

‘Large’ Data Methods-An Introduction to Scalable Statistical Genomics in Linux

John McCallum

john.mccallum@plantandfood.co.nz

Marcus Davy

Samantha Baldwin

Statistical Genomics in 2012+

- Larger is getting cheaper and larger
- Context-specific SNP marker discovery
- Population pool methods
- Genotyping by sequencing
- Global methods
- Reference-based methods
- Largeness requires scalable server-based computing *cloud, cluster etc*
- Diverse tools, scripts available for UNIX (ie OSX, Linux)

Learning Goals

- **To stop you being frightened by Unix, and see it as an highly accessible and powerful science tool**
- Introduce you to a modern linux desktop
- Unix essentials
- Accessing servers and moving files
- Documentation, help and formats
- Getting and using third-party scripts and executables
- Using bwa read mapper+samtools for variant analysis
- Visualization
- R and Galaxy interfaces to same tools

Why Unix?

- It makes the web and our phones work
- OSX and current Linux desktops match or exceed Windows functionality
- Leading platform for scientific computing
- Secure
- Scalable *phone/desktop/server/cluster/cloud*
- Many flavours to suit your needs

*freedom to work how you want, where you want, with
(b)leading-edge tools*

Outline

An intro to the Unix CL -Exploring data and formats (1 h+)

Running Analyses -Variant detection & analysis(1 h)

Visualization with R and IGV (1 h)

Formats and data manipulation in Galaxy (1/2 h)

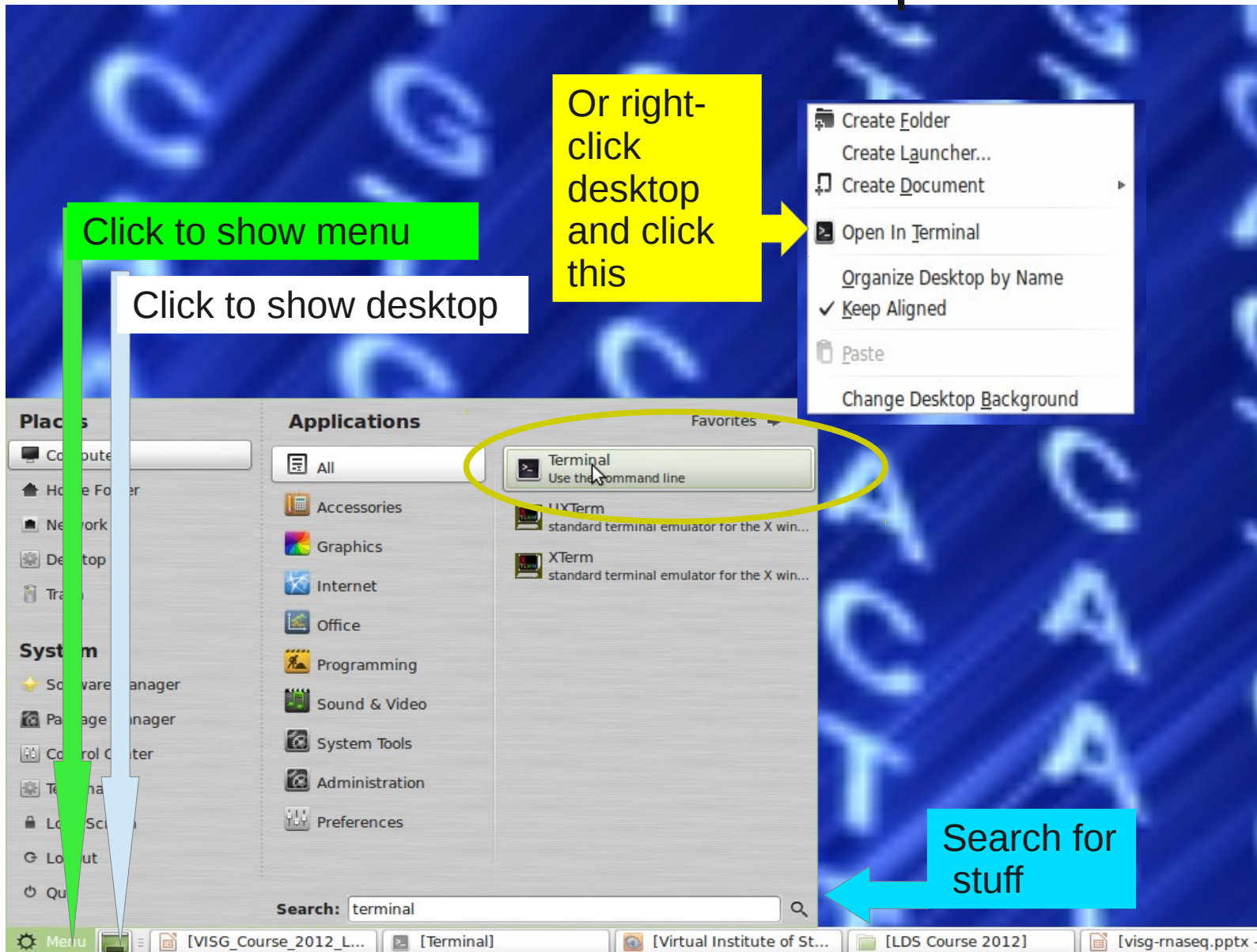
The Data

- Fragmented, barcoded PCR amplicons from flowering candidate genes from 7 populations, 2x 1/16 454 Ti plate
- Parallel-Tagged Sequencing Meyer et al 2008 -now rendered obsolete but basically same workflow from eg Nextera XT
- Raw data
 - Genomic Reference fasta sequence (homozygous reference)
 - Annotation of reference sequences from gmap (gff3)
 - Fastq files, one per population per plate segment
- Workflow
 - Read mapping with BWA SW
 - Manipulation and SNP calling with Samtools
 - Population genetic analyses with PoPoolation2
 - Visualization with IGV and R (reshape/ggplot2)

Biological Question

Which SNPs in these genes
show evidence of strong
population differentiation?

Accessing Terminal from Debian Linux Mint Desktop



> The bash shell

- 'Bourne-again-shell'
- A command-line (CL) interface to operating system
- a command interpreter
- **Command** **-option** *<value>* **argument(s)**
- Inputs and outputs from files or stdin/out

<http://manuals.bioinformatics.ucr.edu/home/linux-basics>

Exercise-Shell Orientation

whoami

Who are you?

pwd

Where are you?

ls

List the files

ls -l

Long listing

cd /

Go to root directory

ls

List the files

ls -la ~

Long listing of home plus hidden files

cd

Go Home

ls -l .

List files in this dir

ls ..

List the files in parent dir

Get the Workflow and Data

- **Browse:** <https://github.com/cfljam/VISG-course-2012>

- `git clone https://github.com/cfljam/VISG-course-2012.git`
`## get the archive`

- `ls -l ##get a directory listing`

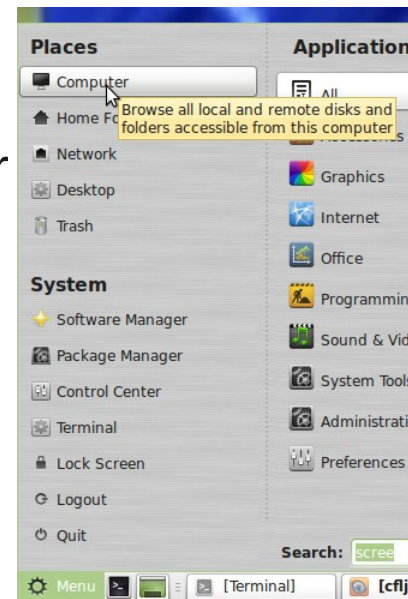
- `ls -R * ##recursively list directory contents`

- `tar -tvf VISG_LDS.tgz ##check the archive`


- `tar -xvf VISG_LDS.tgz ##unpack it`

- `ls -l ls -l 00.raw/ ##check there are data files`

- Browse the directories and archives with the file manager



Security==Permissions

- Permissions protect you and the system from you and others
-  your sysadmin :-)
- If you must....on this install
 - **sudo su ##to become admin**
 - **sudo <do-something requiring admin rights> ## one-off**
- Each file/directory has
 - Owner
 - Group membership
 - Permissions
 -

Permissions links owner group size date name

```
-rw-r--r-- 1 visg_user visg      875 Oct 17 11:36 README.md
drwxr-xr-x 2 visg_user visg    4096 Oct 17 11:36 supplementary_QC
drwxr-xr-x 2 visg_user visg    4096 Oct 17 11:36 supplementary_samtools_usage
-rw-r--r-- 1 visg_user visg 22160684 Oct 17 11:37 VISG_LDS.tgz
```

Help!

`help` #list shell commands

`help cd` #help for cd command

`man ls` #read the man pages for ls, q to exit

`<command>` #may give help eg `bwa`

`<command> --help` # for some programs

Google..

Exercise CL Navigation

`cd VISG-course-2012/00.raw`

`ls P<tab>`

`<tab> <tab>`

`ls Pool1_BARCODE*`

`ls Pool1_BARCODE?.fastq`

`history`

`history | tail`

`<up arrow>`

`<down arrow>`

`< mouse double click>`

`<shift ins>`

Move to

List with *filename completion*

All the options

List the Pool1 files

List the Pool1 fastq files

see all the history

Last few items

Back in history

Forward in history

Copy

paste

Exercise-Explore

`ls *.fastq` *#list to stdout*

`ls *.fastq > somefile` *#redirect to file*

`cat somefile` *#to stdout*

`cat somefile | head` *# pipe file to head*

`cat > somefile` *#read from stdin/ctrl d save*

In Unix...everything is a File

Regular files-human readable text

Directories

Executable files

Compiled

Special text files

Symbolic links -'shortcuts'

Exercise -View, Browse and Search

- › `cd /VISG/00.raw` ##move to raw data dir
- › `head Pool1_BARCODE2.fastq` ##see top of file
- › `tail Pool1_BARCODE2.fastq` ##see bottom of file
- › `less Pool1_BARCODE2.fastq` ##view with the less pager
 - **h** help screen
 - **g** *top of file*
 - **G** bottom of file
 - **/***<pattern> search for pattern*
 - **q** quit less
- › `grep @GYSS Pool1_BARCODE2.fastq | head` ##get readnames
- › `grep -c @GYSS Pool1_BARCODE2.fastq` ##count reads

Gotchas-Symbols, Whitespace , Names

- Stick to `A-Za-z0-9_` for naming files
- Non-printing characters
 - Spaces and tabs
 - Line endings: Unix=LF, Win =CR/LF
- In shell environment many characters have special meaning e.g..

`#` comment

`#!` *shebang*

`>` *redirect to*

`<` *input from*

`|` *pipe*

`$` variable expression

`/` path delimiter

`\` quote next character

`"` strong quote

`“”` weak quote

`` `` evaluate

Formats

- Input and outputs should stick to standard common formats
- Read the specifications!!!!
- Fasta – Raw nucleotide/peptide format
- Fastq – Raw sequence information + quality
- Sam/Bam format – Sequence Alignment/Map format
- GFF – General feature format-annotations
- Tools for format conversion & filtering
 - Unix tr, awk, sed, perl
 - Programming Libraries Python, Perl, R etc
 - Galaxy

Exercise-SSH and SCP

```
ifconfig | grep 'inet addr'
```

```
ping <their IP address>
```

```
ssh visg_user@<IP address>
```

```
scp visg_user@<IP  
address>:/VISG/00.raw/refer  
ence.fasta ~
```

- Get your IP address, swap with a partner
- Check you can reach their IP address
- SSH to each others machine as VISG_USER
- Copy a file to your home dir

Fasta format

- Text-based format for storing nucleotide/peptide sequence(s)
- Restricted to IUPAC *alphabet* letters

No spaces!

```
>gi|63055|emb|V00385.1| Part of the chicken ovalbumin X gene  
ACTGTGTCTTAGCACTCACTGCTTTGCTTCCTTCTTACAGGACAGATCAAAGATTTGCTTGTATCAAGCT  
CCACTGATCTTGATACAACGCTGGTCCTTGTTAATGCCATCTACTTCAAAGGGATGTGGAAGACAGCATT  
TAATGCAGAAGACACTCGAGAAATGCCCTTCCATGTAACAAAGGTAGGGGACGTAGTCACCGCTTCTGGG  
...
```



Header



Content

*Newline wrap usually at 60 - 80
characters*

http://en.wikipedia.org/wiki/FASTA_format

Fastq format

```
@HWUSI-EAS582_157:6:1:1:1501/1
NCACAGACACACACGAACACACAAAGACATGCCCATATGAAGAT
+
%.7786867:778556858746575058873/347777476035
@HWUSI-EAS582_157:6:1:1:1606/1
NCTGGCACCTTGATTTTGGACTTCCCAGCCTCCAGAACTGTGAG
+
%1948988888798988366898888648998788898888588
@HWUSI-EAS582_157:6:1:1:453/1
NCTGCTTGCACCCCTGAAGTCACTGATCACATTTCAAGGTCACC
+
%/868998988888867668888986644788988413488885
@HWUSI-EAS582_157:6:1:1:1844/1
NGATTGACATTGGCAAAGAGGACAACTGATTGCAAACCTTCACAC
+
%-7;:::;;86499;75574586::635:62687666887879
@HWUSI-EAS582_157:6:1:1:1707/1
```

<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2847217/> (Detail)
http://wiki.genomequest.com/index.php/NGS_Reads

Fastq format header

```
@HWI-EAS209_0025_FC427:6:1:1041:14884#ACAGTG/2
AATTTGTTTGTGTTGTTTATTTTTTTGTTAGTTTCGTTTGTGTTTGGATTCCTCTGTGTTGAGTATTT
+HWI-EAS209_0025_FC427:6:1:1041:14884#ACAGTG/2
_____QZNUSUISNQW__U^BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
```

*Header
Sequence
Header
Quality*

- Illumina header contains several fields

HWI-EAS209	Unique machine identifier
0025	Run number
FC427	Unique flowcell identifier
6	Lane number
1	Tile number
1041	X coordinate within tile
14884	Y coordinate within tile
#ACAGTG	illumina barcode multiplexing index tag
/2	Pair number (1 or 2)

Sequence Alignment Map (SAM) format

- sam – text format
- Bam – binary version

```
- @HD VN:1.3 S0:coordinate
- @SQ SN:ref LN:45
- r001 163 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
- r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
- r003 0 ref 9 30 5H6M * 0 0 AGCTAA * NM:i:1
- r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
- r003 16 ref 29 30 6H5M * 0 0 TAGGC * NM:i:0
- r001 83 ref 37 30 9M = 7 -39 CAGCGCCAT *
```

<http://samtools.sourceforge.net/SAM1.pdf>

<http://en.wikipedia.org/wiki/SAMtools>

<http://samtools.sourceforge.net/>

<http://samtools.sourceforge.net/samtools.shtml>

Fastq format quality

- *Quality scores encoded as ASCII characters*
- *Format has been evolving*

[illegible]

```
. . . . . XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX . . . . .
```

[illegible][illegible]

```

LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.....
.....

```


```
!"#$%&'()*+,-./0123456789:;<=>?  
@ABCDEFGHIJKLMNPOQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
```

Age Group	Gender	Should take action (%)	Should not take action (%)
18-29	Male	~75	~25
18-29	Female	~85	~15
30-49	Male	~65	~35
30-49	Female	~75	~25
50-69	Male	~55	~45
50-69	Female	~65	~35
70+	Male	~45	~55
70+	Female	~55	~45

Sequence alignment map format

- sam – text format
- Bam – binary version

```
• @HD VN:1.3 SO:coordinate
• @SQ SN:ref LN:45
• r001 163 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
• r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
• r003 0 ref 9 30 5H6M * 0 0 AGCTAA * NM:i:1
• r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
• r003 16 ref 29 30 6H5M * 0 0 TAGGC * NM:i:0
• r001 83 ref 37 30 9M = 7 -39 CAGCGCCAT *
```

 Header

<http://samtools.sourceforge.net/SAM1.pdf>

<http://en.wikipedia.org/wiki/SAMtools>

<http://samtools.sourceforge.net/>

<http://samtools.sourceforge.net/samtools.shtml>

Alignments

VISG

Virtual Institute of Statistical Genetics

GFF3 Format

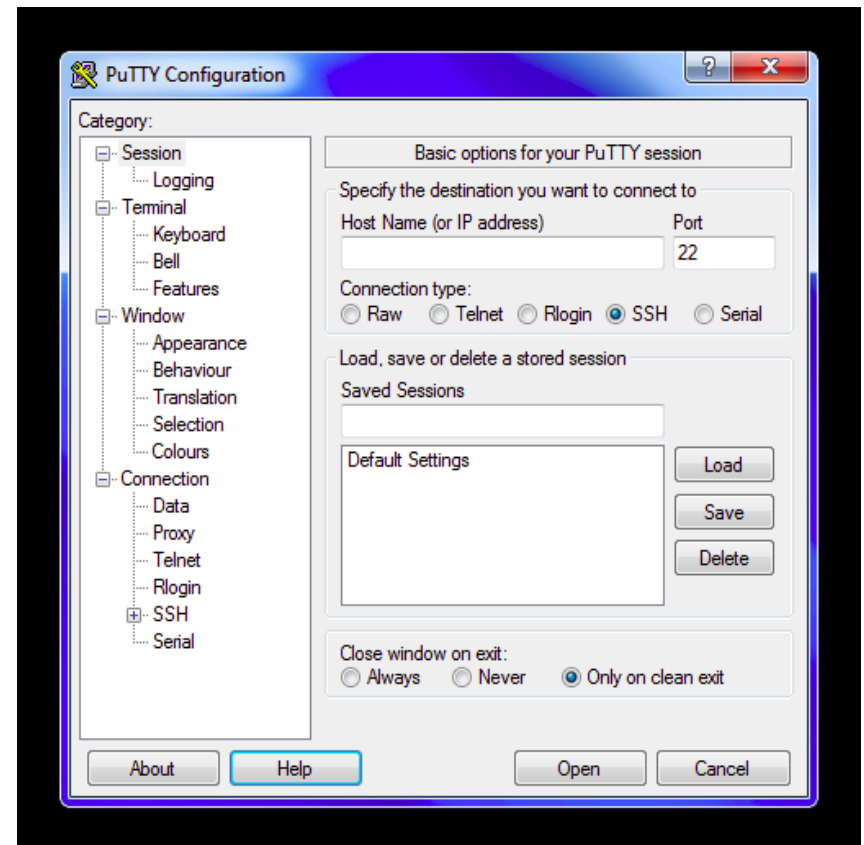
- Generic feature format
- <http://www.sequenceontology.org/gff3.shtml>
- Stream format-one liners of information about a feature

CO_Pool1_contig00004 720;Gap=M221	reference	cDNA_match	355	575	100	-	.	ID=JR851107.path1;Name=JR851107;Target=JR851107 500
SOC1_Pool2_contig00003 242;Gap=M241	reference	cDNA_match	325	565	100	+	.	ID=JR848637.path1;Name=JR848637;Target=JR848637 2
SOC1_Pool2_contig00003 325;Gap=M42 I1 M40	reference	cDNA_match	1320	1401	98	+	.	ID=JR848637.path1;Name=JR848637;Target=JR848637 243
VIN3-like_Pool2_contig00005 613;Gap=M613	reference	cDNA_match	321	933	100	+	.	ID=JR853510.path1;Name=JR853510;Target=JR853510 1
SOC1_Pool2_contig00003 242;Gap=M241	reference	cDNA_match	325	565	100	+	.	ID=JR848637.path1;Name=JR848637;Target=JR848637 2

•

Important Freeware Tools

- Linux -Live CDs/DVDs/USBs-use as installers
- All platforms-Oracle VirtualBox
- OSX
 - Terminal
- Windows
 - Putty
 - Xming
 - Wincp
 - Notepad++



Feeling Overwhelmed Yet?

To keep sane, use
approaches that are

- Scalable
- Open source
- Reproducible
- Documented
- Identifiable
- Disciplined

Reproducibility Questions

Where did these files come from?

What commands and options did I use?

What was I thinking?

How can I re-use this pipeline?

<http://reproducibleresearch.net>

<http://cran.r-project.org/web/views/ReproducibleResearch.html>

Reproducibility -A Simple Approach

- Use one directory per atomic step, with an informative name
- Prefix directory names with numeric order
- Keep filenames consistent and informative
- Paste step commands into an executable shell script file that will enable re-creation
- Document stuff in
 - In-line comments *##some comment*
 - Plain text README , with formatting in [Markdown](#) if desired
- Version control using [git](#)

Scripts and Executables

- `echo $PATH` ##where to look for executables
- `which bwa` ##where is the bwa prog?
- `ls -l /usr/bin/bwa` ##note x in the permissions
- `cd 05.reference/` ##move into a dir with run.sh
- `ls -l run.sh` ##note x in the permissions
- `cat run.sh` ##note shebang, denoting sh ie bash as the interpreter
- Important script interpreters
 - sh (normally bash)
 - Rscript (R)
 - Python
 - Perl
 -

Scripts, Executable Files and Where They Live

visg@mint ~ \$ echo \$PATH *where it looks for executables*

/usr/local/bin:/usr/bin:/bin:/usr/local/games:/usr/games

visg_user@mint ~/visg/00.raw \$ which bwa *where's bwa?*

/usr/bin/bwa

visg_user@mint ~/visg/00.raw \$ ls -l /usr/bin/bwa

-rwxr-xr-x 1 root root 276124 Dec 12 2011 /usr/bin/bwa *executable*

visg_user@mint ~/visg/00.raw \$ cat run.sh *dump file*

#!/bin/sh *uses the sh interpreter ie bash*

Run qa.R script as a batch file *a comment*

Rscript qc.R 2> err > log

visg_user@mint ~/visg/00.raw \$ head qc.R *look at top of script*

#!/bin/env Rscript *uses the Rscript interpreter*

require(ShortRead)

Exercise-Make a Shell Script

- `cd ~`
- `mkdir test`
- `cd test`
- `cat > hello_unix.sh`

```
#!/bin/sh
```

```
echo "hello "
```

```
whoami
```

```
echo "number of lines in file  
listing is:"
```

```
ls -l .. | wc -l
```

```
<ctrl-d>
```

- `ls -l`

- ▢ • Move to HOME
 - Make a dir
 - Move into it
 - Redirect to file (or use editor)
 - Enter each line, then return
 - Ctrl-d to finish
-
- Check you have created a file, and its permissions

Exercise-Run/edit a Shell Script

```
cat hello_unix.sh
sh hello_unix.sh
./hello_unix.sh
chmod +x hello_unix.sh
./hello_unix.sh
cat >> hello_unix.sh
echo "another command"
<ctrl-d>
nano hello_unix.sh
cd ..
rm -r test
```

- View the contents
 - Run using sh
 - Wont work
 - Make it executable
 - Should work
 - Append to the file
-
- edit using nano
 - Move up a level
 - Delete the directory

BWA Aligner

BWA= Burrows-Wheeler Aligner

Produces gapped alignment to reference

<http://bio-bwa.sourceforge.net/>

BWA-SW for reads > 200 bp

Need to index reference first

Produces output in SAM format

<http://samtools.sourceforge.net/>

Exercise-Running BWA

```
which bwa
```

```
bwa
```

```
bwa bwasw
```

```
bwa index <file>
```

```
bwa
```

```
bwasw ../05.reference
```

```
/pool1.fasta ../00.ra
```

```
w/Pool1_BARCODE2.fast
```

```
q | head
```

See where it is

Read the options

Read bwasw
options

Index reference
file

Pipe output into
head

Popoolation

- a collection of tools to facilitate population genetic studies of next generation sequencing data from pooled individuals
- Popoolation
 - A pipeline for analyzing pooled next generation sequencing data for single populations.
 - Tajima's π , Watterson's θ and Tajima's D
 - <http://code.google.com/p/popoolation/>
- Popoolation2
 - Allows analyzing the population frequencies of SNPs from two or more populations.
 - F_{st} , Fisher's exact test, Cochran-Mantel-Haenszel test
 - <http://code.google.com/p/popoolation2/>

Getting PoPoolation

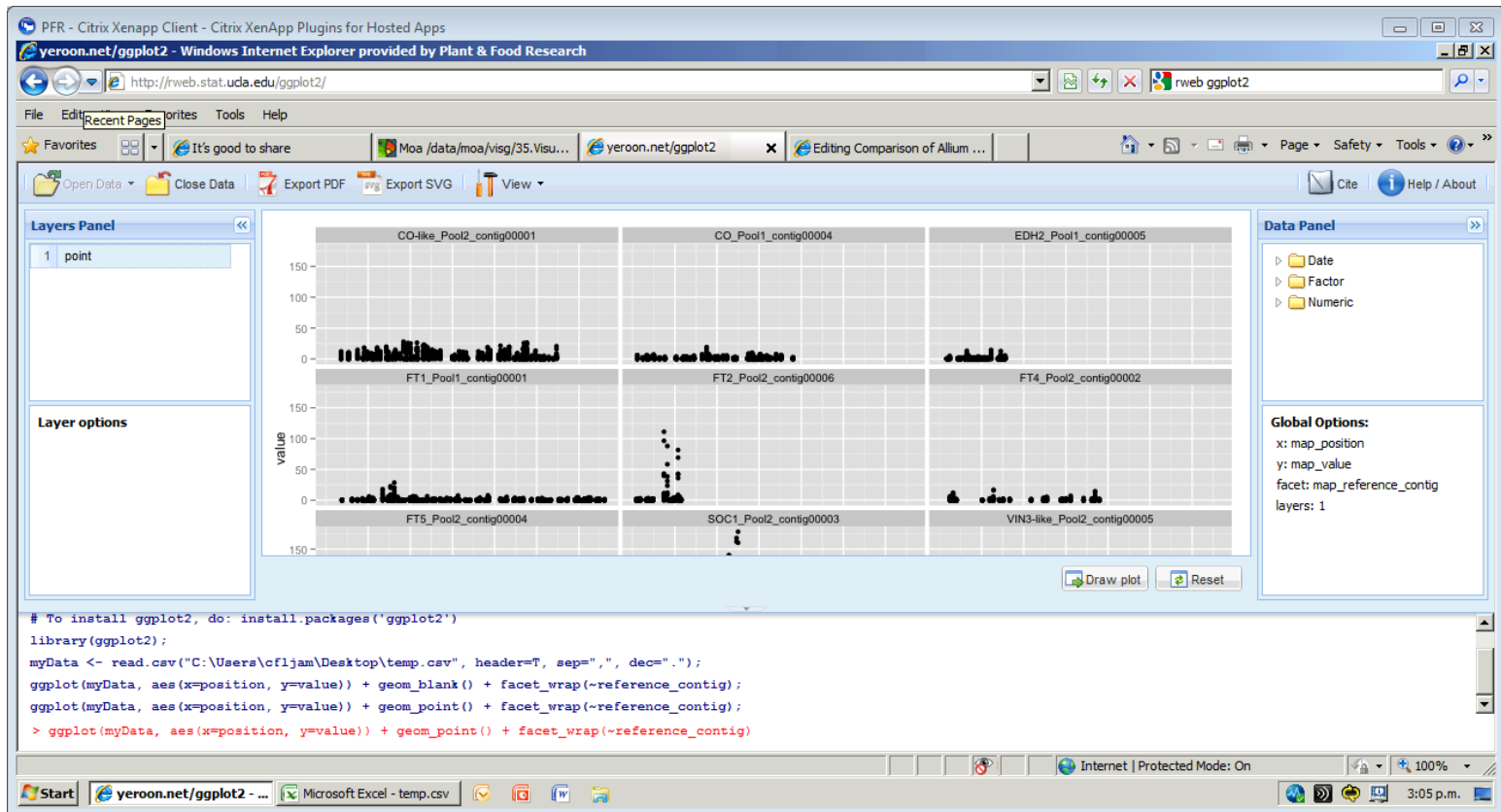
- Browse, download, unpack with archive manager
- or....from CL
- wget http://popoolation2.googlecode.com/files/popoolation2_1201.zip
- unzip popoolation2_1201.zip
- or.....from CL
 - apt-get update ##update package information
 - apt-get install svn ##install SVN
 - svn checkout <http://popoolation2.googlecode.com/svn/trunk/popoolation2> ##check out copy

reShape2

- Flexible rshaping of data
- Especially valuable for turning 'wide' into 'long' (stream-oriented) data
- <http://had.co.nz/reshape/>
- <http://www.jstatsoft.org/v21/i12>

ggplot2

<http://www.stat.ucla.edu/~jeroen/ggplot2/>



<http://had.co.nz/ggplot2/> <http://docs.ggplot2.org/current/>

Getting & Compiling Software-Github

```
cd ~/Downloads/  
git clone https://github.com/lh3/wgsim.git  
cd wgsim  
less README ##read the instructions  
gcc -g -O2 -Wall -o wgsim wgsim.c -lz -lm  
#compile  
echo $PATH ##check your path  
cp wgsim /usr/local/sbin ##copy to PATH  
wgsim ##check it works, read help
```

Plant & Animal Genome Conference ASIA 2013

March 17-19, 2013 SINGAPORE

www.intlpagasia.org

