# COSC2637/2633 Big Data Processing
# Lab 1: Tutorial/Lab class

## Objectives

- Access your individual jumphost that has already been setup for each student
- Create, access and then terminate an AWS EMR cluster
- Responsible use of cloud resources

## Introduction

In order to use AWS EMR, RMIT ITS has created a relatively inexpensive jumphost on AWS for each student enrolled in the course. From this jumphost you will be to create, access and manage an EMR cluster. Creating a cluster requires allocation of additional physical machines in the AWS datacentre (which takes minutes to happen, depending on the size of cluster). However, as the EMR clusters are more expensive machines (and cost by time that the cluster is left set up), it is important that clusters are not left running after you have finished using them (and saved any output). So, at the end of a lab class, it is **very important** you **terminate** your cluster.
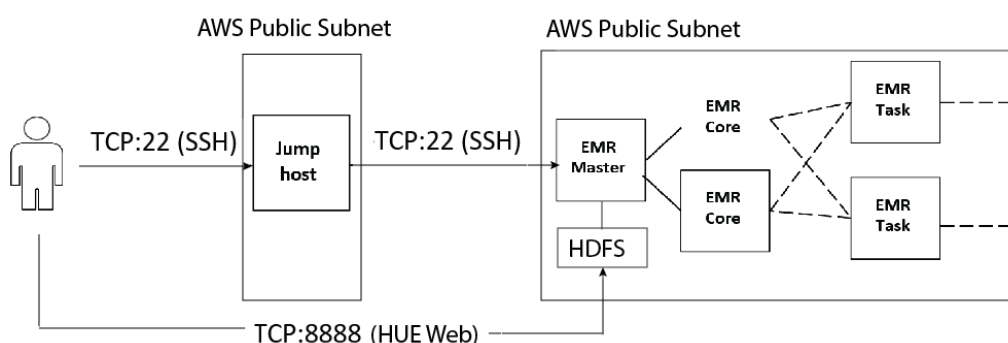
The jumphost and any EMR clusters launched will all be in the AWS US East Region (N.Virginia) (Also known as "us-east-1" or "Standard"). This means any clusters you create will have access to the Public Data Sets hosted in S3 that are provided by AWS in that region.

Example of student jumphost DNS: `sXXXXXXX.jump.cosc2637.route53.aws.rmit.edu.au`

Each student will have a personal SSH key and it will provide access to your jumphost and EMR cluster.

You can run `./create_cluster.sh` from your jumphost to launch your EMR cluster. After the script finishes running, there will be further instructions on how to access that specific cluster (Hue and Hadoop Master node).

Please **do not forget** to run `./terminate_cluster.sh` (from your jumphost) each time you have finished using your cluster (e.g., at the end of the lab class).

# The steps to get you on to the AWS environment
# (For Windows OS users, please see another document)

Any feedback or suggestions about these instructions should be posted on the Course Canvas Discussion Board "Technical questions concerning AMS EMR" via myRMIT.

## Keys (these should be emailed to you at or prior to your practical class)

Save the key location on your device
- For MacOSx/Linux, change permissions:

      `$chmod 600 sxxxxxx-cosc2637.pem`

(replace *sXXXXXXX* with your student number, also shown in the key sent to you from lecturer)

Links for more information:

*http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/AccessingInstances.html*

*http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/AccessingInstancesLinux.html*

*http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/putty.html*

## Configuration required to prevent the inactive shell phenomenon

- **For MacOSx/Linux:** add the following settings into *~/.ssh/config*

      *ServerAliveInterval 15*
      *TCPKeepAlive yes*

## SSH into your jumpbox

`ssh ec2-user@sXXXXXXX.jump.cosc2637.route53.aws.rmit.edu.au -i sxxxxxx-cosc2637.pem`

Accept the host, now you are in the jumpbox if you see something like:

      `[ec2-user@ip-192-168-10-47 ~]$`

From here you can run the start and terminate script as needed. There are several scripts and other files in the directory on the jumphost that should not be changed, otherwise you may not be able to create and manage your cluster.

## Start a cluster

Start the cluster using either of the following commands:

    *$ ./create_cluster.sh*
    or
    *$ sh create_cluster.sh*

The script begins by checking whether a cluster already exists, the following message indicates it does not:

    `An error occurred (ValidationError) when calling the DescribeStacks`
    `operation: Stack with id e23270-emr does not exist`

Then, wait 15 min (it can be longer) for the cluster to start up, you will see when the cluster is ready

## Access the cluster

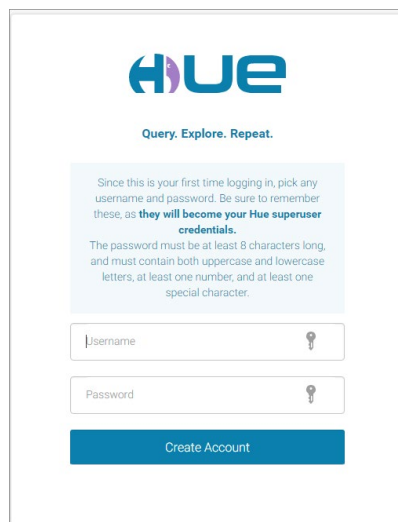You will be shown the DNS for the cluster and the web link now.

```
Checking if EMR Cluster exists ...

An error occurred (ValidationError) when calling the DescribeStacks operation: S
tack with id e20925-emr does not exist
Creating EMR Cluster. This will take about 15 minutes...
{
    "StackId": "arn:aws:cloudformation:us-east-1:089160896324:stack/e20925-emr/7
aac9600-9a7a-11ea-b25a-0e765a9edd1f"
}
Finished creating EMR cluster
Access Hue via a browser here: http://e20925.hue.cosc2637.route53.aws.rmit.edu.a
u:8888

Access EMR Master node via SSH from your jumphost
ie ssh hadoop@e20925.emr.cosc2637.route53.aws.rmit.edu.au -i e20925-cosc2637.pem
[ec2-user@ip-192-168-10-47 ~]$
```

- Use your web browser to use the Hadoop environment via Hue. You will need to setup an account each time as the cluster is new.

  http://sXXXXXXX.hue.cosc2637.route53.aws.rmit.edu.au:8888



For example, my Hue account is "e20925". After logging in HUE, click Files, I will see the HDFS file system like below:



The HDFS file system of the AWS EMR cluster.
- You can directly upload file from your machine to HDFS by click "Upload".
- You can also upload file from AWS EMR cluster master node to HDFS.
More details of the two ways can be found in "AWS EMR - Cheat sheet" in Page 5.

- Connect to the cluster via SSH

`ssh hadoop@e20925.emr.cosc2637.route53.aws.rmit.edu.au -i sxxxxxxx-cosc2637.pem`

> Note before running above ssh, please make sure *sxxxxxxx-cosc2637.pem* in the jumphost. If it does not, you need run the following line on your own machine (replace e20925 by your student id sxxxxxxx).
>
> ```
> [e20925@csitprdap03 BDP]$ scp -i e20925-cosc2637.pem e20925-cosc2637.pem ec2-use
> r@e20925.jump.cosc2637.route53.aws.rmit.edu.au:/home/ec2-user/ █
> ```

You are in the cluster, if you can see something like



Leave the cluster by inputting "exit", then you will return back to Jumphost.

Don't forget to **shutdown the cluster using terminate_cluster.sh**. The shutting down will take about 5 minutes during which you'll still be able to access it via the browser.

# *********Useful Information************

It is a good idea to do just prior to practical class so that it is up and running. However, do NOT leave it running for days at a time as it is expensive and a waste of resources to have idle clusters that are not being used. When finished, don't forget to **shutdown the cluster** to avoid charges. Note that there is no indication when the cluster is terminated.

### AWS EMR - Cheat sheet

**A common issue of cluster login**

You may see the following error when ssh from your jumphost to the hadoop master node,

```
@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
@       WARNING: REMOTE HOST IDENTIFICATION HAS CHANGED!       @
@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
IT IS POSSIBLE THAT SOMEONE IS DOING SOMETHING NASTY!
Someone could be eavesdropping on you right now (man-in-the-middle attack)!
It is also possible that a host key has just been changed.
The fingerprint for the ECDSA key sent by the remote host is
SHA256:e7nHrwQqyHxRsHaHCq+v7+VWNlOzL3D+ZQQTC+yGYjI.
Please contact your system administrator.
Add correct host key in /home/ec2-user/.ssh/known_hosts to get rid of this message.
Offending ECDSA key in /home/ec2-user/.ssh/known_hosts 1
ECDSA host key for s▮▮▮▮▮▮▮.emr.cosc2637.route53.aws.rmit.edu.au has changed and you have
requested strict checking.
Host key verification failed.
```

Every time a new cluster is created it is a new set of hosts so the ssh host key changes, and ssh gives you a warning that it has been changed. To fix the problem, please run the following:

`$ssh-keygen -R sxxxxxxx.emr.cosc2637.route53.aws.rmit.edu.au`