# Bioinformatics tutorials Documentation

*Release 1*

**Jiarong**

**May 18, 2017**

# Contents

Contents:

# Welcome to 2014 MSU QB summer microbial genomics tutorial

| Day | Schedule |
|-----|----------|
| Wed 8/13 | <ul><li>2-3pm: *Get started with Amazon EC2*</li><li>3-5pm: *Read mapping and variant calling*</li></ul> |
| Thu 8/14 | <ul><li>2-5pm: *Genome assembly*</li></ul> |

## Sequencing technologes

- 454
- Illumina
- PacBio
- NanoPore

## Resources

- Biostar

  A high quality question & answer Web site.

- SEQanswers

  A discussion and information site for next-generation sequencing.

- Software Carpentry lessons

  A large number of open and reusable tutorials on the shell, programming, version control, etc.

# Table Of Contents

## Amazon Web Services instructions

### Start up an EC2 instance

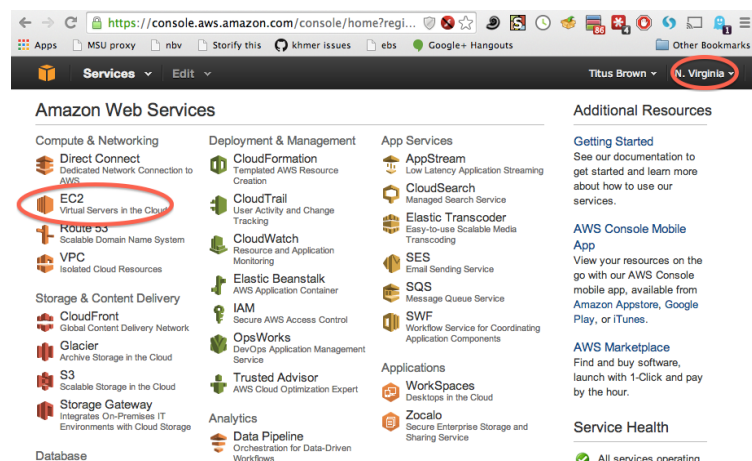Here, we're going to startup an Amazon Web Services (AWS) Elastic Cloud Computing (EC2) "instance", or computer.

---

Go to 'https://aws.amazon.com' in a Web browser.

Select 'My Account/Console' menu option 'AWS Management Console."
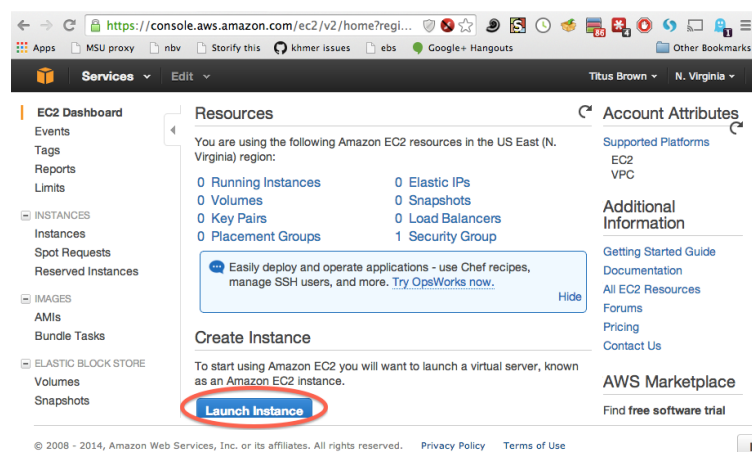
Log in with the following:

- **username: qb2014msu@gmail.com**
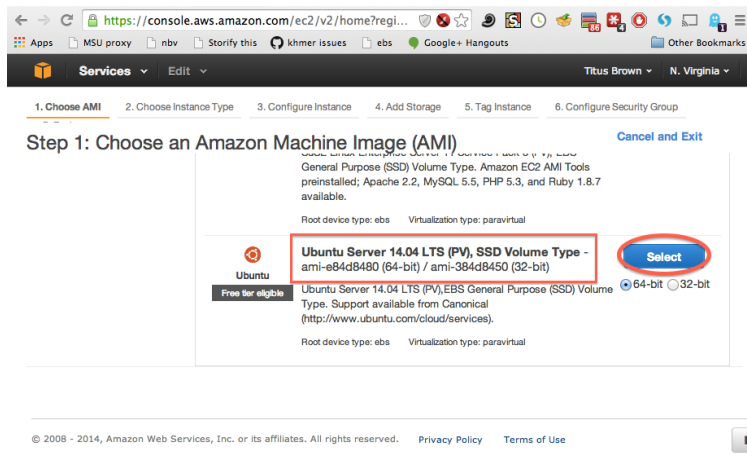
- **password: TemporaryQB2014**

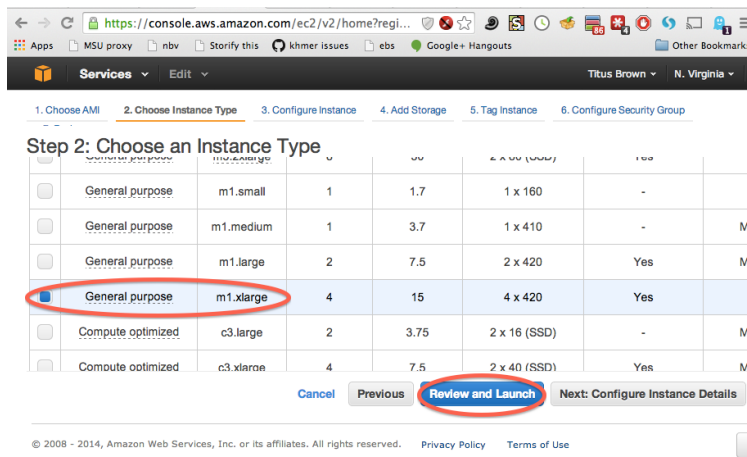Make sure it says North Virginia in the upper right, then select EC2 (upper left).



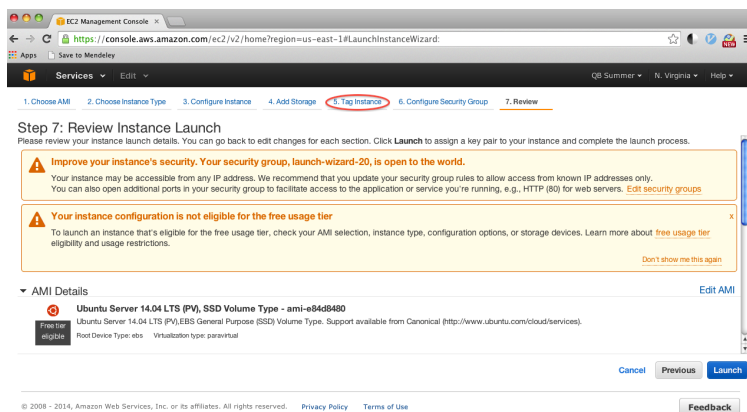Select "Launch Instance" (midway down the page).



Next, scroll down the list of operating system types until you find Ubuntu 14.04 LTS (PV) – it should be at the very bottom. Click 'select'. (See *Starting up a custom operating system* if you want to start up a custom operating system instead of Ubuntu 14.04.)
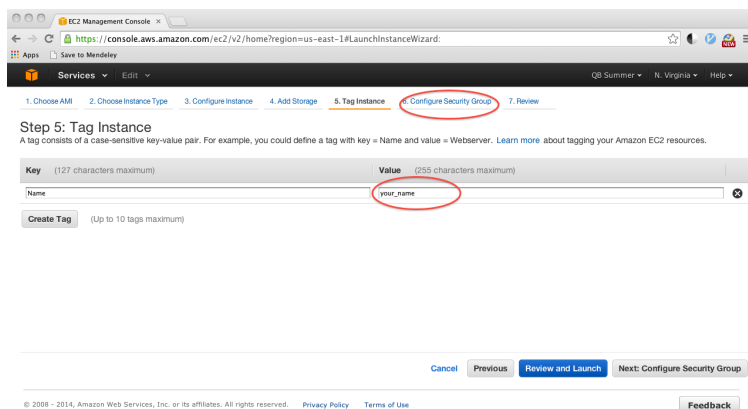
---

Scroll down the list of instance types until you find "m1.xlarge". Select the box to the left, and then click "Review and Launch."



Before the launch, we need to name our intances, so we can tell ours from others'. Click on "5.Tag Instance" at the top of that page.



Add an unique value at the "Value" column, e.g. your name. Then click on "6.Configure Security Group".

Choose "Select an existing security group", select "for_class" under column "Name" and click "Review and Launch":



Ignore the warning, double check that it says "Ubuntu 14.04 LTS (PV)" at AMI Details, and cick "Launch".



Normally you will need to "create a new key pair" the first time through. However, for the purpose of easier management of our shared accout,

**Select "Choose an existing key pair" and Select a key pair "QB2014"**.

Select "Launch Instance."

Click on the link of instance you just started.



Then you should see a "pending" line in the menu.



Wait until it turns green, then make a note of the "Public DNS" (we suggest copying and pasting it into a text notepad somewhere). This is your machine name, which you will need for logging in.

Then, go to *Logging into your new instance "in the cloud" (Windows version)* or *Logging into your new instance "in the cloud" (Mac or Linux version)*

You might also want to read about *Terminating (shutting down) your EC2 instance*.

## Logging into your new instance "in the cloud" (Mac or Linux version)

OK, so you've created a running computer. How do you get to it?

The main thing you'll need is the network name of your new computer. To retrieve this, go to the instance view and click on the instance, and find the "Public DNS". This is the public name of your computer on the Internet.

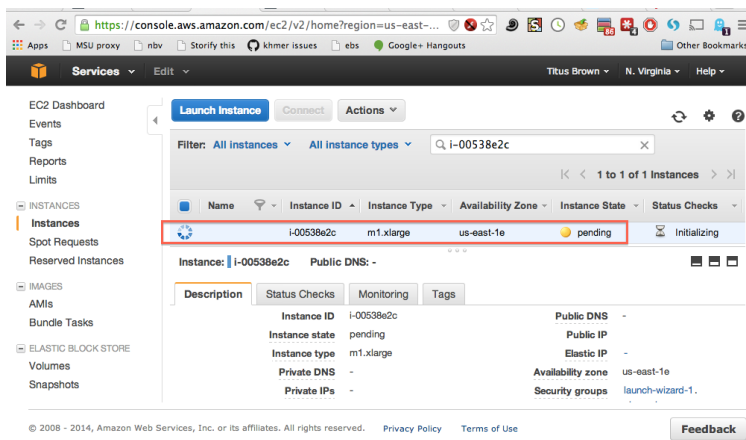Copy this name, and connect to that computer with ssh under the username 'ubuntu', as follows:

First we need a private key file to log in; it is the .pem file we already created during the setup of amazon instance. For convience of managing the shared account, we want you all to use the same .pem file, which can be downloaded from:

**http://lyorn.idyll.org/~gjr/public/QB2014/QB2014.pem**

If your browser is firefox with default settings, it should be downloaded into your "Downloads" folder. Move it onto "Desktop".

If your browser is chrome with default settings, the above link will be shown as a page with text. right click and select "Save As". In the pop up window, change the directory to "Desktop" and the file name is still "QB2014.pem"

Other browsers should be similar to either of the above two case.

Next, start Terminal:

- in Applications... Utilities... for mac
- use shortcut "CTRL+ATL+T" for linux

and then type:

```
chmod og-rwx ~/Desktop/QB2014.pem
```

to set the permissions on the private key file to "closed to all evildoers".

Then type:

```
ssh -i ~/Desktop/QB2014.pem ubuntu@ec2-???-???-???-???.compute-1.amazonaws.com
```

(but you have to replace the stuff after the '@' sign with the name of the host).

Here, you're logging in as user 'ubuntu' to the machine 'ec2-174-129-122-189.compute-1.amazonaws.com' using the authentication key located in 'QB2014.pem' on your Desktop.

You should now see a text line that starts with something like `ubuntu@ip-10-235-34-223:~$`. You're in! Now type:

```
sudo bash
cd /root
```

to switch into superuser mode (see: http://xkcd.com/149/) and go to your home directory.

This is where the rest of the tutorials will start!

You might also want to read about *Terminating (shutting down) your EC2 instance*.

To log out, type:

```
exit
logout
```

or just close the window.

### Logging into your new instance "in the cloud" (Windows version)

### Download .pem file

First we need a private key file to log in; it is the .pem file we already created during the setup of amazon instance. For convience of managing the shared account, we want you all to use the same .pem file, which can be downloaded from:

**http://lyorn.idyll.org/~gjr/public/QB2014/QB2014.pem**

If your browser is firefox with default settings, it should be downloaded into your "Downloads" folder. Move it onto "Desktop".

If your browser is chrome with default settings, the above link will be shown as a page with text. right click and select "Save As". In the pop up window, change the directory to "Desktop" and the file name is still "QB2014.pem"
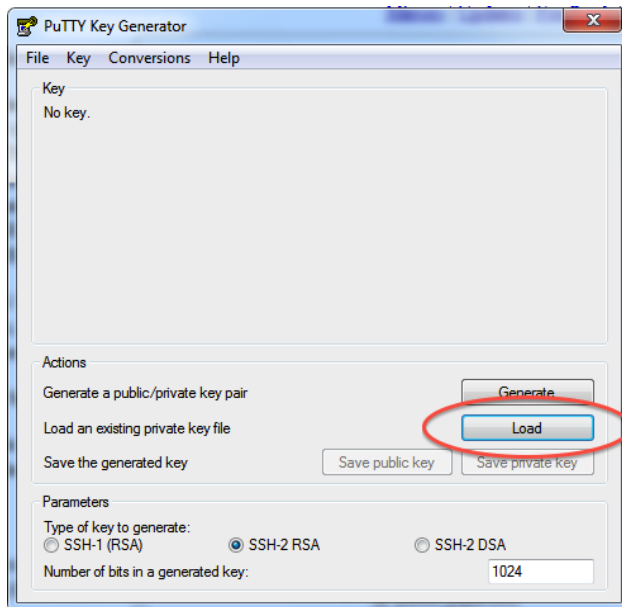
Other browsers should be similar to either of the above two case.

Download Putty and Puttygen from here: http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html
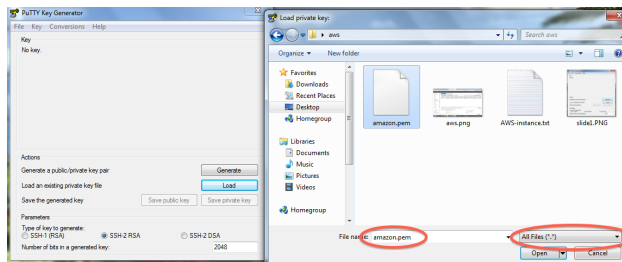
### Generate a ppk file from your pem file
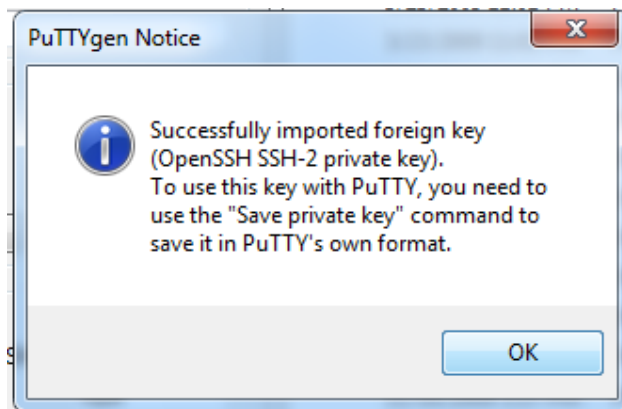
(You only need to do this once!)
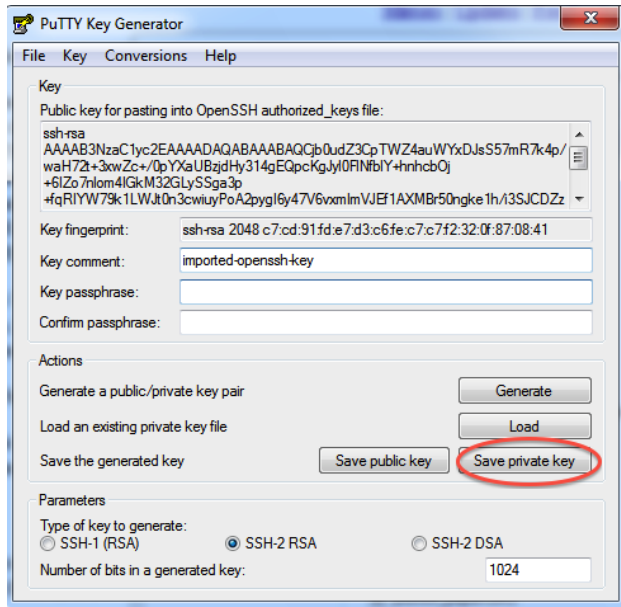
Open puttygen; select "Load".

Find and load your '.pem' file; it's probably in your Downloads folder. Note, you have to select 'All files' on the bottom.
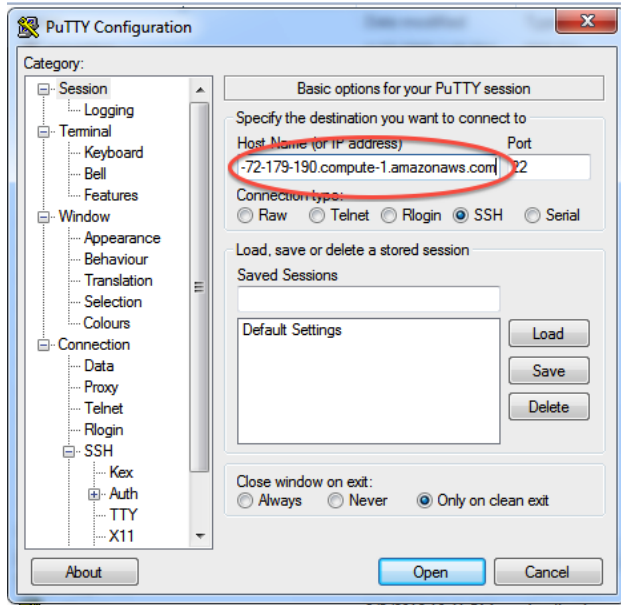


Load it.



Now, "save private key". Put it somewhere easy to find.
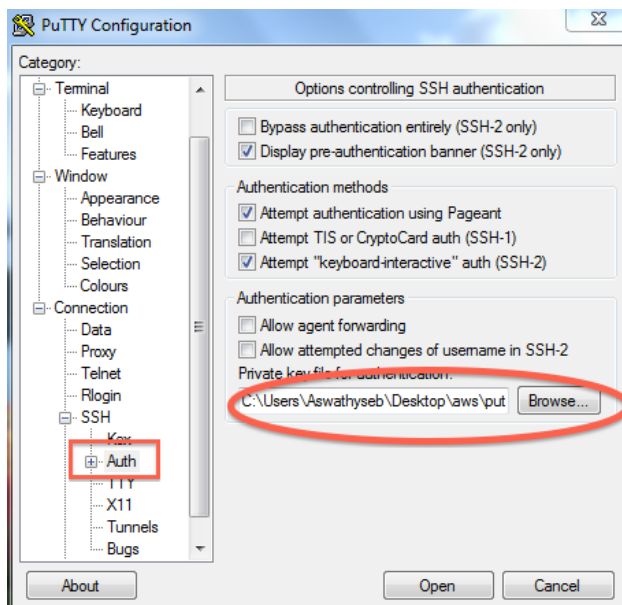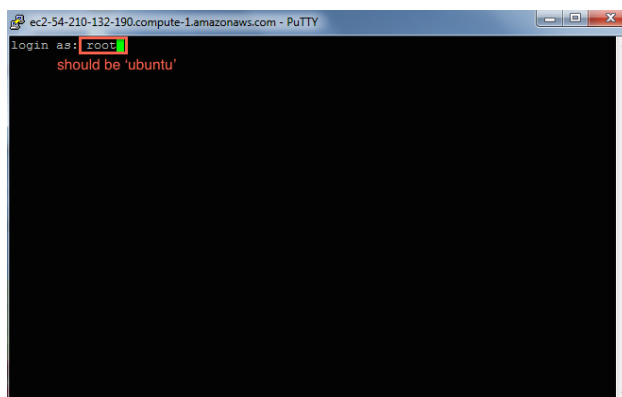
## Logging into your EC2 instance with Putty

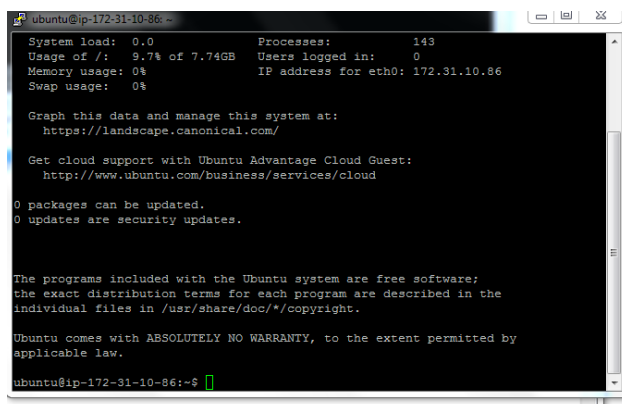Open up putty, and enter your hostname into the Host Name box.



Now, go find the 'SSH' section and enter your ppk file (generated above by puttygen). Then select 'Open'.

Log in as "ubuntu".



Declare victory!



Here, you're logging in as user 'ubuntu' to the machine 'ec2-174-129-122-189.compute-1.amazonaws.com' using the authentication key located in 'amazon.pem' on your Desktop.

You should now see a text line that starts with something like `ubuntu@ip-10-235-34-223:~$`. You're in! Now type:

```
sudo bash
cd /root
```

to switch into superuser mode (see: http://xkcd.com/149/) and go to your home directory.

This is where the rest of the tutorials will start!

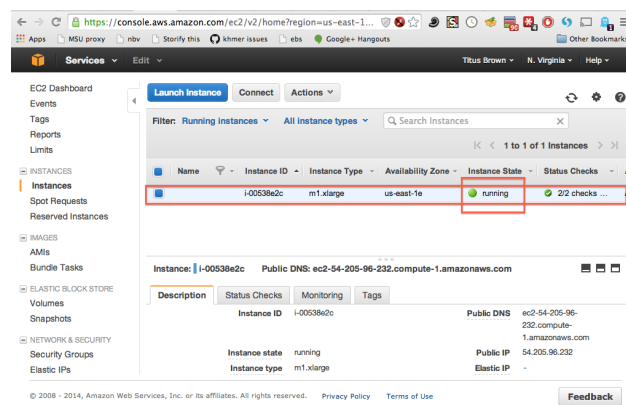You might also want to read about *Terminating (shutting down) your EC2 instance*.

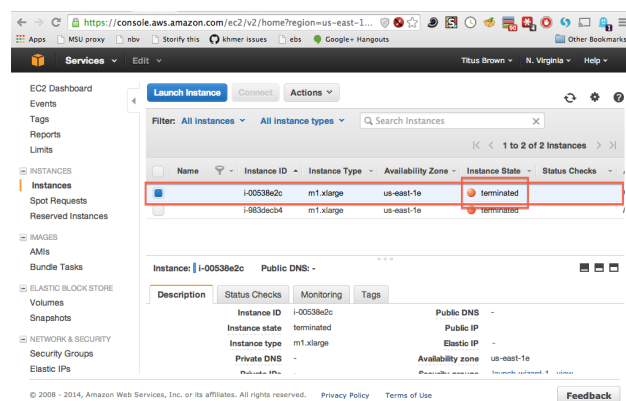To log out, type:

```
exit
logout
```

or just close the window.

## Terminating (shutting down) your EC2 instance

While your instance is running, Amazon will happily charge you on a per-hour basis – check out the pricing for more information. In general, you will want to shut down your instance when you're done with it; to do that, go to your EC2 console and find your running instances (in green).

Then, select one or all of them, and go to the 'Actions...' menu, and then select 'Terminate'. Agree.

After a minute or two, the console should show the instance as "terminated".

### Starting up a custom operating system

The instructions in *Start up an EC2 instance* tell you how to start up a machine with Ubuntu Linux version 14.04 on it, but that machine comes with very little software installed. For anything where you are executing actual analyses, you're going to want to have a bunch of basic software installed.

Therefore, we make custom versions of Ubuntu available as well, that come with some software pre-installed.
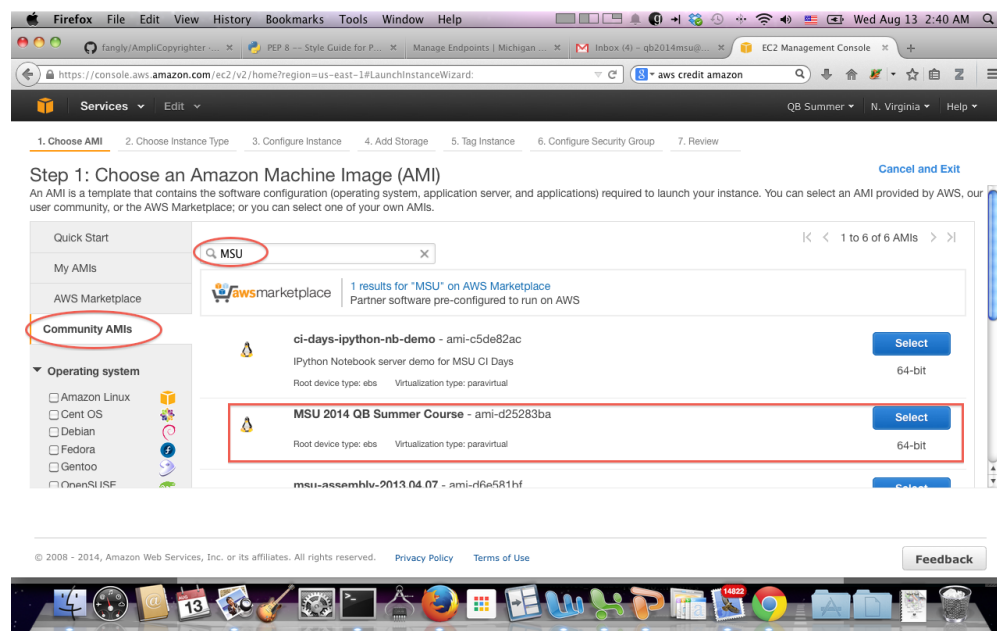
To boot these, go to EC2/Launch and select "Community AMIs" instead of the default "Quick Start";

Then type in the AMI number or name in the search box. We need to use:

**ami-d25283ba**

for this whole workship. Or you can just type "MSU" in the search box and hit ENTER. **"MSU 2014 QB Summer Course"** is the one needed.

Then proceed with the rest of *Start up an EC2 instance*.



### Variant calling

The goal of this tutorial is to show you the basics of variant calling using Samtools.

We'll be using data from one of Rich Lenski's LTEE papers, the one on the evolution of citrate consumption in LTEE.

### Booting an Amazon AMI

Start up an Amazon computer (m1.large or m1.xlarge) using AMI ami-d25283ba (MSU 2014 QB Summer Course) (see *Start up an EC2 instance* and *Starting up a custom operating system*).

Log in the cloud machine from our local computer with Windows or from Mac OS X.

### Logging in

Log in and type:

```
sudo bash
```

to change into superuser mode.

### Software

Softwares are already installed in the ami loaded. Tools used in this tutorial are:

- bwa (mapping reads to reference genome)
- samtools (a fast and versatile tool for processing DNA alignment in SAM/BAM format.)

### Download data

Data are already download in "/mnt" directory:

```
ls /mnt
```

to check the data files:

- assembly (directory of data for tomorrow's workshop)
- REL606.fa (reference genome)
- SRR098038.fastq.gz (reads from illumina sequencing)

### Do the mapping

Now let's map all of the reads to the reference. Start by indexing the reference genome:

```
cd /mnt

bwa index REL606.fa
```

Now, do the mapping of the raw reads to the reference genome:

```
bwa aln REL606.fa SRR098038.fastq.gz  > SRR098038.sai
```

Make a SAM file (this would be done with 'sampe' if these were paired-end reads):

```
bwa samse REL606.fa SRR098038.sai SRR098038.fastq.gz > SRR098038.sam
```

This file contains all of the information about where each read hits on the reference.

Next, index the reference genome with samtools:

```
samtools faidx REL606.fa
```

Convert the SAM into a BAM file:

```
samtools import REL606.fa.fai SRR098038.sam SRR098038.bam
```

Sort the BAM file:

```
samtools sort SRR098038.bam SRR098038.sorted
```

And index the sorted BAM file:

```
samtools index SRR098038.sorted.bam
```

### Visualizing alignments

'samtools tview' is a text interface that you use from the command line; run it like so:

```
samtools tview SRR098038.sorted.bam REL606.fa
```

The '.'s are places where the reads align perfectly in the forward direction, and the ','s are places where the reads align perfectly in the reverse direction. Mismatches are indicated as A, T, C, G, etc.

You can scroll around using left and right arrows; to go to a specific coordinate, use 'g' and then type in the contig name and the position. For example, type 'g' and then 'rel606:553093<ENTER>' to go to position 553093 in the BAM file.

Use 'q' to quit.

### Counting alignments

This command:

```
samtools view -c -f 4 SRR098038.bam
```

will count how many reads DID NOT align to the reference (214518).

This command:

```
samtools view -c -F 4 SRR098038.bam
```

will count how many reads DID align to the reference (6832113).

And this command:

```
gunzip -c SRR098038.fastq.gz | wc
```

will tell you how many lines there are in the FASTQ file (28186524). Reminder: there are four lines for each sequence.

### Calling SNPs

You can use samtools to call SNPs like so:

```
samtools mpileup -uD -f REL606.fa SRR098038.sorted.bam | bcftools view -bvcg - >␣
↪SRR098038.raw.bcf
```

(See the 'mpileup' docs here.)

Now convert the BCF into VCF:

```
bcftools view SRR098038.raw.bcf > SRR098038.vcf
```

You can check out the VCF file by using 'tail' to look at the bottom:

```
tail *.vcf
```

Each variant call line consists of the chromosome name (for E. coli REL606, there's only one chromosome - rel606); the position within the reference; an ID (here always '.'); the reference call; the variant call; and a bunch of additional information.

Again, you can use 'samtools tview' and then type (for example) 'g' 'rel606:4616538' to go visit one of the positions. The format for the address to go to with 'g' is 'chr:position'.

You can read more about the VCF file format here.

### Questions/discussion items

Why so many steps?

## Assembling E. coli sequences with Velvet

The goal of this tutorial is to show you the basics of assembly using the Velvet assembler.

We'll be using data from Efficient de novo assembly of single-cell bacterial genomes from short-read data sets, Chitsaz et al., 2011.

### Booting an Amazon AMI

Start up an Amazon computer (m1.large or m1.xlarge) using AMI ami-d25283ba (see *Start up an EC2 instance* and *Starting up a custom operating system*).

Log in with Windows or from Mac OS X.

### Logging in

Log in and type:

```
sudo bash
```

to change into superuser mode.

### Packages to install

Softwares and packages are already installed:

- velvet (An assmbler)
- khmer (A package for preprocess and assemble large data)
- quast (A tools for evaluating assembly with known genomes)

### Data

Data is already downloaded in /mnt/assembly:

```
ls /mnt/assembly
```

to see data in "/mnt/assembly":

- ecoli_ref-5m-trim.fastq.gz: quality trimmed pair end data sets from E. coli genome sequence data.

- ecoli-reads-5m-dn-paired.fa.gz: data set is a specially processed data set using digital normalization that will assemble quickly.

### Running an assembly

Go inside the data directory:

```
cd /mnt/assembly
```

Now... assemble the small, fast data sets, using the Velvet assembler. Here we will set the required parameter k=21:

```
velveth ecoli.21 21 -shortPaired -fasta.gz ecoli-reads-5m-dn-paired.fa.gz
velvetg ecoli.21 -exp_cov auto
```

Check out the stats for the assembled contigs for a cutoff of 1000:

```
python /usr/local/share/khmer/sandbox/assemstats3.py 1000 ecoli.*/contigs.fa
```

Also try assembling with k=23 and k=25:

```
velveth ecoli.23 23 -shortPaired -fasta.gz ecoli-reads-5m-dn-paired.fa.gz
velvetg ecoli.23 -exp_cov auto

velveth ecoli.25 25 -shortPaired -fasta.gz ecoli-reads-5m-dn-paired.fa.gz
velvetg ecoli.25 -exp_cov auto
```

Now check out the stats for the assembled contigs for a cutoff of 1000:

```
python /usr/local/share/khmer/sandbox/assemstats3.py 1000 ecoli.*/contigs.fa
```

(Also read: What does k control in de Bruijn graph assemblers?.)

### Comparing and evaluating assemblies - QUAST

Run QUAST to mapping contigs (ecoli.23/contigs.fa, and ecoli.25/contigs.fa) reference genome (ecoliMG1655.fa):

```
gunzip ecoliMG1655.fa.gz
/root/quast-2.3/quast.py -R ecoliMG1655.fa ecoli.*/contigs.fa
```

Note that here we're looking at *all* the assemblies we've generated.

Now look at the results:

```
more quast_results/latest/report.txt
```

The first bits to look at are Genome fraction (%) and # misassembled contigs, I think.

### Searching assemblies – BLAST

Build BLAST databases for the assemblies you've done:

```
cd /mnt/assembly

for i in 21 23 25
do
   extract-long-sequences.py -o ecoli-$i.fa -l 500 ecoli.$i/contigs.fa
   formatdb -i ecoli-$i.fa -o T -p F
done
```

BLAST is already installed. Let's search for a specific gene (CRP, a transcription regulator).

The commandline for BLAST is:

```
blastall -i crp.fa -d ecoli-21.fa -p tblastn -b 1 -v 1
```

### Questions and Discussion Points

Why do we use a lower cutoff of 1kb for the assemstats3 links, above? Why not 0?

# Indices and tables

- genindex
- modindex
- search