

Lab Introduction and an Introduction to Unix

Anna Karnkowska, Eric Salomaki



Many slides courtesy Dr. Sophie Shaw

Overview

- Introduction
- Using the Amazon Machine Images (AMI)
- Introduction to Unix

These Slides!

http://evomics.org/workshops/2019-workshop-on-

phylogenomics-cesky-krumlov/

WORKSHOPS	LEARNING	PEOPLE	APPLY	INFORMATIC	л
Workshop on Ger	nomics				
Workshop on Pop	oulation and Spec	iation Genomi	cs		
Workshop on Phy	logenomics				2019 Workshop on Phylogenomics, Cesky Krumlov
Workshop on Mol	lecular Evolution				2017 Workshop on Phylogenomics, Cesky Krumlov
Harvard Universit	y Workshops				
Workshop on Mic	robiome and Tran	iscriptome Ana	alysis, Durbar	n, South Africa	
Advanced Topics					
	WORKSHOPS Workshop on Ger Workshop on Pop Workshop on Phy Workshop on Mo Harvard Universit Workshop on Mic Advanced Topics	WORKSHOPS LEARNING Workshop on Genomics Workshop on Population and Spec Workshop on Phylogenomics Workshop on Molecular Evolution Harvard University Workshops Workshop on Microbiome and Tran Advanced Topics	WORKSHOPSLEARNINGPEOPLEWorkshop on GenomicsWorkshop on Population and Speciation GenomicsWorkshop on PhylogenomicsWorkshop on Molecular EvolutionHarvard University WorkshopsWorkshop on Microbiome and Transcriptome AnaAdvanced Topics	WORKSHOPSLEARNINGPEOPLEAPPLYWorkshop on Genomics	WORKSHOPS LEARNING PEOPLE APPLY INFORMATION Workshop on Genomics Workshop on Population and Speciation Genomics Image: Comparis and Speciation Genomics Workshop on Phylogenomics Image: Comparis and Speciation Genomics Image: Comparis and Speciation Genomics Workshop on Phylogenomics Image: Comparis and Speciation Genomics Image: Comparis and Speciation Genomics Workshop on Phylogenomics Image: Comparis and Speciation Genomics Image: Comparis and Speciation Genomics Workshop on Molecular Evolution Image: Comparis and Transcriptome Analysis, Durban, South Africa Advanced Topics Image: Comparis and Transcriptome Analysis, Durban, South Africa

21 Jan	9a – 12p Scott Handley & Toni Gabaldón Introduction and Orie		Introduction and Orientation	Theater
	2p – 5p	Workshop Team	Lab Introduction	House of Prelate
	7p – 10p	Everyone	Scientific Speed Networking	TBD
22 Jan	9a – 12p	Antonis Rokas	Introduction to Phylogenomics	Theater
	2p – 5p	Workshop Team	Alignment and Alignment Trimming	House of Prelate
	7p – 10p	Workshop Team	Tree Visualization	House of Prelate
23 Jan	9a – 12p	Toni Gabaldón	Introduction to Phylogenetics, and Orthology and Paralogy	Theater
	2p – 5p	Toni Gabaldón	Orthology and Paralogy Prediction lab	House of Prelate
	7p – 10p	Workshop Team	Partitioning and Concatentation Laboratory	House of Prelate

Background on Unix, Cloud computing, and in working in an instance Introduction to Terminal and Unix

21 Jan	9a – 12p	Scott Handley & Toni Gabaldón	Scott Handley & Toni Gabaldón Introduction and Orientation	
	2p – 5p	Workshop Team	Lab Introduction	House of Prelate
	7p – 10p	Everyone	Scientific Speed Networking	TBD
22 Jan	9a – 12p	Antonis Rokas	Introduction to Phylogenomics	Theater
	2p – 5p	Workshop Team	Alignment and Alignment Trimming	House of Prelate
	7p – 10p	Workshop Team	Tree Visualization	House of Prelate
23 Jan	9a – 12p	Toni Gabaldón	Introduction to Phylogenetics, and Orthology and Paralogy	Theater
	2p – 5p	Toni Gabaldón	Orthology and Paralogy Prediction lab	House of Prelate
	7p – 10p	Workshop Team	Partitioning and Concatentation Laboratory	House of Prelate

How and why to build a multiple sequence alignment Is your alignment correct? How to trim your alignment for phylogenetics

21 Jan	9a – 12p	Scott Handley & Toni Gabaldón	n Introduction and Orientation T	
	2p – 5p	Workshop Team	Lab Introduction	House of Prelate
	7p – 10p	Everyone	Scientific Speed Networking	TBD
22 Jan	9a – 12p	Antonis Rokas	Introduction to Phylogenomics	Theater
	2p – 5p	Workshop Team	Alignment and Alignment Trimming	House of Prelate
	7p – 10p	Workshop Team	Tree Visualization	House of Prelate
23 Jan	9a – 12p	Toni Gabaldón	Introduction to Phylogenetics, and Orthology and Paralogy	Theater
	2p – 5p	Toni Gabaldón	Orthology and Paralogy Prediction lab	House of Prelate
	7p – 10p	Workshop Team	Partitioning and Concatentation Laboratory	House of Prelate

Important concepts for making your data accessible Exploring tools for looking at your tree(s) Web Based Software Toolkits

21 Jan	9a – 12p	Scott Handley & Toni Gabaldón	Introduction and Orientation	Theater
	2p – 5p	Workshop Team	Lab Introduction	House of Prelate
	7p – 10p	Everyone	Scientific Speed Networking	TBD
22 Jan	9a – 12p	Antonis Rokas	Introduction to Phylogenomics	Theater
	2p – 5p	Workshop Team	Alignment and Alignment Trimming	House of Prelate
	7p – 10p	Workshop Team	Tree Visualization	House of Prelate
23 Jan	9a – 12p	Toni Gabaldón	Introduction to Phylogenetics, and Orthology and Paralogy	Theater
	2p – 5p	Toni Gabaldón	Orthology and Paralogy Prediction lab	House of Prelate
	7p – 10p	Workshop Team	Partitioning and Concatentation Laboratory	House of Prelate

What is Partitioning and the why/when/how to do it Tools for partitioning your data How to concatenate your data

Good Workshop Practice

- PowerPoint interspersed with Challenges
- Ask us lots of questions!
- Work together
- Take breaks
- Use a cheat sheet, google, etc.
- Have Fun!

Unix/Linux Command Reference



File Commands	System Info
ls - 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 - snow kernel information
rm -f file - force remove file	cat /proc/cpuinto - cpuintormation
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 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 file1 to file2	whereis app - show possible locations of app
if file2 is an existing directory, moves file1 into	which app - show which app will be run by default
directory file2	Compression
In -s file link - create symbolic link link to file	tar of file tar files greate a tar named
touch file - create or update file	file tar containing files
cat > file - places standard input into 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
Brocoss Management	compression
Process Management	tar xjf file.tar.bz2 - extract a tar using Bzip2
ton - display your currency acuve processes	gzip file - compresses file and renames it to
kill nid - kill process id nid	file.gz
killall proc - kill all processes named proc *	gzip - a file.gz - decompresses file.gz back to
bq - lists stopped or background jobs: resume a	jiie
stopped job in the background	Network
fg - brings the most recent job to foreground	ping host - ping host and output results
fg n - brings job 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 octal, 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
changed 777 read write execute for all	make
child 777 - read, write, execute for an	make install
For more options, see man chmod	dpkg -1 pkg.deb - install a package (Debian)
CCU CCU	rpm -Uvh pkg.rpm - install a package (RPM)
SSN	Shortcuts
ssh user@nost - connect to host as user	Ctrl+C - halts the current command
nort 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
Soorching	Ctrl+W - erases one word in the current line
aren nattern files search for nattern in files	Ctrl+U - erases the whole line
aren -r nattern dir - search recursively for	trepeate the last command
pattern in dir	exit - log out of current session
command grep pattern - search for pattern in the	exact - log out of current session
output of command	
locate file - find all instances of file	* use with extreme caution.

https://files.fosswire.com/2007/08/fwunixref.pdf

Watching vs Doing



Watch and Listen



Try it for yourself

What is Unix?

Operating System



Why Unix?



 Bioinformatics software designed to run on Unix platforms.

Beast2

- Large amounts of data.
- Much faster than Windows PC.



A tool for automated alignment trimming



Bayesian evolutionary analysis by sampling trees



How Can We Use Unix?

- Linux computers or servers
- Computer clusters
- The cloud

• What we're going to use this week









AWS "Availability Zones" and Data Centres

How it Works



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



How it Works











Own copy of the AMI = Instance (Virtual Machine or 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.

Connecting to Your Instance







Secure Shell – "SSH" e.g. SSH or PuTTY



http://evomics.org/workshops/2019-workshop-onphylogenomics-cesky-krumlov/

Anna Karnkowska, University of Warsaw

Antonis Rokas, Vanderbilt University

Miscellaneous

- Instance Addresses Check Daily
- Faculty Lunches Sign Up!
- T-shirt Competition

SCHEDULE

Week 1: 20-26 January 2019

E	Worksho File Edit	p on Phyloge View Insert	enomics 2019 Format Data	Instance List Tools Add-ons	t ☆ Help	
F	~ 8 7	100% - \$	% .0, .00 1	23 - Arial	-	
fx						
	A	В	С	D	E	
1	Name	Public Domain	Public Domain			
2	Eric	ec2-34-227-31-223.compute-1.amazonaws.com				
3	Marina	ec2-3-82-175-48.compute-1.amazonaws.com				
4	Jacob	ec2-3-83-102-107.compute-1.amazonaws.com				
5	Ania	ec2-52-91-1-245	c2-52-91-1-245.compute-1.amazonaws.com			



Find your name and copy your public domain

Open your internet browser (e.g. Google Chrome)



Paste the public domain followed by :8080/guacamole

Example:

ec2-34-227-31-223.compute-1.amazonaws.com:8080/guacamole









Select Desktop

ALL CONNECTIONS

- Desktop
- >_ Terminal



Enter the username "phylogenomics" and password again

Login to ip-172-	30-3-212
8	Just connecting
Session	Xorg
username	genomics
password	
	OK Cancel

, Open a 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

	Worksho	p on Phyloge	nomics 2019	Instance List	$\overline{\lambda}$	
Ш	File Edit	View Insert	Format Data	Tools Add-ons	Help	
2	~ ē 7	100% - \$	% .000_ 1	23 - Arial	Ŧ	
fx						
	A	В	С	D	E	
1	Name	Public Domain				
2	Eric	ec2-34-227-31-2	ec2-34-227-31-223.compute-1.amazonaws.com			
3	Marina	ec2-3-82-175-48.compute-1.amazonaws.com				
4	Jacob	ec2-3-83-102-107.compute-1.amazonaws.com				
5	Ania	ec2-52-91-1-245	c2-52-91-1-245.compute-1.amazonaws.com			



Copying and Pasting



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 -> 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 (~), backslash (\), pipe (|), and carat (^) in the terminal works.

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





Introduction to Unix

Eric Salomaki

- BA, The Evergreen State College, 2005
- MS, Ohio University, 2012
 - Molecular systematics of freshwater red algae
- PhD, University of Rhode Island, 2017
 - Red algal parasite evolution
- Postdoc, University of Rhode Island, 2017-18
 - Functional ecology of microeukaryotes in the North Atlantic
- Postdoc, Institute of Parasitology, CAS 2018-Present
 - Impact of transitioning to and from parasitism on genome evolution





The Terminal



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?



phylogenomics@ip-172-30-3-100:~\$ pwd /home/phylogenomics phylogenomics@ip-172-30-3-100:~\$

This is your "present working directory".






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



Change into this directory



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.





Now let's create some directories and files



Make an empty file



And another two

~/Data\$ touch Heaven Earth

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



phylogenomics@ip-172-30-3-100:~/Data\$ touch rags phylogenomics@ip-172-30-3-100:~/Data\$ touch Heaven Earth phylogenomics@ip-172-30-3-100:~/Data\$ ls Earth Heaven rags









Now list ALL of the files





phylogenomics@ip-172-30-3-100:~/Data\$ ls -a Earth Heaven rags phylogenomics@ip-172-30-3-100:~/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





. and .. are used for specifying location

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

. and .. are part of RELATIVE paths







Create a directory called <u>New</u> within the phylogenomics directory using the RELATIVE PATH



Move from Data to New

RELATIVE PATH







Moving from New to home

RELATIVE PATH





Moving 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.













A Note About Dot!



. means "In your pwd ..." This command means "List everything that's in the pwd"

\$ ls ./

This command means "List everything that's in the pwd 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 pwd, the command will work.

Let's put this to practice



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



Change to the directory above



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 complete 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 tab twice quickly



This shows the contents of the root directory

Work smarter, not harder!



Now type:



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

\$ ls /home/

Now type:

\$ ls /home/p

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



Finally type:

\$ ls /home/phylogenomics/w

Followed by tab once to finish the path, and then enter. 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!

* Represents any character For example:

\$ ls /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 these commands listed, simply type



Any Questions???



Binary programs



These are all programs installed on the Unix machine.

They can be found in /bin

\$ ls /bin

phylogenomics@ip-172-30-3-100:~/Data\$ ls /bin							
bash	bzmore	fgconsole	lesspipe	netcat	ping	stty	uncompress
btrfs	cat	fgrep	ln	netstat	ping4		unicode_start
btrfs-debug-tree	chacl	findmnt	loadkeys	networkctl	ping6	sync	vdir
btrfs-find-root	chgrp	fsck.btrfs	login	nisdomainname	plymouth	systemctl	wdctl
btrfs-image	chmod	fuser	loginctl	ntfs-3g	ps	systemd	which
btrfs-map-logical	chown	fusermount	lowntfs-3g	ntfs-3g.probe	pwd	systemd-ask-password	whiptail
btrfs-select-super	chvt	getfacl	ls	ntfscat	rbash	systemd-escape	ypdomainname
btrfs-zero-log	CD.	grep	lsblk	ntfscluster	readlink	systemd-hwdb	zcat
btrfsck	cpio	gunzip	lsmod	ntfscmp	red	systemd-inhibit	ZCMD
btrfstune	dash	gzexe	mkdir	ntfsfallocate	rm.	systemd-machine-id-setup	zdiff
bunzip2	date	gzip	mkfs.btrfs	ntfsfix	rmdir	systemd-notify	zegrep
busybox	dd	hciconfig	mknod	ntfsinfo	rnano	systemd-sysusers	zfgrep
bzcat	df	hostname	mktemp	ntfsls	run-parts	systemd-tmpfiles	zforce
bzcmp	dir	ip	моге	ntfsmove	sed	systemd-tty-ask-password-agent	zgrep
bzdiff	dmesg	journalctl	mount	ntfsrecover	setfacl	tar	zless
bzegrep	dnsdomainname	kbd_mode	mountpoint	ntfssecaudit	setfont	tempfile	zmore
bzexe	domainname	kill	mt	ntfstruncate	setupcon	touch	znew
bzfgrep	dumpkeys	kmod	mt-gnu	ntfsusermap	sh	true	
bzgrep	echo	less	mv	ntfswipe	sh.distrib	udevadm	
bzip2	ed	lessecho	nano	open	sleep	ulockmgr_server	
bzip2recover	egrep	lessfile	nc	openvt	SS	umount	
bzless	false	lesskev	nc.openbsd	pidof	static-sh	uname	

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 <PROGRAM>

Open the manual page for Is

\$ man ls

Scroll through (enter) and find the options for:

long listing format (-I), human-readable 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



PATH



The computer needs to know where a program is so that it will run

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 will run

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














Therefore try to ALWAYS use rm –i



riches









Typical File Sizes





On the Illumina NextSeq 3,000,000 reads = 1 Gb

But typically your experimental design/sequencing strategies will result in much more than 3,000,000 reads. 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 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



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
ср	mv
rm —i	mkdir
cd	







Navigate to the workshop_materials directory

\$ cd ~/workshop_materials

Unarchive the Blast_Out.tar.gz

\$ tar -xzvf Blast_Out.tar.gz

phylogenomics@ip-172-30-3-100:~/workshop_materials\$ tar -xzvf Blast_Out.tar.gz Blast_Out/ Blast_Out/LP_blast_seqs.fna Blast_Out/CP_blast_seqs.fna phylogenomics@ip-172-30-3-100:~/workshop_materials\$



View the first 10 lines of a file

\$ head CP_blast_seqs.fna



\$ head – 30 CP_blast_seqs.fna



CP_Blast_seqs.fna LP_Blast_seqs.fna



TTGCCCTTCTT

View the last 10 lines of a file

\$ tail CP_blast_seqs.fna



Now view the last 30 lines of the file

\$ tail -30 CP blast seqs.fna

CP Blast seqs.fna LP Blast seqs.fna

Blast Out

phylogenomics



Print the entire file

\$ cat CP_blast_seqs.fna



>AY237287.1 Osmundaria spiralis 18S ribosomal RNA gene, partial sequence CCCCCACCTGGTTGATCCTGCCAGTGGTATATGCTTGTCTCAAAGACTAAGCCATGCAAGTGCCTGTATGAGTGTTTATA CAACGAAACTGCGAATGGCTCGGTAAAACAGCAATAATTTCTTCGGGGGTAATACTACTCGGATAACCG raatacgtgctacaaaggcgagatcgctctcgtggtactaagtattaggtacaagccagt ATCTTCTGATCGCACTCTGTGCGACGCCCCGATCAAATTTCTGACCTATCAACTATGATGGTAAGGTAG TTGTGACGGGTAACGGACCGTGGGTGCGGGATTCCGGAGAGGGAGCCTGAGAAACGGCTACCA INNNNNNNNACCCAATCCGAACTTCGGGAGGTAGTGACAAAAATATCAATAGGGGGGA GAATGAAAACAATGTAAACACCTTATTGAGGAACCAGCAGAGGGCAAG IGTAAGCGTATACCAAAGTTGTTGCAGTTAAAACGCTCGTAGTC ATTAGAGTGTTCAAAGCGAGCGATTGCCATGAATACAT GGTTTGTCGGTATCAGGTAATGATTAAGAGGGACGGTCGGGGGGCAT GGAAGACAAACAGCTGCGAAAGCGTCTGCCAAGGACG CGATCAGATACCGTCGTAGTCTTTACTATAAACGATGAG TEGGEACEETEEGGAAAEEAAAGTGTTTGETTTEEGGGGGGGGGTATGGTEGCAAGTETGAAAETTAAAGGAA AAGGGCATCACCGGGTGTGGAGCCTGCGGCTTAATTTGACTCAACACGGGAAAACTTACCAGGTCCAGA TTGACAGATTGAGAGCTCTTTCTTGATTCTATGGTTGGTGGTGCATGGCCGTTCTTAGTTGGTGGAGTGATC AACGAGCGAGACCTGGGCGTGCTAAATAGGGTATGTTAT TTTAGTCTATGGAAGCTCCAGGCAATAACAGGTC ACTGAACGGTTCAATGGGTGAGGTAGTGCGAAAGCATGACCCAATCTCTAAAT ACGAGGAATACCTTGTAAGCGCGAGTCATCATCT ΔΔΟΤ AGGATCAAGCTAA

phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$





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

CP_Blast_seqs.fna LP_Blast_seqs.fna





CP_Blast_seqs.fna LP_Blast_seqs.fna



CP_Blast_seqs.fna LP_Blast_seqs.fna

Use 'grep' to print occurrences of a pattern

\$ grep ">" CP_blast_seqs.fna

MG272245.1 Neochondria nidifica voucher UC2026095 185 ribosomal RNA gene, partial sequence AF427537.1 Womersleyella setacea small subunit ribosomal RNA gene, complete sequence MF093929.1 Gredgaria maugeana isolate PD1230 18S ribosomal RNA gene, partial sequence HM560626.1 Digenea simplex voucher Dig sim 18S ribosomal RNA gene, partial sequence MG272238.1 Chondria crassicaulis voucher SAP115362 18S ribosomal RNA gene, partial sequence AF339897.1 Lenormandia sp. MELU_000065 18S ribosomal RNA gene, complete sequence MF093941.1 Periphykon beckeri isolate JH1427 185 ribosomal RNA gene, partial sequence MF093947.1 Polysiphonia scopulorum isolate PD0899 185 ribosomal RNA gene, partial sequence JX828186.1 Pterosiphonia pennata voucher CH816 small subunit ribosomal RNA gene, partial sequence AF427526.1 Boergeseniella fruticulosa small subunit ribosomal RNA gene, complete sequence MF093935.1 Lophocladia kuetzingii isolate PD1509 185 ribosomal RNA gene, partial sequence JX828165.1 Boergeseniella thuvoides voucher CH824 small subunit ribosomal RNA gene. partial seguence HM582914.1 Osmundaria obtusiloba voucher 03151 small subunit ribosomal RNA gene, partial sequence AF427527.1 Enelittosiphonia stimpsonii small subunit ribosomal RNA gene, complete sequence JX828192.1 Polysiphonia brodiei voucher CH410 small subunit ribosomal RNA gene, partial sequence AF203890.1 Heterocladia australis small subunit ribosomal RNA gene, complete sequence KX499570.1 Brongniartella australis isolate PD931 small subunit ribosomal RNA gene, partial sequence FJ153773.2 Uncultured eukaryote clone 53 18S ribosomal RNA gene, partial sequence AB219915.1 Polysiphonia sp. SNI07 gene for 18S rRNA, complete sequence MF093961.1 Vertebrata thuyoides isolate PD0546 185 ribosomal RNA gene, partial sequence JX828170.1 Leptosiphonia schousboei voucher CH826 small subunit ribosomal RNA gene, partial sequence JX828189.1 Symphyocladia linearis voucher CH419 small subunit ribosomal RNA gene, partial sequence JX828171.1 Neosiphonia elongella voucher CH415 small subunit ribosomal RNA gene, partial sequence AF427530.1 Polysiphonia fucoides small subunit ribosomal RNA gene, complete sequence AF203892.1 Heterocladia umbellata small subunit ribosomal RNA gene, complete sequence AB219910.1 Polysiphonia sp. BRI50 gene for 185 rRNA, complete sequence AF203885.1 Neosiphonia savatieri small subunit ribosomal RNA gene, complete sequence JX828179.1 Polysiphonia atlantica voucher CH1268 small subunit ribosomal RNA gene, partial sequence JX828175.1 Neosiphonia yendoi voucher CH420 small subunit ribosomal RNA gene, partial sequence JX828181.1 Polysiphonia elongata voucher CH1416 small subunit ribosomal RNA gene, partial sequence MF093943.1 Polysiphonia brodiei isolate PD0516 18S ribosomal RNA gene, partial sequence AF427534.1 Polysiphonia nigra small subunit ribosomal RNA gene, complete sequence JX828168.1 Ceramiales sp. HGC-2012 voucher CH1414 small subunit ribosomal RNA gene, partial sequence MF093937.1 Melanothamnus harveyi isolate PD0890 185 ribosomal RNA gene, partial sequence AF427529.1 Polysiphonia elongata small subunit ribosomal RNA gene, complete sequence >AY237287.1 Osmundaria spiralis 18S ribosomal RNA gene, partial sequence phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$





Use 'grep' to print occurrences of a pattern

\$ grep ">" CP blast seqs.fna

Create a new file of the fasta headers

\$ grep ">" CP blast seqs.fna > CP blast headers.txt

phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ grep "^>" CP_blast_seqs.fna > CP_blast_headers.txt phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ head CP_blast_headers.txt >AY617126.1 Choreocolax polysiphoniae 18S ribosomal RNA gene, partial sequence >AY617142.1 Odonthalia washintoniensis 18S ribosomal RNA gene, partial sequence >AY617140.1 Neorhodomela larix 185 ribosomal RNA gene, partial sequence >AY617141.1 Odonthalia floccosa 185 ribosomal RNA gene, partial sequence >MF093958.1 Thaumatella adunca isolate PD1388 185 ribosomal RNA gene, partial sequence >MF093926.1 Digenea simplex isolate PD1820 18S ribosomal RNA gene, partial sequence >JX828176.1 Odonthalia corymbifera voucher OK230 small subunit ribosomal RNA gene, partial sequence >AY617145.1 Rhodomela confervoides 18S ribosomal RNA gene, partial sequence >MF093954.1 Rhodomela confervoides isolate PD508 18S ribosomal RNA gene, partial sequence >L26203.1 Rhodomela confervoides DNA fragment

LP blast seqs.fna

CP blast headers.txt





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

\$ grep -c ">" CP_Blast_seqs.fna

phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ grep -c "^>" CP_blast_seqs.fna 100 phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ grep -c ">" CP_blast_seqs.fna 100 phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$





Quotation marks are important

\$ grep –c ">" CP_blast_seqs.fna

phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ grep -c "^>" CP_blast_seqs.fna 100 phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ grep -c ">" CP_blast_seqs.fna 100 phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$

\$ grep -c > CP_blast_seqs.fna

phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ grep -c > CP_blast_seqs.fna
Usage: grep [OPTION]... PATTERN [FILE]...
Try 'grep __belp' for more information

phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ grep -c ">" CP_blast_seqs.fna

phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ head CP_blast_seqs.fna
phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$

CP_Blast_seqs.fna is empty





Search for headers that are not partial sequences

\$ grep -v "partial" CP_blast_headers.txt

salami:workshop_files ericsalomaki\$ grep -v "partial" CP_blast_headers.txt >L26203.1 Rhodomela confervoides DNA fragment >AF339893.1 Lenormandia latifolia 18S ribosomal RNA gene, complete sequence >AF203897.1 Lenormandia muelleri small subunit ribosomal RNA gene, complete sequence >AF339899.1 Neurymenia fraxinifolia 18S ribosomal RNA gene, complete sequence >AF203889.1 Melanamansia mamillaris small subunit ribosomal RNA gene, complete sequence >AF339896.1 Lenormandia spectabilis 18S ribosomal RNA gene, complete sequence >AF339901.1 Protokuetzingia australasica 18S ribosomal RNA gene, complete sequence >AF339895.1 Lenormandia smithiae 18S ribosomal RNA gene, complete sequence >AF203895.1 Lenormandia prolifera small subunit ribosomal RNA gene, complete sequence >AF339900.1 Osmundaria prolifera 18S ribosomal RNA gene, complete sequence >AF203896.1 Micropeuce strobiliferum small subunit ribosomal RNA gene, complete sequence >AF251513.1 Halopithys incurva 18S ribosomal RNA gene, complete sequence >AF339892.1 Lenormandia angustifolia 18S ribosomal RNA gene, complete sequence >AF203894.1 Laurencia filiformis small subunit ribosomal RNA gene, complete sequence >AF251512.1 Melanamansia glomerata 18S ribosomal RNA gene, complete sequence >AF203887.1 Murrayella periclados small subunit ribosomal RNA gene, complete sequence >AF203886.1 Polysiphonia lanosa small subunit ribosomal RNA gene, complete sequence >AF339898.1 Lenormandia sp. MELU_000064 18S ribosomal RNA gene, complete sequence >AF427537.1 Womersleyella setacea small subunit ribosomal RNA gene, complete sequence >AF339897.1 Lenormandia sp. MELU_000065 18S ribosomal RNA gene, complete sequence >AF427526.1 Boergeseniella fruticulosa small subunit ribosomal RNA gene, complete sequence >AF427527.1 Enelittosiphonia stimpsonii small subunit ribosomal RNA gene, complete sequence >AF203890.1 Heterocladia australis small subunit ribosomal RNA gene, complete sequence >AB219915.1 Polysiphonia sp. SNI07 gene for 18S rRNA, complete sequence >AF427530.1 Polysiphonia fucoides small subunit ribosomal RNA gene, complete sequence >AF203892.1 Heterocladia umbellata small subunit ribosomal RNA gene, complete sequence >AB219910.1 Polysiphonia sp. BRI50 gene for 18S rRNA, complete sequence >AF203885.1 Neosiphonia savatieri small subunit ribosomal RNA gene, complete sequence >AF427534.1 Polysiphonia nigra small subunit ribosomal RNA gene, complete sequence >AF427529.1 Polysiphonia elongata small subunit ribosomal RNA gene, complete sequence

LP_blast_seqs.fna





Replacing Text in Large files

sed 's/**FIND**/**REPLACE**/g' filename > output_file

\$ sed 's/ /_/g' CP_blast_headers.txt

phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ sed 's/ /_/g' CP_blast_headers.txt > CP_blast_headers_clean.txt phylogenomics@ip-172-30-3-100:~/workshop_materials/Blast_Out\$ head CP_blast_headers_clean.txt >AY617126.1_Choreocolax_polysiphoniae_18S_ribosomal_RNA_gene,_partial_sequence >AY617142.1_Odonthalia_washintoniensis_18S_ribosomal_RNA_gene,_partial_sequence >AY617140.1_Neorhodomela_larix_18S_ribosomal_RNA_gene,_partial_sequence >AY617141.1_Odonthalia_floccosa_18S_ribosomal_RNA_gene,_partial_sequence >MF093958.1_Thaumatella_adunca_isolate_PD1388_18S_ribosomal_RNA_gene,_partial_sequence >MF093926.1_Digenea_simplex_isolate_PD1820_18S_ribosomal_RNA_gene,_partial_sequence >XX828176.1_Odonthalia_corymbifera_voucher_OK230_small_subunit_ribosomal_RNA_gene,_partial_sequence >AY607145.1_Rhodomela_confervoides_18S_ribosomal_RNA_gene,_partial_sequence >MF093954.1_Rhodomela_confervoides_isolate_PD508_18S_ribosomal_RNA_gene,_partial_sequence >L26203.1_Rhodomela_confervoides_DNA_fragment

CP_blast_seqs.fna

LP blast seqs.fna





Replacing Text in Large files

Use the -- i flag to make changes inplace

\$ sed -i 's/,//g' CP_blast_headers.txt
\$ head CP_blast _headers.txt

>AY617126.1 Choreocolax polysiphoniae 18S ribosomal RNA gene partial sequence >AY617142.1 Odonthalia washintoniensis 18S ribosomal RNA gene partial sequence >AY617140.1 Neorhodomela larix 18S ribosomal RNA gene partial sequence >AY617141.1 Odonthalia floccosa 18S ribosomal RNA gene partial sequence >MF093958.1 Thaumatella adunca isolate PD1388 18S ribosomal RNA gene partial sequence >MF093926.1 Digenea simplex isolate PD1820 18S ribosomal RNA gene partial sequence >JX828176.1 Odonthalia corymbifera voucher OK230 small subunit ribosomal RNA gene partial sequence >AY617145.1 Rhodomela confervoides 18S ribosomal RNA gene partial sequence >MF093954.1 Rhodomela confervoides 18S ribosomal RNA gene partial sequence >L26203.1 Rhodomela confervoides DNA fragment

CP_blast_seqs.fna

LP_blast_seqs.fna

Regular Expressions

Encoding	Modern Equivalent	Pattern Type
•		a single character
.+		one or more characters
• *		zero or more characters
.?		Maybe present
^		first on the line
\$		last on the line
[0-9]	\d	digits
[a-zA-Z]	\w	letters
	\s \t	space
{3}		must be exactly 3 characters long
{3,5}		between 3-5 characters long
[ACGT]		a specific set of characters (a class)

Regular Expressions

Use –E with sed and grep to use extended regular expressions

Store pattern in memory using parentheses

Print out only the GenBank accessions



[salami:workshop_files ericsalomaki\$ sed -E 's/(>[A-Z0-9.]+)(.+)/\1/' CP_blast_headers.txt
>AY617126.1
>AY617142.1
>AY617140.1
>AY617141.1
>MF093958.1
>MF092026_1

Regular Expressions

Use –E with sed and grep to use extended regular expressions

Store pattern in memory using parentheses

Print out everything except the GenBank accessions



salami:workshop_files ericsalomaki\$ sed -E 's/(>[A-Z0-9.]+)(.+)/\2/' CP_blast_headers.txt Choreocolax polysiphoniae 18S ribosomal RNA gene, partial sequence Odonthalia washintoniensis 18S ribosomal RNA gene, partial sequence Neorhodomela larix 18S ribosomal RNA gene, partial sequence Odonthalia floccosa 18S ribosomal RNA gene, partial sequence Thaumatella adunca isolate PD1388 18S ribosomal RNA gene, partial sequence Digenea simplex isolate PD1820 18S ribosomal RNA gene, partial sequence Odonthalia corymbifera voucher OK230 small subunit ribosomal RNA gene, partial sequence Rhodomela confervoides 18S ribosomal RNA gene, partial sequence

Pipes



Use "|" to string multiple commands together

Combine the fasta files from the blast output and identify how many sequences there are

cat CP_blast_seqs.fna LP_blast_seqs.fna | grep ">" | wc -l

[salami:workshop_files ericsalomaki\$ cat CP_blast_seqs.fna LP_blast_seqs.fna | grep ">" | wc -l 200

Pipes



Use the pipe "|" to string multiple commands together

Combine the fasta files from the blast output and identify how many sequences there are

cat CP_blast_seqs.fna LP_blast_seqs.fna | grep ">" | wc -l

[salami:workshop_files ericsalomaki\$ cat CP_blast_seqs.fna LP_blast_seqs.fna | grep ">" | wc -l 200

Other built-in unix binaries that are great to use with pipes

'sort' – sort lines of text files'uniq' – report or omit repeated lines (only works on a sorted file)



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



Using a Loop



Using for loops allow you to apply a command to multiple files at once

\$ for f in *.fna; do grep -c "ACT" \$f; done

Assign variable "f" to all files that end in .fna Then carry out command on all values of "f"

[salami:Blast_Out ericsalomaki\$ for f in *.fna; do grep -c "ACT" \$f ; done 1354 1291

I often use loops to make several single gene trees at once



Using a Loop



Use a 'for loop' to find how many times you can find "CAT" in these fasta files

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

1257 1319





Online Downloads

wget will download to your pwd Use the –P option

wget -P /path/to/desired/location/ File_To_Download.txt

Online Downloads



Use wget –P to download these slides to a new directory -P will make a new directory if it does not yet exist

wget –P ./wget_dir/ http://evomics.org/wpcontent/uploads/2019/01/Introduction_to_Unix_Krumlov2019_Jan21_red.pdf

http://evomics.org/wp-content/uploads/2019/01/Introduction_to_Unix_Krumlov2019_Jan21_red.pdf

Then use Is to see the new directory and then again to look inside the directory to ensure the slides have been downloaded
Final Unix Challenge



Download the challenge and view it in terminal

http://evomics.org/wp-content/uploads/2019/01/Final_Unix_Challenge.txt

Save Your Data Everyday



We will launch new instance tonight so everything you have done today will be gone tomorrow

Use scp or rsync to transfer files

scp -r username@IP.address:/Location/On/Instance /Location/On/Your/Computer

\$ scp -r phyogenomics@ec2-34-227-31-223.compute-1.amazonaws.com:~/workshop_materials/Unix ~/Location/On/Your/Computer

Any questions?

