# COSC2637/2633 Big Data Processing
# Lab 1: Tutorial/Lab class (Windows User)

**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.
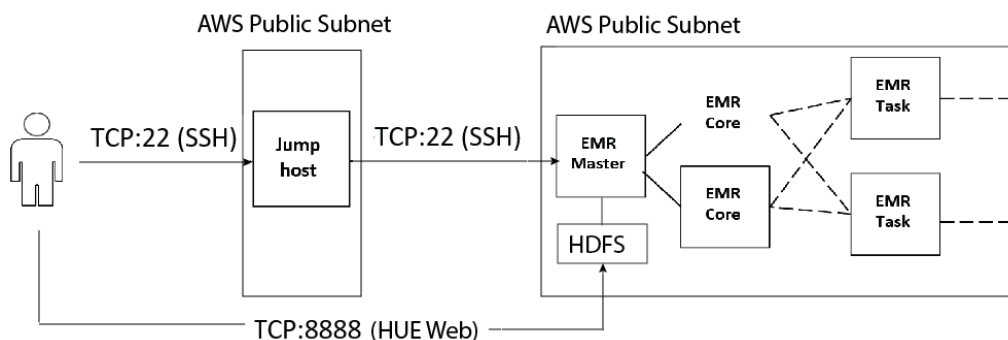
The jumphost and any EMR clusters launched will all be in the AWS US East Region (N.Virginia) (Also known as "us-east-1" or "Standard"). This means any clusters you create will have access to the Public Data Sets hosted in S3 that are provided by AWS in that region.

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 (Windows)

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

1. Save the key location on your device
2. Convert the key using PuTTYGen in Windows OS

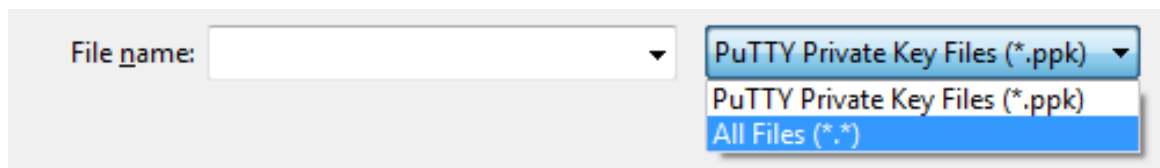**Converting Your Private Key Using PuTTYgen**
PuTTY does not natively support the private key format (.pem) generated by Amazon EC2. PuTTY has a tool named PuTTYgen, which can convert keys to the required PuTTY format (.ppk). You must convert your private key into this format (.ppk) before attempting to connect to your instance using PuTTY.

**To convert your private key**

  a. Start PuTTYgen (for example, from the **Start** menu, choose **All Programs > PuTTY > PuTTYgen**).

  b. Under **Type of key to generate**, choose **RSA**.



  c. Choose **Load**. By default, PuTTYgen displays only files with the extension .ppk. To locate your .pem file, select the option to display files of all types.



  d. Select your .pem file for the key pair that you specified when you launched your instance, and then choose **Open**. Choose **OK** to dismiss the confirmation dialog box.

  e. Choose **Save private key** to save the key in the format that PuTTY can use. PuTTYgen displays a warning about saving the key without a passphrase. Choose **Yes**.

     Note: A passphrase on a private key is an extra layer of protection, so even if your private key is discovered, it can't be used without the passphrase. The downside to using a passphrase is that it makes automation harder because human intervention is needed to log on to an instance or copy files to an instance.

  f. Specify the same name for the key (.pem). PuTTYgen automatically adds the .ppk file extension.

Your private key is now in the correct format for use with PuTTY. You can now connect to your instance using PuTTY's SSH client.
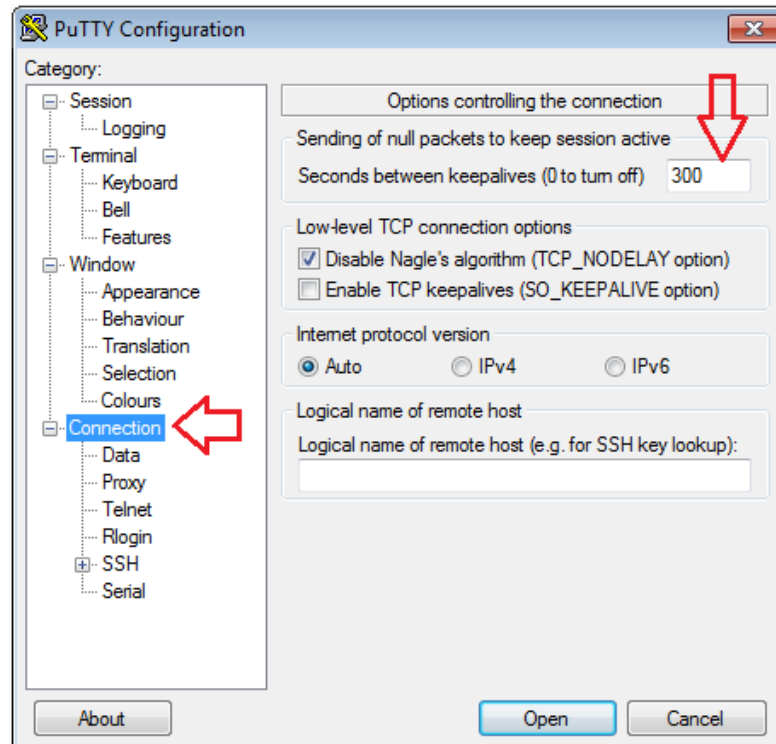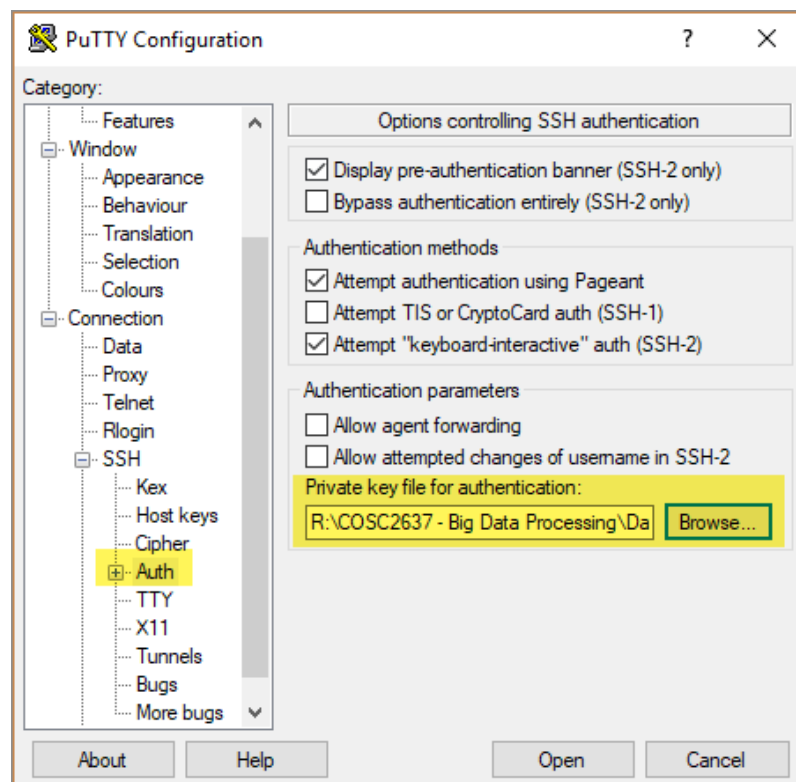
Links for more information:
*http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/AccessingInstances.html*

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

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

### 3. Set up PuTTY connection

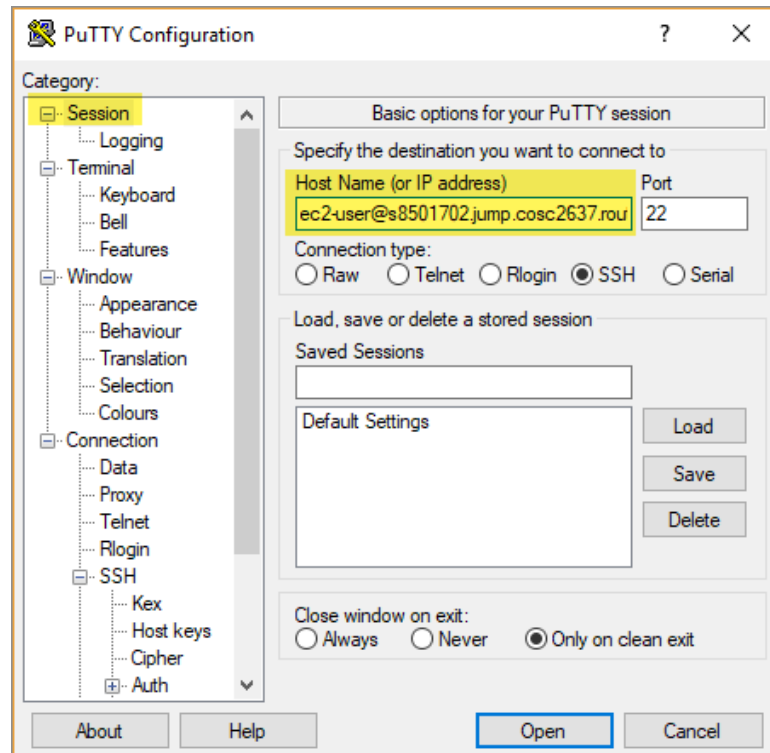a. Configuration required to prevent the inactive shell phenomenon



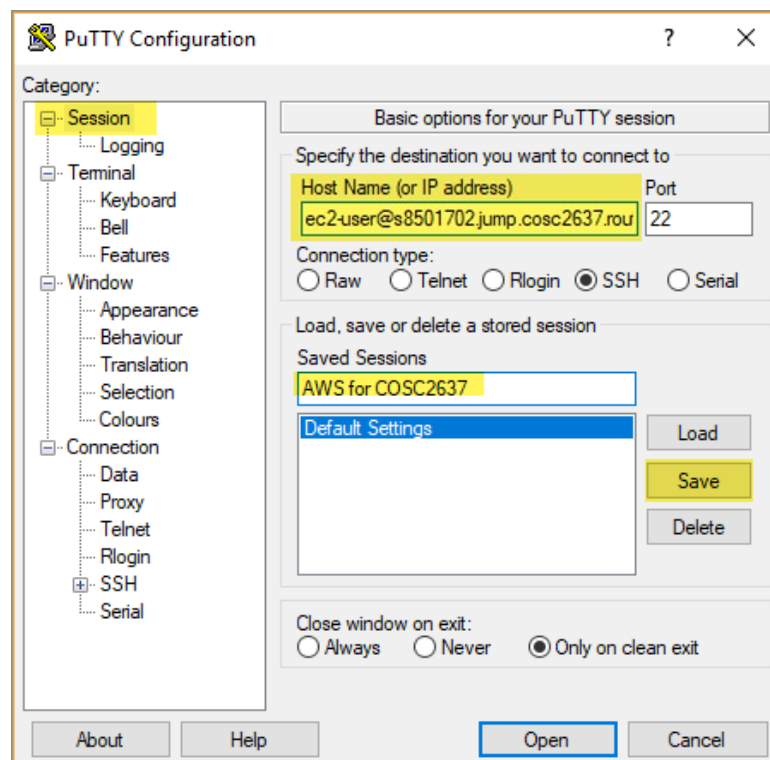b. Add private key you saved earlier (extracted from the key pair file emailed to you)

**c.** Enter jump host name:

`ec2-user@sXXXXXXX.jump.cosc2637.route53.aws.rmit.edu.au`
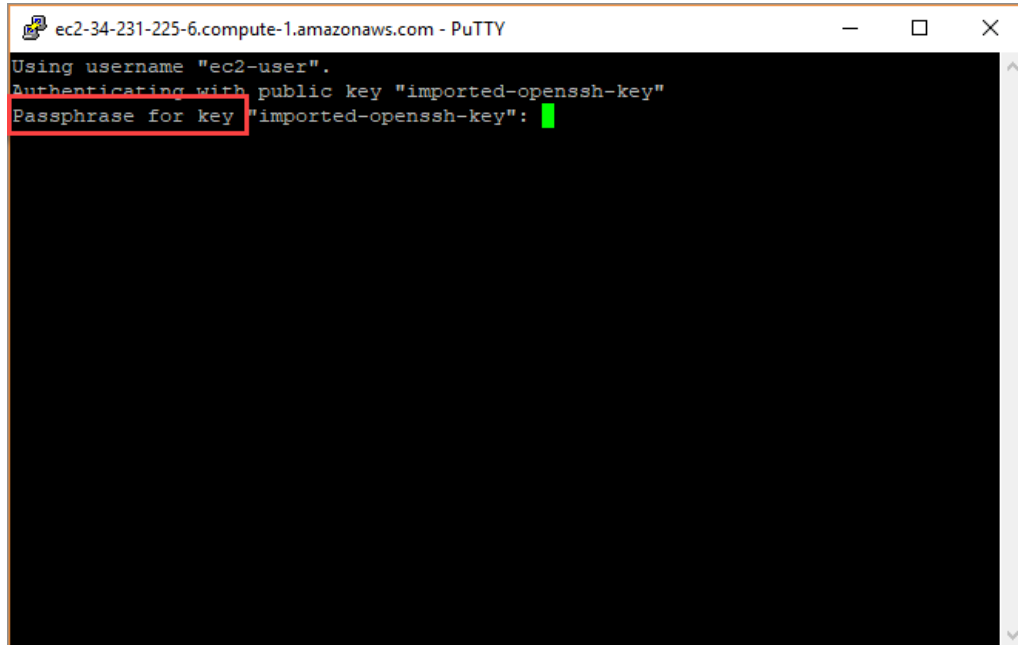Where XXXXXXX is your student number



**d.** Save settings
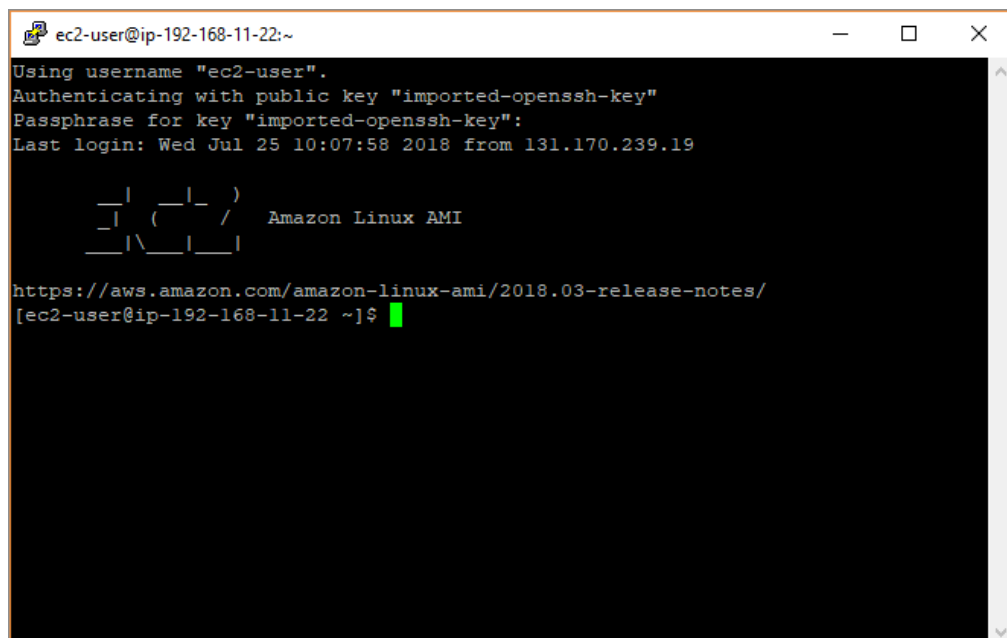
4. Open PuTTY connection

   Open button

5. If you entered a passphrase when saving the private key, you will be prompted for it now.



6. You should now have the AWS $ prompt

7. Enter `./create_cluster.sh`

   Note: `sh create_cluster.sh` also works

```
ec2-user@ip-192-168-11-22:~                                          —  □  ×
Using username "ec2-user".
Authenticating with public key "imported-openssh-key"
Passphrase for key "imported-openssh-key":
Last login: Wed Jul 25 10:07:58 2018 from 131.170.239.19


       __|  __|_  )
       _|  (     /    Amazon Linux AMI
      ___|\___|___|

https://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/
[ec2-user@ip-192-168-11-22 ~]$ ./create_cluster.sh
Checking if EMR Cluster exists ...

An error occurred (ValidationError) when calling the DescribeStacks operation: S
tack with id s8501702-emr does not exist
Creating EMR Cluster. This will take about 15 minutes...
{
    "StackId": "arn:aws:cloudformation:us-east-1:089160896324:stack/s8501702-emr
/49829b70-9080-11e8-bacb-50d501a936b3"
}
```

The EMR cluster will now be created. **It can take 15 minutes or more.** The script begins by checking whether a cluster already exists, so the ValidationError is expected as it is simply saying that the cluster doesn't already exist.

8. Once the cluster is created, it will show the URL you need use to connect to your cluster
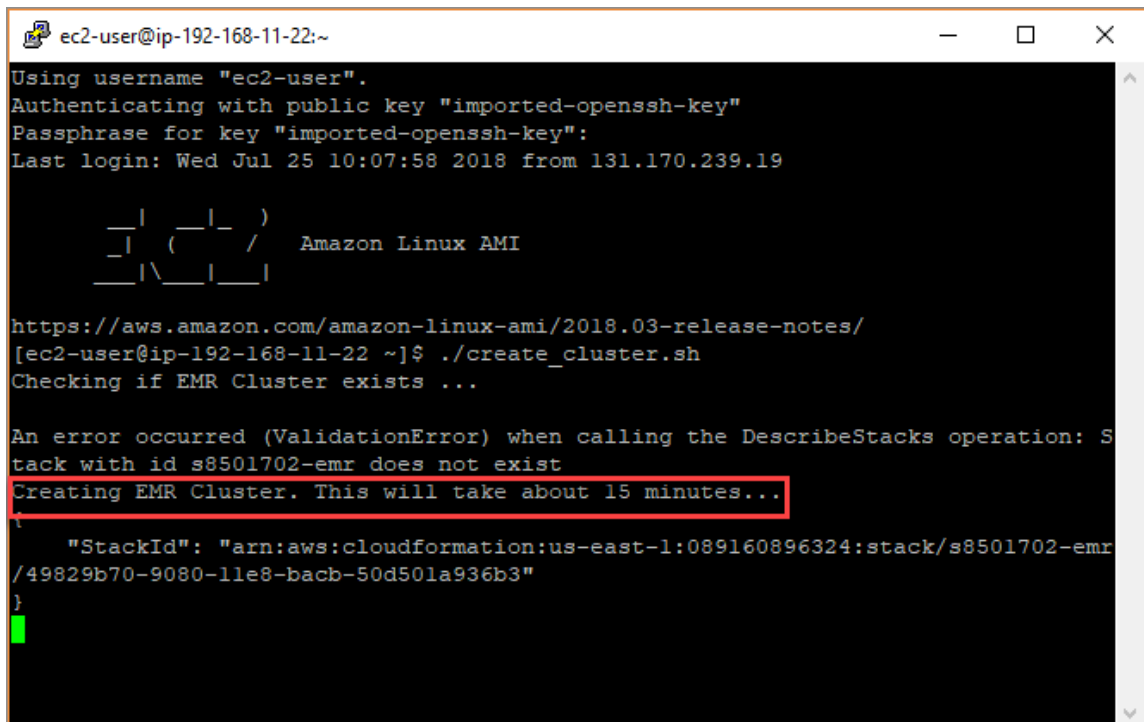
```
ec2-user@ip-192-168-11-22:~                                          —  □  ×
Using username "ec2-user".
Authenticating with public key "imported-openssh-key"
Passphrase for key "imported-openssh-key":
Last login: Wed Jul 25 10:07:58 2018 from 131.170.239.19


       __|  __|_  )
       _|  (     /    Amazon Linux AMI
      ___|\___|___|

https://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/
[ec2-user@ip-192-168-11-22 ~]$ ./create_cluster.sh
Checking if EMR Cluster exists ...

An error occurred (ValidationError) when calling the DescribeStacks operation: S
tack with id s8501702-emr does not exist
Creating EMR Cluster. This will take about 15 minutes...
{
    "StackId": "arn:aws:cloudformation:us-east-1:089160896324:stack/s8501702-emr
/49829b70-9080-11e8-bacb-50d501a936b3"
}
Finished creating EMR cluster
Access Hue via a browser here: http://s8501702.hue.cosc2637.route53.aws.rmit.edu
.au:8888

Access EMR Master node via SSH from your jumphost
ie ssh hadoop@s8501702.emr.cosc2637.route53.aws.rmit.edu.au -i s8501702-cosc2637
.pem
[ec2-user@ip-192-168-11-22 ~]$
```
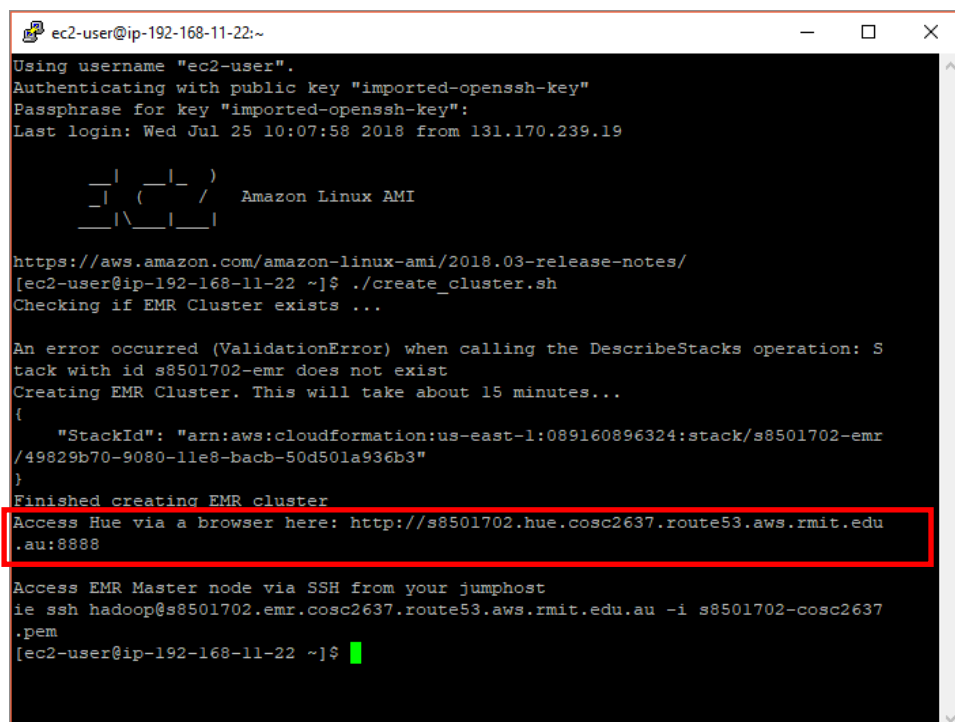
9. Use your browser to use the Hadoop environment via Hue. You will need to setup an account each time as the cluster is new.

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 12.

**Using WinSCP for transferring key to JumpHost home (the key should be the xxx-xxx.pem emailed to you at or prior to your practical class)**

1. Enter jump host name, user and Click Advanced, then select Advanced

2. Under SSH |Authentication, browse to the private key you extracted from the emailed key pair

3. Click OK, and when you're back to the Login screen, save the connection if desired.

4. Login. If you entered a passphrase when saving the private key, you will be prompted for it now



5. In the left panel, browse to the location of the key pair file you were emailed (*.pwm), then drag this to the right panel (the jump host), to copy the key pair file.

Then, go to JumpHost home (as shown in the last picture in page 5 of this document) and you should find xxx-xxx.pem. Next

**$ chmod 400 xxx-xxx.pem**

And login the newly created cluster by ssh hadoop@sxxxx.emr.... as shown below

```
Access EMR Master node via SSH from your jumphost
ie ssh hadoop@s8501702.emr.cosc2637.route53.aws.rmit.edu.au -i s8501702-cosc2637
.pem
```

If logging in successfully, you will see

```
https://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/
40 package(s) needed for security, out of 63 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEE MMMMMMM          MMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::E M::::::::M        M:::::::M R::::::::::::::R
EE::::EEEEEEEEE::::E M:::::::::M      M:::::::M R:::::RRRRR:::::R
  E:::::E       EEEEE M::::::::::M    M:::::::::M RR::::R     R::::R
  E:::::E            M:::::M:::M    M:::M:::::::M   R:::R      R::::R
  E::::EEEEEEEEEE    M:::::M M:::M M:::M M:::::M   R:::RRRRRR:::::R
  E:::::::::::::::E   M:::::M  M:::M:::M  M:::::M   R:::::::::::RR
  E::::EEEEEEEEEE    M:::::M   M:::::M   M:::::M   R:::RRRRRR::::R
  E:::::E            M:::::M    M:::M    M:::::M   R:::R     R::::R
  E:::::E       EEEEE M:::::M     MMM     M:::::M   R:::R      R::::R
EE::::EEEEEEEE::::E M:::::M             M:::::M R:::R      R::::R
E::::::::::::::::::E M:::::M             M:::::M RR::::R     R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMM            MMMMMMM RRRRRRR      RRRRRR


[hadoop@ip-192-168-23-29 ~]$
```

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.

```
ec2-user@ip-192-168-11-22:~
Using username "ec2-user".
Authenticating with public key "imported-openssh-key"
Passphrase for key "imported-openssh-key":
Last login: Wed Jul 25 10:07:58 2018 from 131.170.239.19

     __|  __|_  )
     _|  (     /   Amazon Linux AMI
    ___|\___|___|

https://aws.amazon.com/amazon-linux-ami/2018.03-release-notes/
[ec2-user@ip-192-168-11-22 ~]$ ./create_cluster.sh
Checking if EMR Cluster exists ...

An error occurred (ValidationError) when calling the DescribeStacks operation: S
tack with id s8501702-emr does not exist
Creating EMR Cluster. This will take about 15 minutes...
{
    "StackId": "arn:aws:cloudformation:us-east-1:089160896324:stack/s8501702-emr
/49829b70-9080-11e8-bacb-50d501a936b3"
}
Finished creating EMR cluster
Access Hue via a browser here: http://s8501702.hue.cosc2637.route53.aws.rmit.edu
.au:8888

Access EMR Master node via SSH from your jumphost
ie ssh hadoop@s8501702.emr.cosc2637.route53.aws.rmit.edu.au -i s8501702-cosc2637
.pem
[ec2-user@ip-192-168-11-22 ~]$ sh terminate_cluster.sh
Terminating EMR cluster
[ec2-user@ip-192-168-11-22 ~]$
```
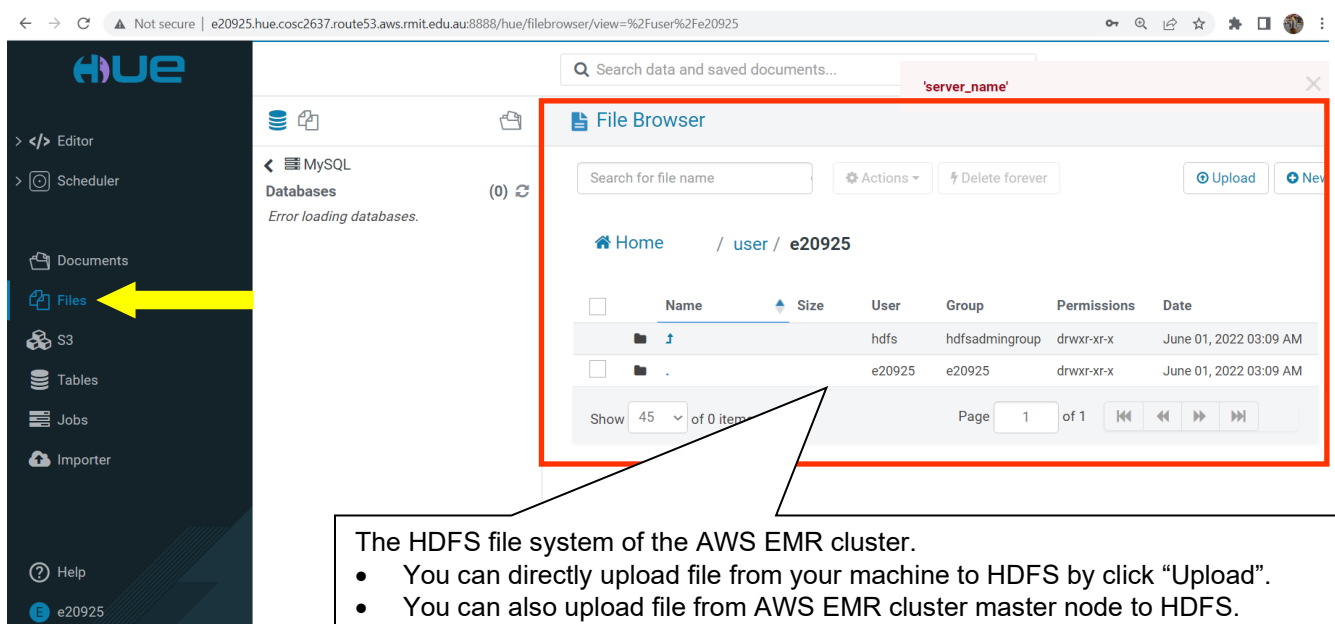
# *********Useful Information************

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.

## A common issue of cluster login

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

```
@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
@     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`

## AWS EMR - Cheat sheet

upload data from you desktop/laptop to HDFS via HUE

**Your desktop/laptop**

1. WinSCP - copy a file to jumphost
2. Putty - login to jumphost

**jumphost**

```
[ec2-user@ip-192-168-10-47 ~]$ ./create_cluster.sh
```

1. ./create_cluster.sh to create the cluster
2. ./terminate_cluster.sh to terminate the cluster
3. ssh-keygen -R sxxx.emr.cosc2637.route53.aws.rmit.edu.au
   You may need run if the cluster is newly created
4. ssh hadoop@sxxx.emr.cosc2637.route53.aws.rmit.edu.au -i e20925-cosc2637.pem
   login the cluster
5. scp -i xxx.pem file hadoop@sxxx.emr.cosc2637.route53.aws.rmit.edu.au:/home/hadoop/
   copy file to cluster from jumphost

**Cluster Master node (.jar)**

```
[hadoop@ip-192-168-16-22 ~]$
```

1. hadoop fs -copyFromLocal ~/file /user/ke/
2. hadoop fs -copyToLocal /user/ke/ ~/
The ~/file above is the file name in the Cluster Master node (known as local relative to HDFS);
/user/ke is a folder in HDFS

**HUE** → **HDFS (data)**

---

AWS Public Subnet          AWS Public Subnet

TCP:22 (SSH) → **Jump host** → TCP:22 (SSH) → **EMR Master**

**EMR Core** — **EMR Task**
**EMR Core** — **EMR Task**

**HDFS**

TCP:8888 (HUE Web)