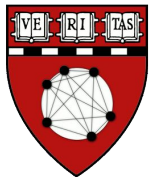# SBGrid RELION Workshop 2017

---

Software installation
Intro to Computing on Linux

Jason Key PhD & Peter Meyer PhD

BCMP Harvard Medical School

SBGrid Consortium

key@hkl.hms.harvard.edu
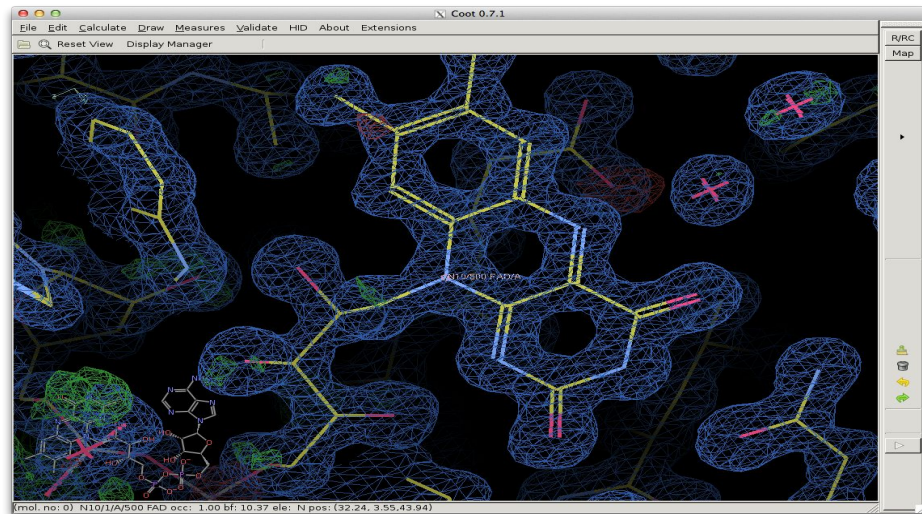meyer@hkl.hms.harvard.edu

# Welcome!

## SBGrid : Structural Biology Research Computing
---

Who we are:

Non-profit Consortium based in
BCMP @HMS focused on Structural
Biology computing

Structural Biologists, IT pros,
software engineers, programmers,
software policy advocates,
postdocs, students

# Welcome!

## SBGrid : Structural Biology Research Computing
___

Pete
Meyer

Mick
Timony

Dimitry
Filonov

Justin
O'Connor

Piotr
Sliz

James
Vincent

Michelle
Ottaviano

Saythyda
Corrado

Carol Herre
Rob DiCamillo

# SBGrid RELION Workshop 2017

---

Software Installation

Intro to Linux Computing

# SBGrid RELION Workshop 2017

We will be using Amazon EC2.
You may want to install RELION and
associated applications:

- RELION

   gctf ctf motioncorr motioncor2 unblur summovie
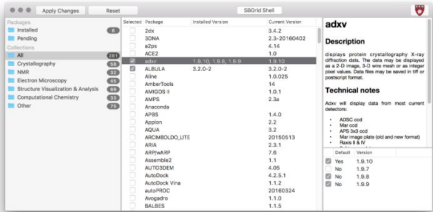
- Chimera

# SBGrid RELION Workshop 2017

## https://sbgrid.org/wiki/client_install

# SBGrid RELION Workshop 2017

---

## Intro to Linux Computing

- ## Introduction to Linux
    Why Linux?

- ## The Linux interface
  (Understanding the Shell, scripting)

- ## Scientific Computing on linux -
    Tips and tools for computing and research

# Linux

---

Linux is an open-source operating systems modeled on UNIX developed by Linux Torvalds in 1991.

Comp.os.minix

*Hello everybody out there using minix - I'm doing a (free) operating system (just a hobby, won't be big and professional like gnu) for 386(486) AT clones. …*

# Linux



---

The GNU (Gnu's Not Unix) was an effort to develop free and open source OS and applications.

Torvalds developed 'kernel' and combined it with software from Richard Stallman @ MIT.



GNU's Not Unix

# Linux

---

**Multi-user, multi-tasking**

Many users on the same machine at once, running many programs

**Multi-platform**

runs on many different processor types

# Linux

___

**Multi-user, multi-tasking**

   Many users on the same machine at once, running many programs

**Multi-platform**

   runs on many different processor types

**Linux is a Unix-like system free of proprietary software for which source code is available and freely distributed**

# Linux

---

**Multi-user, multi-tasking**

Many users on the same machine at once, running many programs

**Multi-platform**

runs on many different processor types

**Linux is a Unix-like system free of proprietary software for which source code is available and freely distributed**

# Linux

___

But why did Linux succeed?
  ( and not HURD, BSD, Minux, etc …)

- Decentralized Development

- Pragmatic (Not an academic or ideological exercise)

- Technological Superiority

- Luck?

# Linux

___

The Filesystem Hierarchy Standard (FHS) defines the directory structure and directory contents.

"Everything is a file"

Defining features of *nix

Resources (documents, directories, keyboards, printers, storage, network communications, etc)

Are simple streams of bytes exposed through the filesystem name space

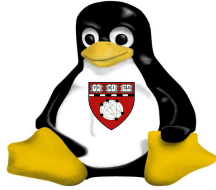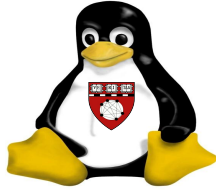| ROOT DIRECTORY OF THE ENTIRE FILE SYSTEM HIERARCHY / PRIMARY HIERARCHY | | |
|---|---|---|
| /bin/ | ESSENTIAL USER COMMAND BINARIES | |
| /boot/ | STATIC FILES OF THE BOOT LOADER | |
| /dev/ | DEVICE FILES | |
| /etc/ | HOST-SPECIFIC SYSTEM CONFIGURATION REQUIRED DIRECTORIES: OPT, X11, SGML, XML | |
| /home/ | USER HOME DIRECTORIES | /home/student/ → /home/student/dir |
| /lib/ | ESSENTIAL SHARED LIBRARIES AND KERNEL MODULES | /home/linuxgym |
| /media/ | MOUNT POINT FOR REMOVABLE MEDIA | |
| /mnt/ | MOUNT POINT FOR A TEMPORARILY MOUNTED FILESYSTEMS | |
| /opt/ | ADD-ON APPLICATION SOFTWARE PACKAGES | FILESYSTEM HIERARCHY STANDARD ( FHS ) |
| /sbin/ | SYSTEM BINARIES | |
| /srv/ | DATA FOR SERVICES PROVIDED BY THIS SYSTEM | |
| /tmp/ | TEMPORARY FILES | /usr/local/bin |
| /usr/ | (MULTI-)USER UTILITIES AND APPLICATIONS SECONDARY HIERARCHY REQUIRED DIRECTORIES: BIN, INCLUDE, LIB, LOCAL, SBIN, SHARE | /usr/local |
| /var/ | VARIABLE FILES | /usr/local/games |
| /root/ | HOME DIRECTORY FOR THE ROOT USER | |
| /proc/ | VIRTUAL FILESYSTEM DOCUMENTING KERNEL AND PROCESS STATUS AS TEXT FILES | LINUXCONFIG.ORG |

# Hardware and Workflows

# The Shell

———

## The 'shell' is the Command Line Interface for Linux

This is an program that interprets what you type, keeps track of programs on the system, etc.

## Common Shells:

**tcsh** : exTended C SHell
**bash** : Bourne Again SHell
**ksh**  : Korn SHell
**csh**  : C SHell (early popular shell)
**sh**   : the original shell, often a synonym for bash now

# The Shell : Commands

———

Import / include (source)

Navigation (cd, pwd, ls)

Manipulating Files (cp, rm, mv, mkdir)

Search (grep, find)

Permissions (chmod, chown, chgrp)

Job Control (jobs, ps, fg, bg, nohup)

http://linuxcommand.org/lc3_learning_the_shell.php

# The Shell : The Environment

— — —

The shell environment is configured globally per user in files and startup scripts

- Settings for variables
- Function definitions
- Aliases

Except for the reserved Shell special parameters variable names can be set by the user

Quotes remove special meaning from one or multiple characters

# The Shell : The Environment

___
The 'printenv' command

# The Shell : The Environment

___

The 'printenv' command

Shell variables:

     PATH
        Where executables can be found
     HOME
        User's home directory
     USER
        User's username
     SHELL
        Default shell setting

# The Shell : The Environment

———

The 'printenv' command

Shell variables:

   PS1
      Shell prompt settings

   LD_LIBRARY_PATH
      Primary search path for library directories
 ...

# The Shell : The Environment

———
The 'alias' command

An alias is a shortcut or abbreviation.

Great for avoiding typing a long command sequences

Aliases do NOT get passed to scripts (sub-shells)

# The Shell : The Environment

———
Functions: The 'declare -f' command


Functions are subroutines :
a code block (list of commands) that implements a set
of operations.

# The Shell : stdout stdin stderr

———

stdin :
    Input for commands
    usually come from the keyboard

stdout :
    Output from commands
    written to the screen

stderr :
    Error messages from processes
    usually written to the screen

# The Shell : stdout stdin stderr

———

Pipe (|):
    stdout of one command to stdin of another command

Output Redirection (>):
    stdout of a command to a file

Output Appending (>>):
    stdout of a command appending to a file

Input Redirection (<):
    stdin of a command from a file

    Use "-" to read this from standard input

# The Shell : stdout stdin stderr
___


Stderr redirection

For tcsh
    &> filename

For bash
    2>&1 filename

# The Shell : stdout stdin stderr

———

Most Linux (*nix) commands can be strung together

Example:
    How many image files do I have?

```
ls -l *img | wc
```

    How many image files do I have that are not have 'native' in the name?

```
ls -l *img | grep -v "native" | wc
```

# The Shell : stdout stdin stderr

---

Most Linux (*nix) commands can be strung together

Example:
  A list of all my image files :

```
ls *img > my_images.txt
```

  A list of all my images sorted in reverse numerical order?

```
ls -l *img | sort -rn -k 9 > sorted_files.txt
```

# The Shell script

———

Shell scripts are text files of variables, functions, and commands

A 'shebang' (#!/bin/bash, …) is required to indicate which interpreter the OS programs loader should use

Conditional expressions: if/else, case
Loops: for, while, until

http://tldp.org/HOWTO/Bash-Prog-Intro-HOWTO.html

# Scripting and Scientific Computing

———

Why *wouldn't* you want to script data processing?

- Dealing with intrinsically visual data (density interpretation, particle picking, etc)
- Using programs that are GUI only

# Scripting and Scientific Computing

___

Why would you want to script data processing?

- Documenting *how* a dataset was processed, and *why* particular options were used
- Easier to process several datasets identically (e.g. - comparing apo structure and complex structure)
- Easier to explore alternative ways to process (e.g. - MR/density fitting for 50 models

# Scripting and Scientific Computing

———

```bash
#!/bin/bash



process_my_data input*.mrc output.mrc
```

# Scripting and Scientific Computing

---

```bash
#!/usr/bin/env bash



# first stage of my data processing



process_my_data input*.mrc output.mrc > processing.log
```

# Scripting and Scientific Computing

———

```bash
#!/usr/bin/env bash
# first stage of my data processing
job=ProcessingStage01
inp="input*.mrc"
opf="${job}-output.mrc"
process_my_data $inp $opf > ${job}.log
```

# Scripting and Scientific Computing

— — —

```bash
#!/usr/bin/env bash
# first stage of my data processing
#tuning parameter
tune_param="0.2"
# pipeline flag
flag="EvalReconstruct"
job=ProcessingStage01
inp="input*.mrc"
opf="${job}-output.mrc"
process_my_data --input $inp --output $opf --tune $tune_param << eof > ${job}.log
PIPELINE_FLAG $flag
eof
```

# Scripting and Scientific Computing

— — —

```bash
#!/usr/bin/env bash

job=TestScan

input_dir="models/"
output_dir="results/"
map=input.mrc

for model in `ls $input_dir/*.pdb | awk -F. '{print $1}'`
do
    search_density_for_model --map $map --search_model ${input_dir}/${model}.pdb --output
${output_dir}/${model}_results.out
done
```

# Scientific Computing: tools, tips and tricks
---

- Getting there and moving data

- What resources does this computer have (and what is it doing)?

- Reproducibility and collaboration

# Scientific Computing: tools, tips and tricks

## SSH

provides a secure channel (encrypted) over an unsecured network in a client-server architecture

- remote command-line login
- remote command execution
- any network service can be secured with SSH.

# Scientific Computing: SSH

———

SSH public-key authentication allows login and command execution without passwords based on a public/private key pair

Setup: create keys, set a password

```
ssh-keygen -t rsa
```

Public key goes on the remote server in your .ssh directory in the file $HOME/.ssh/authorized_keys

Private key stays in $HOME/.ssh

Ssh-agent manages keys - typically running by default on most Linuxes

Use ssh-add command to add key, **No more passwords**!

# Scientific Computing: SSH

———

SSH public-key authentication allows login and command execution without passwords based on a public/private key pair

Use ssh -X to forward X11 for graphics access

Execute code remotely with a single command

# Scientific Computing: Moving data with rsync

———

## rsync

- Transfers only changes in a file tree
- Local and remote synchronization of data file and directories

```
rsync -rv /my/files/here/ /my/files/over_there/
```

- Can be run over ssh for secure transfer

```
rsync -rv /my/files/here/ remote.server.org:/my/files/over_there/
```

- Ideal for data backup

```
man rsync
```

For more info

# Scientific Computing: Hardware

---

CPU, Storage,memory,usb and pci

- CPU : cpuinfo or cat /proc/cpu
- DISK : df or lsblk
- MEMORY : free
- USB : lsusb
- PCI (internal cards, etc) : lspci

# Scientific Computing: What is running

---

CPU and memory use, jobs, IO

- top, uptime
- ps
- sar

# Scientific Computing: history

`___`

# history

    The shell records all commands.
This record can be accessed with the 'history' command.


Some relevant variables:
    HISTSIZE
        Define number of commands

    HISTFILE
        Define file

    HISTCONTROL=ignoredups
        Ignore duplicates

# Scientific Computing: Terminal multiplexer

———

## Tmux (or screen)

A terminal multiplexer is terminal-based program that gives the user

- Ability to detach and reattach sessions from a terminal
  - Sessions persist on the remote machine
  - A terminal session can be accessed from multiple machines
  - Persist through network disconnection

- Multiple separate login sessions inside a single terminal window

# Scientific Computing: Version Control

---

## VCS

Version control systems are designed for software development, are are great for scientific computing projects



Version control software keeps track of every modification to the code

Earlier versions of code are retained and can be accessed

# Scientific Computing: Getting started with git

---

More GIT
Gitlab, github bitbucket, RELION is in GIT

https://git-scm.com/book/en/v1/Getting-Started

# Scientific Computing: Install Client CLI

# Scientific Computing: Install client GUI