
user to enable a keyed or passwordless login	Ctrl+D - log out of current session, similar to exit
aser to enable a keyeu of passwordless logili	Ctrl+W – erases one word in the current line
Searching	Ctrl+U – erases the whole line

grep pattern files - search for pattern in files grep -r pattern dir - search recursively for pattern in dir command | grep pattern - search for pattern in the output of command

locate file - find all instances of file

^k use with extreme caution.

!! - repeats the last command

exit - log out of current session

Ctrl+R - type to bring up a recent command

What is UNIX?

Operating system



What is UNIX?

Operating system



Why do we use it?

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

What is UNIX?

Operating system



Why do we use it?

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

... And how?

• Linux computers or servers

aws

- Computer clusters
- The cloud



AWS availability zones



How it works

AMI ("Amazon Machine Image")

Base computer with all data and software





Own copy of AMI = Instance (Virtual machine, VM)

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.

Watch & listen

Follow along on your computer





Connecting to your instance



Windows: remote desktop software

Guacamole, X2GO



Linux/Mac: Terminal SSH ("Secure shell")



Instance addresses <u>http://evomics.org/2024-workshop-</u> <u>on-phylogenomics-cesky-krumlov/</u>

FACULTY AND GENOMICS

- Click here to see our Faculty arrival and departure dates
- Faculty Bios and Teaching Assistant Bios
- Faculty Lunches: sign-up sheet
- Instance IPs
- Cheatsheet copying files from/to the instances



Find your name and copy your IP address



,						о щи 🗖 "А у	
\leftarrow	\rightarrow G \bigcirc	ocs.google.com/s	preadsheets/d/1-vNj1erBvvqzKp3c1k	다. C Search		ව 🛄 🔽 🗸	» =
	Workshop on Phyloenomics Instance List $4 \otimes 2$ File Edit View Insert Format Data Tools Extensions Help $5 \otimes 3$ $5 \otimes 3$ $5 \otimes 3$ $5 \otimes 3$ $5 \otimes 3$ $5 \otimes 3 \otimes 3 \otimes 3 \otimes 3$ $5 \otimes 3 \otimes 3 \otimes 3 \otimes 3$ $5 \otimes 3 \otimes 3 \otimes 3 \otimes 3 \otimes 3$ $5 \otimes 3 \otimes $						
C	오 승 금 뭄 100% ▼ Kč % .0 .0 123 Calibri ▼ - 12 + B I 중 A ▷ 田 문 ▼ 토 ▼ 남 ▼ A ▼ I: ^ 3						
A3	A3 - fx Phylo-Karin						
	А	В	С	D	E	F	
1	Name	IP address	Guacamole connection	ssh connection	RStudio server connection	Username: phyloge	
2	auv	3 2 1// 5	3 2 144 5:8080/quacamole	ssh genomics@3.2.144.5	3.2.144.5:8787	Date:	Ø
3	Phylo-Karin	52.91.211.21	52.91.211.21:8080/guacamole	ssh genomics@52.91.211.21	52.91.211.21:8787		
4	Phylo-Ivlarina	44.212.5.42	44.212.5.42:8080/guacamole	ssh genomics@44.212.5.42	44.212.5.42:8787		
5	Phylo-Jacob	3.89.92.242	3.89.92.242:8080/guacamole	ssh genomics@3.89.92.242	3.89.92.242:8787		•
6	Phylo-Gemma	54.205.120.236	54.205.120.236:8080/guacamole	ssh genomics@54.205.120.236	54.205.120.236:8787		
7	Phylo-Michal	44.207.6.114	44.207.6.114:8080/guacamole	ssh genomics@44.207.6.114	44.207.6.114:8787		
8	Phylo-Rosa	3.94.82.239	3.94.82.239:8080/guacamole	ssh genomics@3.94.82.239	3.94.82.239:8787		
9			:8080/guacamole	ssh genomics@	:8787		
10			:8080/guacamole	ssh genomics@	:8787		
11			:8080/guacamole	ssh genomics@	:8787		+
12			:8080/guacamole	ssh genomics@	:8787		
13			:8080/guacamole	ssh genomics@	:8787		
14			·8080/quacamole	ssh genomics@	·8787		



- Open your internet browser (e.g. Google Chrome)
- Paste the IP address followed by ':8080/guacamole'
- 52.91.211.21:8080/guacamole

$\leftarrow \rightarrow \mathbf{G}$	Q 52.91.211.21:8080/guacamole/	\rightarrow Q Search	▽ 👱 🐵 Z
🌣 Most Visited 🛛 🧐 Gett	ing Started		



Enter the username "phylogenomics" and password



Select Desktop, enter the same user name and pw again







Open terminal window using this icon



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

But First...

- The domain address will change every day after we stop and re-start the instances.
- Each morning, you will need to return to the "Instance List" webpage, retrieve your new address and log in again

\leftarrow	ightarrow C $ ightarrow$ https://do	ocs.google.com/s	preadsheets/d/1-vNj1eYBvvqzKp3c1	公式 Q Search	♡ 🛃 🤷 📄 🤇	ම 🗰 🔽 දු 💈	» ≡
E	■ Workshop on Phyloenomics Instance List ☆ ☆ ☆ File Edit View Insert Format Data Tools Extensions Help ▼ K						
C	えちさ骨骨100% 🕶	Kč % .0 .0	0 123 Calibri - 12	+ B I ÷ A è 🖽 🗄	· ≣ • ↓ • ÷ • A •	: ^	31
A3	✓ fx Phylo-Karin						
	А	В	С	D	E	F	
1	Name	IP address	Guacamole connection	ssh connection	RStudio server connection	Username: phyloge	
2	guy	3.2.144.5	3.2.144.5:8080/guacamole	ssh genomics@3.2.144.5	3.2.144.5:8787	Date:	${\bf \odot}$
3	Phylo-Karin	52.91.211.21	52.91.211.21:8080/guacamole	ssh genomics@52.91.211.21	52.91.211.21:8787		
4	Phylo-Marina	44.212.5.42	44.212.5.42:8080/guacamole	ssh genomics@44.212.5.42	44.212.5.42:8787		•
5	Phylo-Jacob	3.89.92.242	3.89.92.242:8080/guacamole	ssh genomics@3.89.92.242	3.89.92.242:8787		•
6	Phylo-Gemma	54.205.120.236	54.205.120.236:8080/guacamole	ssh genomics@54.205.120.236	54.205.120.236:8787		
7	Phylo-Michal	44.207.6.114	44.207.6.114:8080/guacamole	ssh genomics@44.207.6.114	44.207.6.114:8787		Q
8	Phylo-Rosa	3.94.82.239	3.94.82.239:8080/guacamole	ssh genomics@3.94.82.239	3.94.82.239:8787		
9			:8080/guacamole	ssh genomics@	:8787		
10			:8080/guacamole	ssh genomics@	:8787		
11			:8080/guacamole	ssh genomics@	:8787		+



Copy & Paste

AVOID COPYING AND PASTING WHEREVER POSSIBLE! But if you do need to...

Press Ctrl+Alt+Shift

Paste the text into the box with right click \rightarrow Paste

Press Ctrl+Alt+Shift again

You can now paste into the instance using right click







Your final task before we get started!

Make sure that typing tilde (\sim), backslash (\), pipe (|), and carat ($^$) in the terminal works.

Google search to find these on your computer if you don't know where they are.

Questions?



Commands overview



pwd	print working directory
~	home
•	here
• •	one directory up
mkdir	make new directory
cd	change directory
touch	create file
ls	list
man	manual
mv	move
rm	remove

ср	сору
gunzip	unzip
tar	unarchive
head	first (n=10) lines
tail	last (n=10) lines
cat	concatenate
WC	word count
grep	pattern search
	pipeline
sed	stream editor
chmod	change file modes

The terminal (command line, shell, 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?



This is your working directory, i.e. where you currently are.















Create a new directory called "Data" in your current directory

\$ mkdir ./Data

Change into the new directory

\$ cd Data

Where are you now? What is your present working directory?

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





Make an empty file "rags"

\$ touch rags

And another two "Heaven" and "Earth"

\$ touch Heaven Earth

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

\$ ls




Now list **all** of the files in the directory

\$ ls -a





Now list **all** of the files in the directory

\$ ls -a





Now list **all** of the files in the directory **\$ 1s -a**

. points to the current directory

.. points to one directory above



Now list **all** of the files in the directory \$ 1s -a

. points to the current directory.. points to one directory above



. and .. are used for specifying location

Whenever you do anything on Unix (move around, move a file, rename a file, run a program or script, etc...) you have to tell the system where that thing is using a path.

. and .. are part of RELATIVE paths





Create a directory called New within the phylogenomics directory using the RELATIVE PATH

\$ mkdir ../New





Move from Data to New RELATIVE PATH

\$ cd ../New





Move from New to home RELATIVE PATH



Relative paths will always change depending on your location. The alternative is ABSOLUTE paths.

These always start from root and will never change.





Move from your home directory to New ABSOLUTE PATH

\$ cd /home/phylogenomics/New





Move from your home directory to New ABSOLUTE PATH

\$ cd /home/phylogenomics/New

Move from Data to New ABSOLUTE PATH

\$ cd /home/phylogenomics/New

A note about . dot



. means in your (present) working directory

This command means "List everything that's in the (present) working directory" \$ 1s ./

This command means "List everything that's in the working directory within a subdirectory called Data"

\$ ls ./Data/

In most cases, people don't use ./ at the beginning of a path. As long as the file/directory is within your working directory, the command will work.

\$ ls ./Data/ = \$ ls Data

Let's practice



Where am I right now? (Should be the Data directory)\$ pwdChange 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.

Work smarter, not harder!



Tab completion is a nice trick to save you typing paths

For this example we are going to list everything in the directory /home/phylogenomics/workshop_materials/. Start by typing followed by pressing tab twice quickly. This shows the contents of the root directory:

			phylo	genomics@ip-	-172-31-91-145: ~/\	vorkshop_mate	erials			\odot
File Edit View	Search Termin	al Help								
hylogenomics	akrumlov:[~/	workshop_ma	terials]\$ ls							
s l sattr l hylogenomics@	İsb_release Isblk @krumlov:[~/	lscpu lsdiff workshop_ma	lshw lsinitramfs terials]\$ ls	lsipc lslocks	lslogins lsmem	lsmod lsns	lsof lspci	lspcmcia lspgpot	lss3 lsusb	

Work smarter, not harder!



Now type: **\$ 1s /h** followed by tab once.

The path to the /home/ directory has filled in.

Now type: **\$ 1s /home/p** followed by tab once.

The path to the /home/phylogenomics/ directory has filled in.

Finally type: **\$ 1s /home/phylogenomics/w** finish the pain, and then enter.

followed by tab once to

You've now listed that directory contents.

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

Work smarter, not harder!



Two more tricks for less typing! The * (asterisk) represents any character For example: **\$ 1s /home/phylogenomics/*.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 your previous commands listed, simply type **\$ history**

Questions?



Binary programs



These are all programs installed on the Unix machine. They can be found in /bin **\$ 1s /bin**

	\sim \sim		
File Edit View Search Terminal Help			
File Edit View Search Terminal Help h2xs h5c++ h5cc h5fc hardlink hbf2gf hbpldecode hciattach hciconfig hcitool hd head helixturnhelix helpztags hex2hcd hex2hcd hipercdecode hiperc	<pre>pod2html pod2man pod2markdown pod2readme pod2texi pod2text pod2usage podchecker podebconf-display-po podf pollinate polydot pon pooltype post-grohtml podc ppdtml ppdi ppdi ppdmerge ppdpo pobs</pre>	<pre>zegrep zeisstopnm zenity zfgrep zforce zgrep zip zipcloak zipdetails zipgrep zipinfo zipnote zipsplit zjsdecode zless zmore znew zoom2sam.pl zstd zstdgrep zstdless zstdgrep zstdless zstdmt</pre>	
<pre>hmmemit hmmfetch phylogenomics@krumlov:[~]\$ ls /bin</pre>	ppltotf ppm3d		

These include pwd, mkdir, ls ...

Binary programs have manuals



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

Open the manual page for Is \$ man 1s

Scroll through (enter) and find the options for:

long listing format (-I), human-readable file sizes (-h) and sort by modification time (-t).

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

\$ ls -l -t -h ./Data = \$ ls -lth ./Data

PATH



The computer needs to know where a program is so that it can access the code to run the program.

The PATH environment variable is a list of locations your computer looks for programs.

You can either provide the path to the program you want to run

\$ /usr/bin/mkdir

PATH



The computer needs to know where a program is so that it can access the code to run the program.

The PATH environment variable is a list of locations your computer looks for programs.

You can either provide the path to the program you want to run

\$ /usr/bin/mkdir

Or make sure the program is in your PATH environment variable

To view locations in your PATH environment variable: **\$ echo \$PATH**

There are ways to add new locations to your PATH, but that is for another time.

Let's practice some more home lib bin phyloubuntu main genomics Working Data

Earth

Heaven

rags



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

Afterwards your file structure should look like this.







Now move Earth, too.

Remember to Tab complete!





Moving Files

mv can also be used to rename files.

Let's change rags to riches.



Let's practice some more lib bin home phyloubuntu main genomics Working Data riches Earth



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.

Let's practice some more home lib bin phyloubuntu main \sim genomics Working Data riches



Deleting Files

Now let's delete the entire 'Working' directory including Earth.

\$ rm -i ../Working/Heaven

Let's practice some more home lib bin phyloubuntu main \sim genomics Working Data riches



Deleting Files

There is no 'Trash' or 'Recycle Bin' in Unix!

Once gone, files are gone forever!

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.)



Let's practice some more lib home bin phyloubuntu main genomics riches Backup Data riches riches



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

Data management



Some files can become quite big so people will archive directories and 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 archived directory
- .gz means that it is gzipped file

These can be used alone or in combination

To uncompress a tar archive (x = extract, v = verbose, f = all files)

\$ tar -xvf <filename>

A Gzipped file

\$ gunzip <filename>

A Gzipped Tar archive

\$ tar -xzvf <filename>

Challenge 2



1. Change to the workshop_materials directory at the following path: ~/workshop_materials/unix

You should find a compressed directory: Sequences.tar.gz

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

tar gunzip cp mv rm -i cd gunzip mv mkdir

Questions?







Navigate to the workshop_materials directory
\$ cd ~/workshop_materials
Unarchive the Blast_Out.tar.gz
\$ tar -xzvf Blast_Out.tar.gz

Blast_Out.tar.gz





CP_Blast_seqs.fna LP_Blast_seqs.fna

LINUX TERMINAL FOR BEGINNERS





CP_Blast_seqs.fna LP_Blast_seqs.fna

View the first 10 lines of a file

\$ head CP_Blast_seqs.fna

To view the first 30 lines of the file

\$ head -n 30 CP_Blast_seqs.fna

LINUX TERMINAL FOR BEGINNERS







CP_Blast_seqs.fna LP_Blast_seqs.fna

View the last 10 lines of a file \$ tail CP_Blast_seqs.fna

To view the last 30 lines of the file

\$ tail -n 30 CP_Blast_seqs.fna

LINUX TERMINAL FOR BEGINNERS






CP_Blast_seqs.fna LP_Blast_seqs.fna

Print the entire file

\$ cat CP_Blast_seqs.fna





Many files are too large to meaningfully view in terminal or to edit in a unix text editor.







'cat' can also combine multiple files

\$ cat CP_Blast_seqs.fna LP_Blast_seqs.fna >
All_blast_seqs.fna

And then count the number of lines in each fasta file to confirm they have been combined.

\$ wc -1 *.fna





Use 'grep' to print occurrences of a pattern
\$ grep ">" CP_Blast_seqs.fna





Use 'grep' to print occurrences of a pattern
\$ grep ">" CP_Blast_seqs.fna

Create a new files of the fasta headers

\$ grep ">" CP_Blast_seqs.fna >
CP_blast_headers.txt

\$ grep ">" LP_Blast_seqs.fna >
LP_blast_headers.txt

blast_seqs.fna CP_blast_headers.txt LP_blast_headers.txt





Use 'grep' for to count the number of times pattern occurs

\$ grep -c ">" CP_Blast_seqs.fna





Use 'grep' for to count the number of times pattern occurs \$ grep -c ">" CP Blast seqs.fna

Quotations marks are vital!

\$ grep -c > CP_Blast_seqs.fna

CP Blast seqs.fna is now empty and this can't be undone.



What can you do to your file to protect them? Change permissions

\$ chmod 444 CP_Blast_seqs.fna



What can you do to your file to protect them? Change permissions

\$ chmod 444 CP_Blast_seqs.fna

You might know and need this when you write scripts you want to execute. Sometimes you need to change permission to allow them to be executed.

\$ chmod +x my_script.sh





Search for headers that are not partial sequences \$ grep -v "partial" LP_headers.txt





Replacing Text in Large files sed 's/**FIND**/**REPLACE**/g' filename > output_file \$ sed "s/ /_/g" CP_blast_headers.txt







Regular expressions

Expression	Modern equivalent	Pattern matched
		a single character
.+		one or more characters
.*		zero or more characters
.?		maybe present
٨		first in the line
\$		last in the line
[0-9]	\d	digits
[a-zA-Z]	\w	letters
	\s \t	space
{2}		exactly 2 characters long
{2,4}		between 2 and 4 characters long
[ACGT]		a specific set of characters

Regular expressions



Use sed -E to use extended regular expressions

\$ sed -E 's/(>[A-Z0-9.]+)(.+)/\1/' CP_blast_headers.txt

Store pattern in memory using parentheses

Print out only the GenBank accessions

Regular expressions



Use sed -E to use extended regular expressions

\$ sed -E 's/(>[A-Z0-9.]+)(.+)/\2/' CP_blast_headers.txt

Store pattern in memory using parentheses

Print out the rest of the headers

Pipes |



Combine the fasta files from the blast output and identify how many sequences they contain

\$ cat CP_Blast_seqs.fna LP_Blast_seqs.fna | grep ">" | wc -1

When building pipelines, it is useful to pipe to head to follow along

\$ cat	CP_	_Blast_seqs.fna	LP	_Blast_	_seqs.fna	head	
\$ cat	CP_	_Blast_seqs.fna	LP	_Blast_	_seqs.fna	grep ">"	head
\$ cat	CP_	_Blast_seqs.fna	LP	_Blast_	_seqs.fna	grep ">"	wc -l

Pipes



String together many commands to count the number of unique accessions from these blast results

\$ cat *.fna Print out all fasta files

\$ cat *.fna | grep ">" Find all the fasta headers

Extract accessions

\$ cat *.fna | grep ">" | sed -E "s/(>[A-Z0-9.]+)(.+)/\1/"

Sort accessions

\$ cat *.fna | grep ">" | sed -E "s/(>[A-Z0-9.]+)(.+)/\1/" | sort

Keep only unique accessions

\$ cat *.fna | grep ">" | sed -E "s/(>[A-Z0-9.]+)(.+)/\1/" | sort | uniq Count unique accessions

Challenge 3



Imagine you want to obtain the GenBank Accessions from the headers and you want to remove the version. Why is this command not working as expected?

\$ sed -E "s/(>[A-Z0-9.]+)(.+)/\1/" CP_blast_headers.txt |
sed "s/.1//g"

What would you need to change to make it work?

Loops



Iterate over e.g. files to execute a command repeatedly.

```
$ for i in *; do echo $i; done
```

Assign variable "i" to all files that end in .fna, then carry out command on all values of "i".

\$ for i in *.fna; do grep -c "CAT" \$i; done

There is usally more than one way to do things. Try the following:

\$ grep -c "CAT" *.fna

Downloading



From an internet URL

wget <url>

wget -P /path/to/where/the/download/should/be <url>

Save your data every day

We will launch new instance daily so everything you have done today will be gone tomorrow. Use scp or rsync to copy the notes you want to save to your own computer,

scp -r phylogenomics@<your.IP.address>:/Location/On/Instance /Local/path
rsync -avz phylogenomics@<your.IP.address>:/Location/On/Instance /Local/path

\$ rsync -avz phylogenomics@ec

Downloading



GUI: Filezilla https://sourceforge.net/projects/filezilla/



Logging in to the instance through your local terminal

ssh -Y phylogenomics@<your.IP.address> Accept and continue the connection with 'yes'

💿 🜔 🍵 🛅 Karin — steffk1@gw344:/nobackup/rokaslab/steffk1/Shen_etal_2020/alternaria_phylo_inference/ma	💿 😑 💼 Karin — phylogenomics@ip-172-31-91-145: ~ — ssh -Y phylogenomics@52.91.211.21 — 104×31
<pre>((anaconda3) IFB-FK0G-Karin:~ Karin\$ ssh -Y genomics@52.91.211.21 The authenticity of host '52.91.211.21 (52.91.211.21)' can't be established. ED25519 key fingerprint is SHA256:eH/xB39SR120qz5q/7nI+tWI04fYpw2IaxFuVB8RMLo. This key is not known by any other names</pre>	
Are you sure you want to continue connecting (yes/no/[tingerprint])? yes	[(anaconda3) IFB-FK0G-Karin:~ Karin\$ ssh -Y phylogenomics@52.91.211.21 [phylogenomics@52.91.211.21's password:

	## Workshop on Phyloenomics 2024 ##
	## Cesky Krumlov ##
	## @evonics #evonics2024 ##
3	
	Welcome to Ubuntu 22.04.3 LTS (6.2.0-1017-aws).
	System information as of Fri Jan 19 16:10:20 CET 2024
	System load: 0.27392578125 Processes: 234
	Usage of /: 32.9% of 484.63GB Users logged in: 0
	Memory usage: 12% IPv4 address for docker0: 172.17.0.1
	Swap usage: 0% IPV4 address for ens5: 172.31.91.145
	Last togin: web Jan 1/ 14:31:44 2024 from 194.228.207.2
	physogenomics and antovi [2] a

Questions?



Quests!

Do these quests in whichever order you like. Whatever sounds interesting to you.

Easy

- open, modify and exit with saving a text file in an editor of your choice: nano, vim, emacs
- explore the commands: tr, cut, less, tree
- How do you append to a file without overwriting it?

You haven't heard about all of this. Use your creativity, neighbor and the internet to tackle these!

Fun

Beautiful figures: <u>Tidy Tuesday</u>, <u>R graph</u> <u>gallery</u>, <u>xkcd style</u>, color choices (<u>Ghibli</u>, <u>Wes</u> <u>Anderson</u>, <u>I want hue</u>, <u>Scico</u>))

Challenging

 what about the programming language 'awk'? :P

What are your favorite one-liners or UNIX tricks? (Collect them in a Google Docs!)