# COMP 330/543 Lab 1: Using Amazon EMR to Run a Hadoop Job

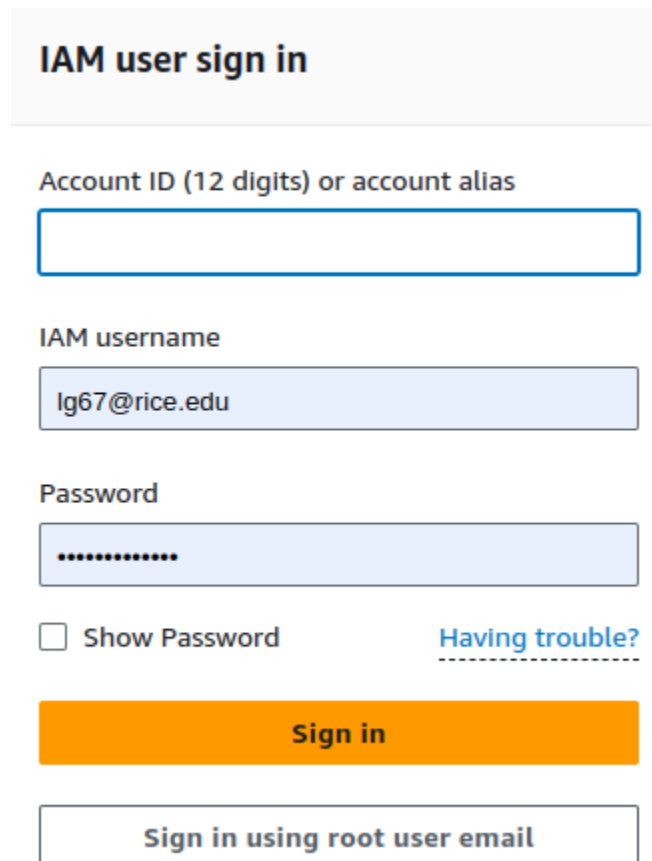**Note:** this assumes you have previously signed up for an Amazon account.

In this lab, you will:

1. Create a Hadoop cluster using Amazon AWS.
2. Compile a Hadoop MapReduce program using the Java compiler (Hadoop is a popular open source MapReduce tool).
3. Connect to your AWS cluster using SSH (Secure Shell).
4. Load data into Hadoop's Distributed File System (HDFS).
5. Run your Hadoop program to process the data.

## Task 1: Start Up a Hadoop Cluster

1. Go to Amazon's AWS website (aws.amazon.com) and click on **"Sign in to console"**. Sign in with your username and password.
    a. If you get a Login page asking for **"IAM user sign in",** click on the **"Sign in using root user mail"** at the bottom.

### IAM user sign in

Account ID (12 digits) or account alias

IAM username

lg67@rice.edu

Password

••••••••••••

☐ Show Password          Having trouble?

**Sign in**

Sign in using root user email

b.  You should then get a different Login form that looks as follows

**Sign in**

● **Root user**
Account owner that performs tasks requiring
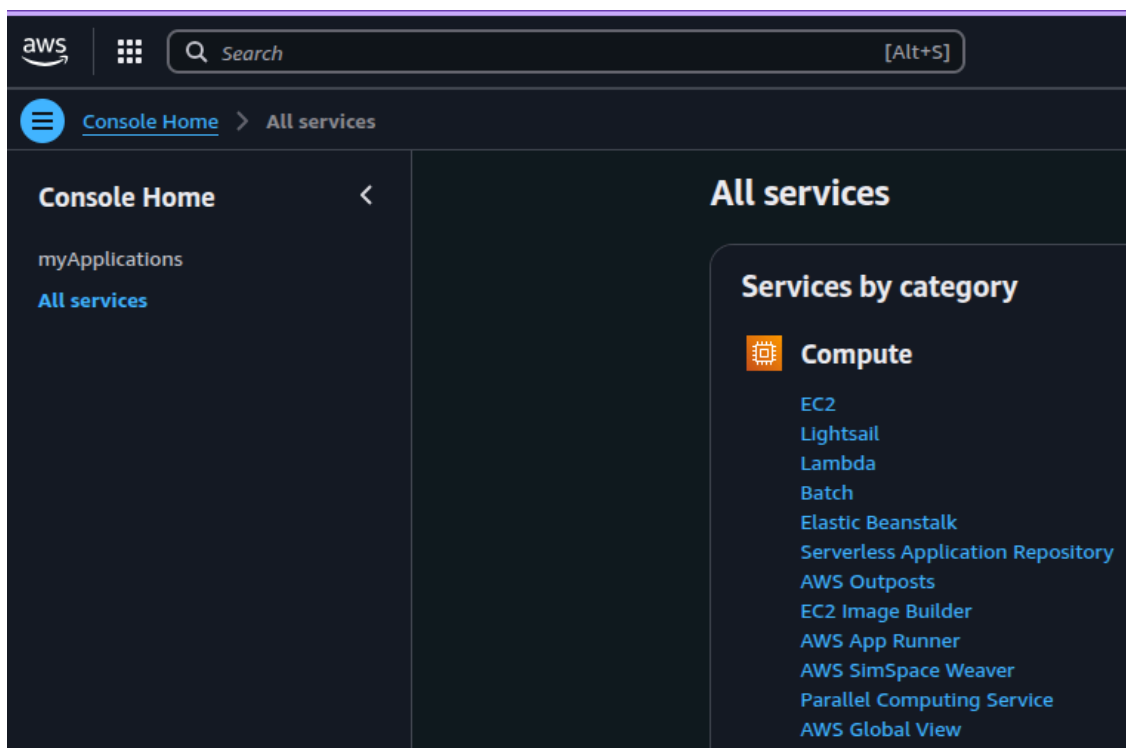unrestricted access. Learn more

○ **IAM user**
User within an account that performs daily tasks.
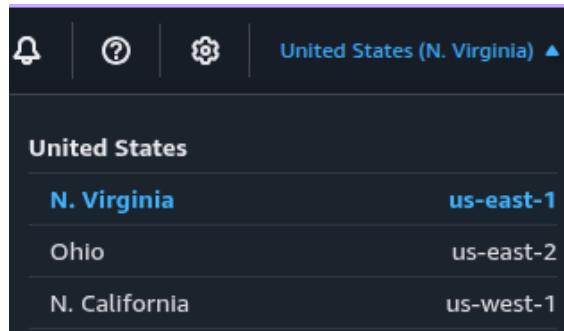Learn more

**Root user email address**

lg67@rice.edu

[ Next ]

2.  Click on the hamburger menu icon in the top left corner, followed by **"All Services"**, and finally on **"EC2"** which can be found under the **"Compute"** section.
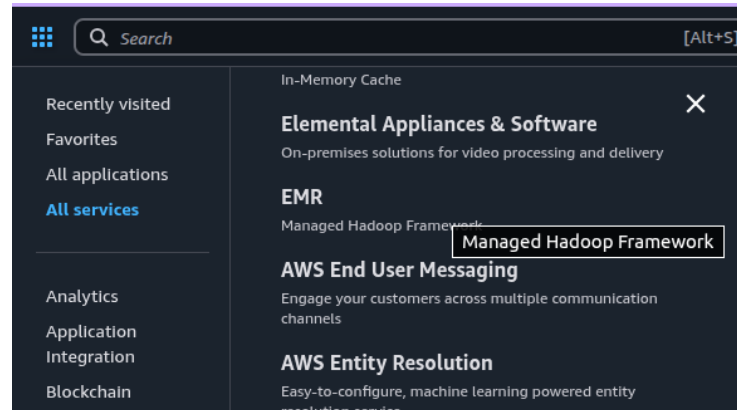
3.  In the top right corner, you will find a dropdown menu to select a region. Make sure you select **"United States (N. Virginia) us-east-1"**.



4.  Next, you will create a **"Key Pair"** that will allow you to connect securely to the cluster. In the left-hand side menu click on **"Key Pairs"** under **"Network & Security"**.

5.  Click the orange "Create key pair" button in the top right corner.
    a.  Pick a name that is likely unique to you.
    b.  Select RSA as the Key pair type.
    c.  Select **".pem"** for the format if using Mac/Linux or **".ppk"** if using Windows.

6. When you click on "Create key pair", a file named with your key will automatically be downloaded into your computer. The file name will be the same name you chose for your key and its extension will be either .pem or .ppk.
   a. **Make sure to save it somewhere you can find it!**

7. Now it's time to create a cluster where you can run Hadoop programs. Again, go to "Services" and "All Services" but this time look for "EMR" (Elastic Map Reduce)



8. Click on the orange "Create cluster" button in the top right corner and select the following options:

9. For the **"Primary"** an **"Core"** instance groups you can choose *m7i.2xlarge*. These machines cost ~40 cents per hour each.



10. If any additional tasks are shown ("Task - 1" in the picture below), you should remove them by clicking the **"Remove instance group button"**.

11. Set your cluster size to have two worker nodes



12. Under **"Networking"**, click on **"Create VPC"** which will open a new tab in your browser.
    a. Here, click on the orange **"Create VPC"** button at the top, select **"VPC and more"** and name your VPC.

b. Then choose **3** under **"Number of Availability Zones"** and **set the number of private subnets to 0.**



c. Finally, click on **"Create VPC"** at the bottom.

13. Go back to the previous tab and select your newly created VPC for your cluster. If your new VPC does not show, simply click on the reload button in the top right corner of the pop-up window.



14. Under **"Cluster logs"**, disable the checkbox "P**ublish cluster-specific logs to Amazon S3"**

15. Under "Security configuration and EC2 key pair", select the key pair created previously.



16. Under **"Identitiy and Access Management (IAM) roles"**
    a.  Select **"EMR_DefaultRole"** for your **"Amazon EMR service role"**. (Note: if that option does not show see below!)



    b.  If "EMR_DefaultRole" is not there, you will need to create it.
        i.  Open a new AWS tab/window and again go to **"All Services"** but this time look for **"IAM"**.

ii. In the left-hand side, under **"Access Management"**, click on **"Roles"** and then click the orange **"Create role"** button in the top right corner.

iii. In the next page, select **"AWS service"** as the entity type and **"EMR"** as the service or use case.



iv. Click "next" a couple times and then name your role as **"EMR_DefaultRole"**.

v. Finish by clicking **"Create role"**. The new **"EMR_DefaultRole"** should now appear in the Roles listing. **Click on it** and it should take you to the following page:

vi. Click on **"Add permissions -> Attach policies"**. Then, in the next page's search bar, search for *"AmazonElasticMapReduce Role"*, click the check box next to it and then click on **"Add permissions"**:



vii. Back in the EMR_DefaultRole page, remove the *"AmazonEMRServicePolicy_v2"* permission policy by clicking the check box next to it and then the **"Remove"** button.

viii.    You can now select the **"EMR_DefaultRole"** for your cluster. (You may need to hit the reload arrow next to the selector).

c.    Under **"EC2 instance profile for Amazon EMR"**, choose **"Create an instance profile"** and give it access to all S3 buckets in your account.



17. Finally, click on the **"Create cluster"** button on the right-hand side. It will take a few minutes for your cluster to be created and the instances to start.
    a.    You will know that your cluster has finished booting up when you see the option to **"Connect to the Primary node using SSH"** under the Cluster management column.



18. In order to connect to your cluster, you must make it so that you can connect via SSH. In your cluster's main page, go to **"Network and Security"** and click on **"EC2 security groups (firewall)"**.

a. Next click on the link below **"Primary node, EMR managed security group"**



b. In the next page, click on the **"Edit inbound rules"** button located on the right-hand side of the screen.
c. Next, click on **"Add rule"**, and select **"SSH"** for the type and **"Anywhere-IPv4"** as the Source.



d. Click on **"Save Rules"** to apply the changes.

19. If you ever want to get to the page that lists all your EMR clusters, just click the **"AWS"** logo in the top left corner and then search for **"EMR"**.



20. **Shutdown your cluster** by selecting it and clicking on the "Terminate" button in the top right. You are ready to move onto Task 2.

21. Some of you **may encounter a not-so-uncommon issue** when attempting to create your clusters with the specifications detailed here. After a few minutes, your cluster will terminate and you will see a message like the following:



The issue stems from the fact that sometimes Amazon assigns new accounts with a smaller quota for standard type (A, C, D, H, I, M, R, T Z) instances.

You can check your assigned quota by using the search bar on top to search for **"Service Quotas"** which should take you to your service quotas dashboard.

Use the search box in the **"Manage quotas"** box to look for **"Amazon EC2"** and then click on **"View quotas"**.



In the next page, search for "demand" in the search bar and then click on **"Running On-Demand Standard (A, C, D, H, I, M, R, T, Z) instances:**



If your assigned quota is less than 24, you will not be able to create the cluster with the specified requirements. You should use the **"Request increase at account level"** button in the top right to request a quota of at least 24 units.

## Task 2: Compile a Hadoop Program

**Note:** This assumes you have installed the IntelliJ IDE in your computer. If you used DataGrip for Assignment 1, IntelliJ is also included with your student license.

1. Start by booting up IntelliJ and creating a new Java project.
   a. **Make sure you are using Java 8** (version 1.8) as other versions can have compatibility issues.



   b. Note that **if you don't have Java version 1.8** you can download it directly from IntelliJ. Open a new Java project. Click on the JDK drop down menu. Then click on "Download JDK". A list of available versions will show (it may take a little bit of time to load), select version 1.8 (near the bottom). Any vendor should be fine (e.g., Amazon Coretto 1.8.0_422). Click download and finish creating the project.
   c. Here are another couple of options to download the Java 1.8 JDK: Option 1, Option 2.

2. Download and unzip all the JAR files included here (make sure there are 7 of them). Put them into a directory in your computer.

3. Add the JAR files to your project by going to **"Files" -> "Project Structure" -> "Modules" -> "Dependencies"** and clicking the **"+"** symbol
   a. Make sure that all 7 files are added, sometimes the *"hadoop-client-2.7.1.jar"* **does not get added** for some reason.
      i. If this happens then repeat the process but add the *"hadoop-client-2.7.1.jar"* file individually. Select **"Classes"** in the small window that pops up.
   b. Make sure you click on the **"Apply"** button.

4. Next, create a new file called **WordCount.java** in the **src** folder of your project.
   a. **Capitalization is important** here so make sure that the name you use is an exact match.

5. In the ZIP file that included all the JAR files, there was a ***WordCount.java***. Copy and paste the contents of this file into the WordCount.java file in your project.

6. Next, we will create a JAR artifact of our WordCount program that can be run in the EMR cluster.
   a. Go to **"File" -> "Project Structure" -> "Artifacts"**, click on the **"+"** and select **"JAR" -> "From modules with dependencies"**.



   b. In the next pop up, select **"WordCount"** as the main class and click **"OK"**. Click **"Apply"** and then **"OK"**.

c. Next, go to **"Build" -> "Build Artifacts"** and select **"Build"** in the pop up window.



d. A new ***[Project Name].jar*** file will be created in the the following folder:
   i. ***[Project Name]/out/artifacts/[Project Name]***

7. You are now ready to move on to Task 3.

## Task 3: Connect to the EMR Cluster through SSH

Time to go back to AWS. Once your cluster is up and running, you will want to connect to the primary node so that you can run Hadoop jobs on it.

1. If you followed the instructions in Task 1, you will have terminated your cluster and won't be able to use it anymore. Fortunately, you do not need to repeat the whole process as you can easily create a new clone of your previous cluster.

   a. Go to your cluster page and clone it by selecting it and clicking on the **"Clone"** button located on the top right corner.



   b. Next, just click on the orange **"Clone cluster"** button on the right side of the screen and your clone will begin booting up.

    c.   Remember that you won't be able to connect into your cluster until you see the option to **"Connect to the Primary node using SSH"**.

    d.   Once the option appears, click on it to obtain the ssh command used to connect to your cluster.

2. The next step will be platform dependent,

    a.   If you are using **MacOS/Linux**, simply open a terminal and paste the ssh command.

        i.   Keep in mind that the way the command is structured by default is to assume that your key pair file (*.pem* file) is **located in your Home directory**. If your *.pem* file is not there, the command will not work. Simply, **change the relative path** used by the command to point to the location of your *.pem* file.

        ii.   **IMPORTANT:** If you get an error message saying that you are using an **"Unprotected Private Key File"**, all you need to do is change the permissions on your .pem file by running the following command:

*chmod 600 [Your Key Name].pem*

        iii.   You should be able to successfully connect to your cluster

```
lguzmann@lguzmann-rice:~$ ssh -i ~/luis_aws_key.pem hadoop@ec2-44-204-127-246.compute-1.ama
zonaws.com
Last login: Fri Sep 27 20:32:57 2024 from 168.5.20.64
    ,       #_
   ~\_  ####_        Amazon Linux 2
  ~~  \_#####\
  ~~     \###|       AL2 End of Life is 2025-06-30.
  ~~       \#/ ___
   ~~       V~' '->
    ~~~         /    A newer version of Amazon Linux is available!
     ~~._.   _/
        _/ _/        Amazon Linux 2023, GA and supported until 2028-03-15.
      _/m/'             https://aws.amazon.com/linux/amazon-linux-2023/

12 package(s) needed for security, out of 18 available
Run "sudo yum update" to apply all updates.

EEEEEEEEEEEEEEEEEEEE MMMMMMMM           MMMMMMMM RRRRRRRRRRRRRRRR
E::::::::::::::::::::E M:::::::M         M:::::::M R::::::::::::::::R
EE:::::EEEEEEEEE:::E M:::::::::M       M:::::::::M R:::::RRRRRR:::::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 MMMMMMMM             MMMMMMM RRRRRRR      RRRRRR

[hadoop@ip-10-0-26-88 ~]$
```

b.  If you are using Windows, we will assume you are using [PuTTy](PuTTy).
      i.    In the left-hand side menu, click on "Connection" -> "SSH" -> "Auth", then click "Browse" to select the private key file (.ppk file).



     ii.    Select your key and click "Open", then go back to the "Session" and enter your machine information ("hadoop@...", you get this from your cluster page in AWS) and click on "Open".

    iii.    You should be able to successfully connect to your cluster.

## Note on PuTTy time out:

If PuTTy is timing out. Do these steps:

1. Go to EC2 console
2. Click Instances on Left
3. Select your instance
4. In the Description tab, locate Security Groups and click the available group link for "Primary Node"
5. Click edit button on Inbound tab
6. Click Add Rule and select SSH for type, Port Range 22, and Source Anywhere

## Task 4: Load Data into Hadoop's Distributed File System (HDFS)

Next we need to transfer the JAR file we created in Task 2 to the EMR cluster.

1. Remember that your JAR file will be located in a folder

   ***[Project Name]/out/artifacts/[Project Name]_jar/[Project Name].jar***

   So, for instance, if your project name is Lab1, the path will be

   ***Lab1/out/artifacts/Lab1_jar/Lab1.jar***

   **Note:** These examples are, of course, relative paths that assume your current working directory is the folder where IntelliJ saves projects.

2. How you transfer the file **into** your cluster will be platform dependent.

   a. **If you are using MacOS/Linux** you can use either the ***scp*** or ***sftp*** commands.

      i. Using scp (Secure CoPy), it is a single-line command as follows:

         *scp -i [path to key] [path to jar] [cluster ip address]:*

To me this looked like:

*scp -i ~/luis_aws_key.pem ~/IdeaProjects/Lab_1/out/artifacts/Lab_1_jar/Lab_1.jar* [*hadoop@ec2-44-204-127-246.compute-1.amazonaws.com*](mailto:hadoop@ec2-44-204-127-246.compute-1.amazonaws.com)*:*

**IMPORTANT:** You should execute the command above **in your local computer's**

**terminal, not in the cluster's terminal**. Remember that you are trying to copy something **FROM your computer INTO your cluster**.

> ii. Using sftp (Secure File-Transfer Protocol), it is a multi-step process:
> > 1. *sftp -i [path to key] [cluster ip address]*
> > 2. *put "[path to jar]"*
> > 3. *exit*
>
> iii. Either way, the JAR file should be copied into your cluster's home directory. You can corroborate this with the **ls** command in your cluster's terminal.
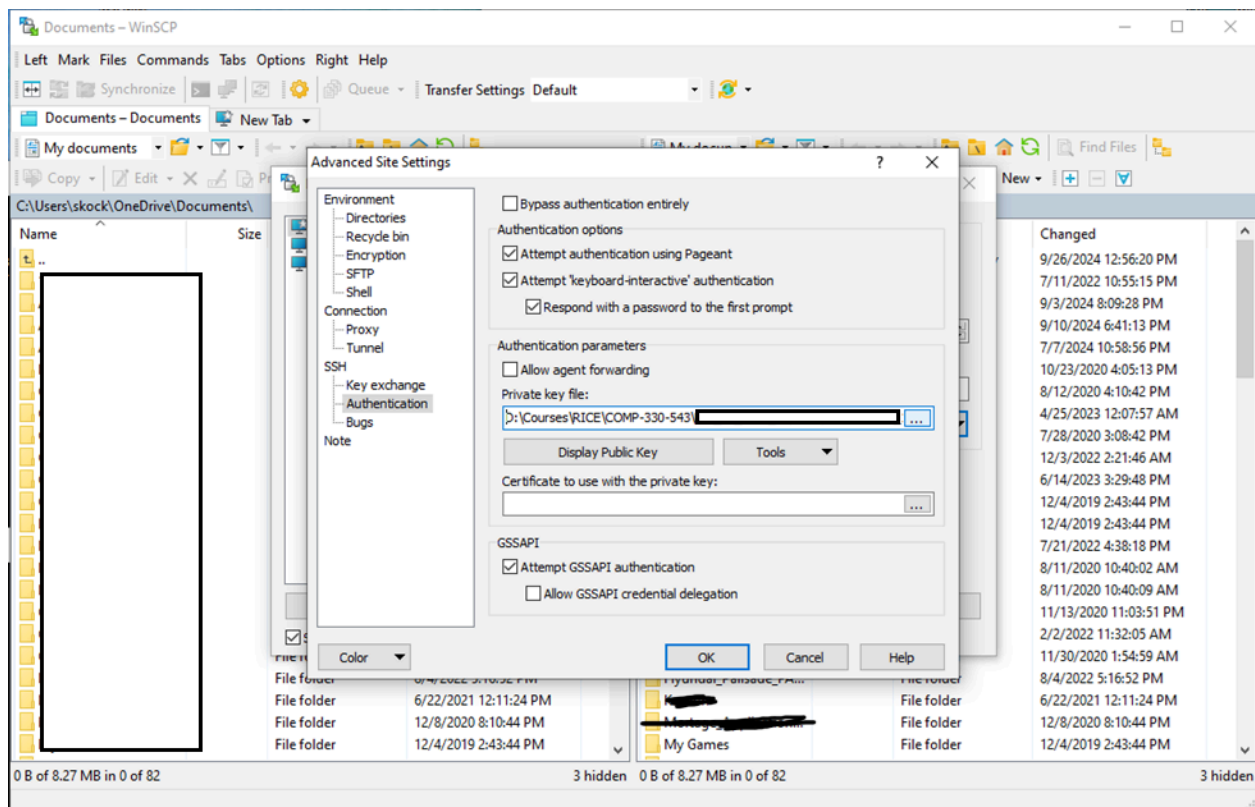
```
EEEEEEEEEEEEEEEEEEEEE MMMMMMMM          MMMMMMMM RRRRRRRRRRRRRRR
E::::::::::::::::::::E M:::::::M        M:::::::M R::::::::::::::R
EE:::::EEEEEEEEE:::E M::::::::M        M::::::::M R:::::RRRRRR::::R
  E:::E        EEEEE M:::::::::M      M:::::::::M RR::::R     R::::R
  E:::E              M::::::M::::M    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:::::EEEEEEEEE:::E M:::::M              M:::::M   R:::R     R::::R
E::::::::::::::::::::E M:::::M              M:::::M RR::::R     R::::R
EEEEEEEEEEEEEEEEEEEEE MMMMMMM              MMMMMMM RRRRRRR     RRRRRR

[hadoop@ip-10-0-26-88 ~]$ ls
Lab_1.jar
[hadoop@ip-10-0-26-88 ~]$ 
```
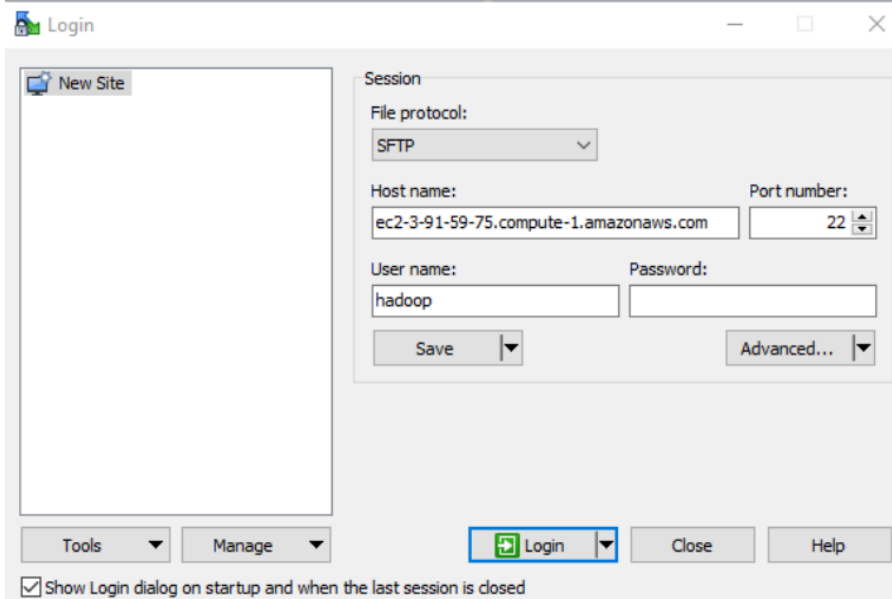
b. If you are using Windows, download [WinSCP]. Start it up and click on the "Advanced" button, then "Authentication" and browse to find your .ppk file.

After installing WinSCP, it may import from PuTTY, click Ok in that case.

1. In WinSCP, enter your EC2 instance public DNS name into **Host** *name* box.

2. Enter **"hadoop"** as a user name

3. Press the **Advanced** button to open advanced site settings dialog and go to **SSH -> Authentication**.

4. In the **Private key** file-box select your .ppk file

5. Submit these settings by clicking the **OK** button

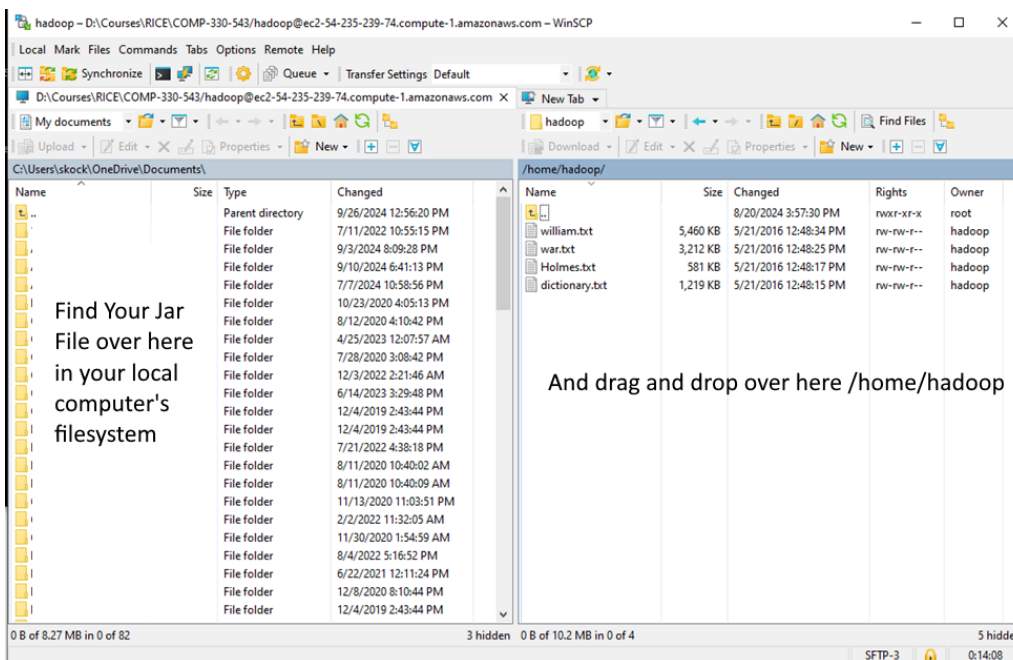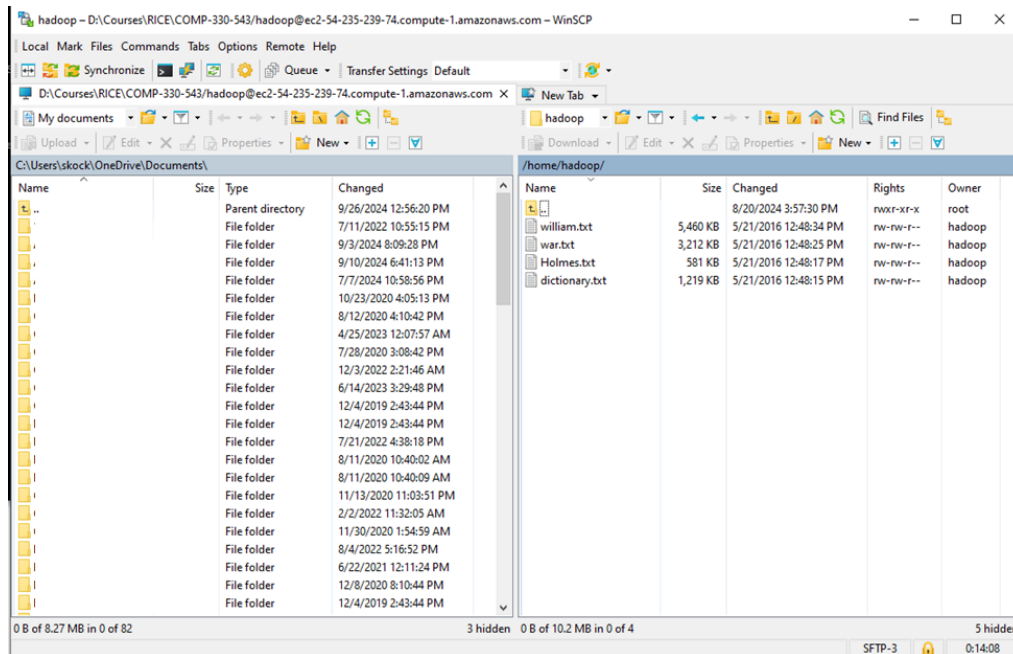6. Save your site settings using **Save** button.

/



c. Enter the primary node address, "hadoop" for the username, and press "Login"

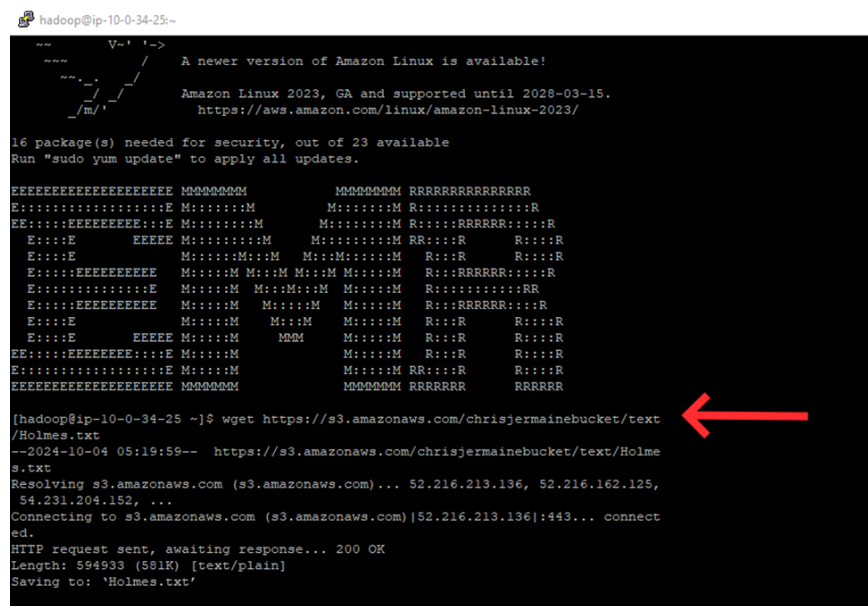d. WinSCP will connect to the cluster and you can use its graphical interface to transfer files to it.

   Notice that WinSCP shows the uploaded files(Task 4.e): under /home/hadoop

e. Next we will load some text files into the cluster to use with our WordCount program.

    i. Simply copy and paste the following commands into your cluster terminal

```
wget https://s3.amazonaws.com/luisguzmannateras/text/Holmes.txt
wget https://s3.amazonaws.com/luisguzmannateras/text/dictionary.txt
wget https://s3.amazonaws.com/luisguzmannateras/text/war.txt
wget https://s3.amazonaws.com/luisguzmannateras/text/william.txt
```

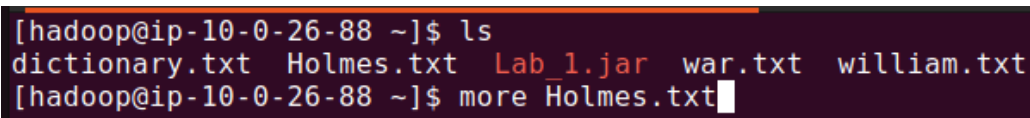If you are working on Windows, right click on your PuTTY terminal to paste the commands.



    ii. This will download four large text files into your cluster. You can look at the contents of these files by using the "more" or "less" commands (q to stop, space bar to see more).



    iii. Then, we need to **load this files into the HDFS**, run the following:

## Task 5: Run a Hadoop Program in the EMR Cluster

1. At last, we can run our WordCount program on the files in the HDFS

```
[hadoop@ip-10-0-26-88 ~]$ hadoop jar "Lab_1.jar" -r 8 words wordsOutput
```

2. To check the results, type:

```
[hadoop@ip-10-0-26-88 ~]$ hadoop fs -ls wordsOutput
Found 9 items
-rw-r--r--   1 hadoop hdfsadmingroup          0 2024-09-27 21:23 wordsOutput/_SUCCESS
-rw-r--r--   1 hadoop hdfsadmingroup      36982 2024-09-27 21:23 wordsOutput/part-r-00000
-rw-r--r--   1 hadoop hdfsadmingroup      36307 2024-09-27 21:23 wordsOutput/part-r-00001
-rw-r--r--   1 hadoop hdfsadmingroup      37143 2024-09-27 21:23 wordsOutput/part-r-00002
-rw-r--r--   1 hadoop hdfsadmingroup      37246 2024-09-27 21:23 wordsOutput/part-r-00003
-rw-r--r--   1 hadoop hdfsadmingroup      36980 2024-09-27 21:23 wordsOutput/part-r-00004
-rw-r--r--   1 hadoop hdfsadmingroup      36118 2024-09-27 21:23 wordsOutput/part-r-00005
-rw-r--r--   1 hadoop hdfsadmingroup      37234 2024-09-27 21:23 wordsOutput/part-r-00006
-rw-r--r--   1 hadoop hdfsadmingroup      37052 2024-09-27 21:23 wordsOutput/part-r-00007
[hadoop@ip-10-0-26-88 ~]$
```

3. Note that it's OK if yours looks a bit different. To copy some of the results from the HDFS into the master node of your cluster, type:

```
[hadoop@ip-10-0-26-88 ~]$ hadoop fs -get wordsOutput/part-r-00001 .
[hadoop@ip-10-0-26-88 ~]$
```

4. Now, you can take a look at the counts using the *more* command again, simply type:

```
[hadoop@ip-10-0-23-186 ~]$ head part-r-00000
a            28760
aaron    20
abate    12
abbreviated      1
abhorred         9
abide    34
abjure   4
absolute         55
absorption       2
abstains         1
```

5. **Copy and paste some of the counts you got to Canvas** in order to get points for the Lab.

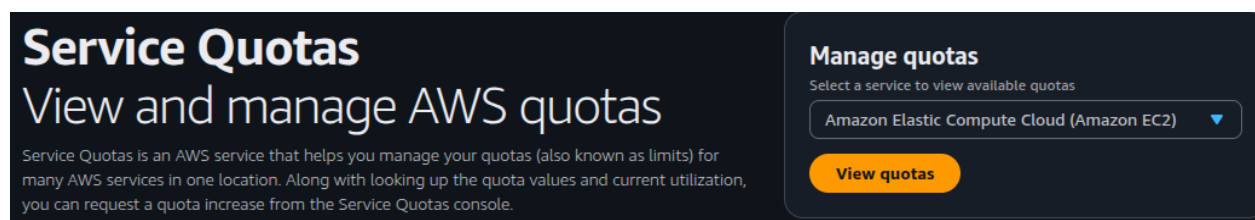# IMPORTANT TASK 6: SHUT DOWN YOUR CLUSTER

1. Never leave your cluster running when you are not using it, **YOU ARE BEING CHARGED!**

2. Remember to shut off your cluster by clicking the **"Terminate"** button in your EMR page.

3. After you kill them, **make sure** that they are dead. Go to **"All services" -> "EC2"** and look for **"Instances (running)"**. There should not be any. If they are still there, click on **"Running Instances"**. Then click the checkbox next to each of your machines, and under **"Actions" -> "Instance State"** choose **"Terminate"**. Only log out after you have verified from the EC2 page that you have no running instance.

4. It is a good idea to set a termination option under **"Cluster termination and node replacement"**. You can set your cluster to terminate after a specific amount of idle time. However, if you are connected through SSH into your cluster, it won't be considered idle and will continue to run.

## Extra Task: Request Quota Increase for G-type Instances

For Assignment 6, we will need to make use of G-type EC2 instances that have GPUs available. However, chances are you are going to have to request for a quota increase for this type of instance because new accounts usually are assigned a default quota of 0.

First, log in to your AWS console and use the search bar on top to search for **"Service Quotas"** which should take you to your service quotas dashboard.

Use the search box in the **"Manage quotas"** box to look for **"Amazon EC2"** and then click on "View quotas".



In the next page, type "G and VT" in the search bar and then click on **"Running On-Demand G and VT instances"** which should take you to a separate page where you can see your current allotted quota for G-type instances.

**Important Note:** Double check that you do **NOT** ask for a quota increase for **SPOT** instances. It should be for **ON-DEMAND** instances.

Your current quota will likely be zero and you will have to use the **"Request increase at account level"** button to make a request. Ask for a quota of 8 and click on Request.



Most of the time, quota-increase requests are approved fairly quickly (2~3 hrs) and without incident.

However, every once in a while Amazon will deny your request and you will have to talk to their customer support in order to get it approved. This process can be slow and take a few interactions with them. As such, we recommend you submit this request early so that you may be already approved by the time Assignment 6 comes along.