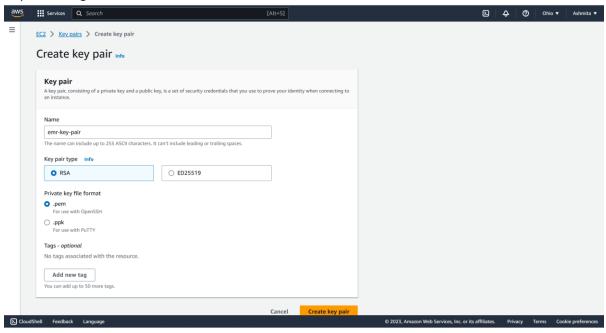# CSP 554

HOMEWORK 2

ASHMITA GUPTA (A20512498)

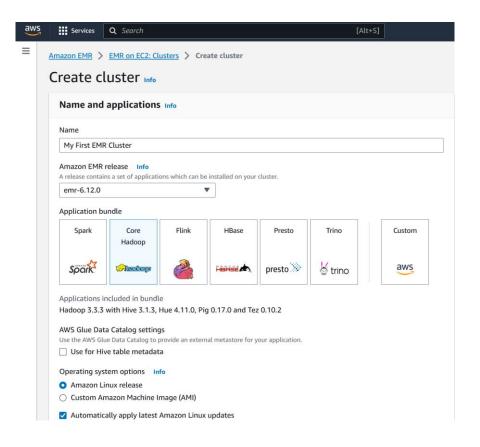**1.** Key Pair using the EC2 service



**2.** Key pair successfully created



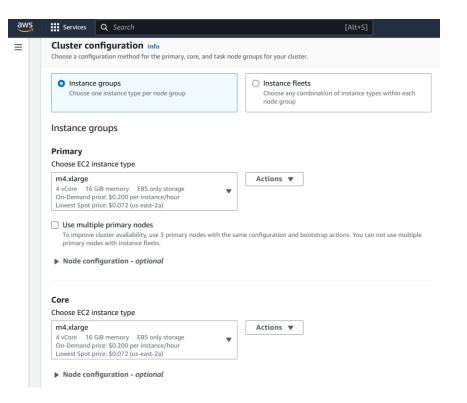**3.** Modifying permission for key pair

**4.** Creating Amazon EMR cluster by using steps given in the instruction

▶ Node configuration - *optional*

Add task instance group

You can add up to 48 more task instance groups.

▶ EBS root volume - *optional*

## Cluster scaling and provisioning option  Info

Amazon EMR console only supports EMR-managed scaling. To create a cluster with auto-scaling, use CLI or SDK.

Choose an option

| ○ Set cluster size manually | ○ Use EMR-managed scaling |
|---|---|
| Use this option if you know your workload patterns in advance. | Monitor key workload metrics so that EMR can optimize the cluster size and resource utilization. |

**Provisioning configuration**

Set the size of your core instance group. Amazon EMR attempts to provision this capacity when you launch your cluster.

| Name | Instance type | Instance(s) size | Use Spot purchasing option |
|---|---|---|---|
| Core | m4.xlarge | 1 | ☐ |

---

aws ::: Services  🔍 Search                                [Alt+S]

≡

Networking resources

We've already added the resources that you configured in the Networking section. Choose the VPC, subnet, and security groups that the service role can access.

Virtual Private Cloud (VPC)

Choose one or more VPCs ▼

-
vpc-0a5aee1c03ba24438  ✕

Subnet

Choose one or more subnets ▼

-
subnet-0017279326b7f90fc  ✕

Security group

Choose one or more security groups ▼

**EC2 instance profile for Amazon EMR**

The instance profile assigns a role to every EC2 instance in a cluster. The instance profile must specify a role that can access the resources for your steps and bootstrap actions.

| ○ Choose an existing instance profile | ● Create an instance profile |
|---|---|
| Select a default role or a custom instance profile with IAM policies attached so that your cluster can interact with your resources in Amazon S3. | Let Amazon EMR create a new instance profile so that you can specify a custom set of resources for it to access in Amazon S3. |

**S3 bucket access**  Info

○ Specific S3 buckets or prefixes in your account  Info
Choose the buckets or prefixes that you want this instance profile to access.

● All S3 buckets in this account with read and write access
Grant the instance profile access to all buckets that have read and write access enabled in your account.
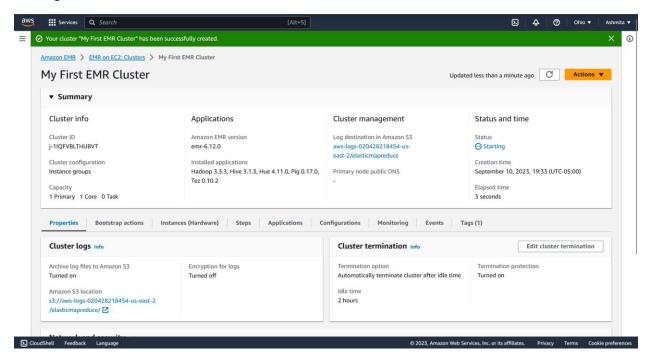
[>_] CloudShell    Feedback    Language

**Cluster created:**

Starting state:



Waiting state:

**Inbound rules**:
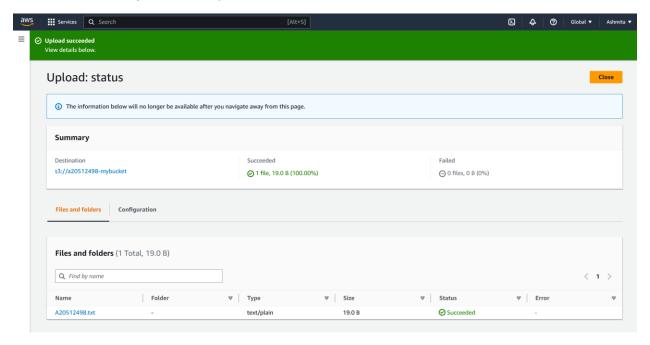


## Connecting to the Primary Node Using SSH

**SCP command** executed to copy file from local machine to the home directory of Hadoop master node account:
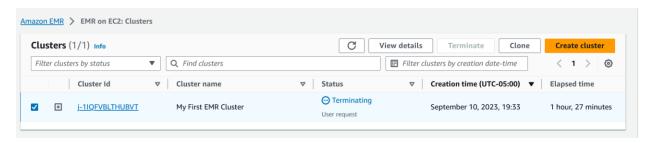
```
gashm@Ashmita MINGW64 ~/OneDrive/Desktop/Big Data
$ scp -i emr-key-pair.pem ashmitagupta.txt hadoop@ec2-52-15-223-17.us-east-2.compute.amazonaws.com:/home/hadoop
ashmitagupta.txt

gashm@Ashmita MINGW64 ~/OneDrive/Desktop/Big Data
$
```

**A20512498.txt file uploaded** to my bucket:



**Cluster created as below:**



9. **(2 points) Execute the following hdfs command to list the files or directories that are listed (also indicating which is a file and which a directory):**

   **hadoop fs –ls /**

Take a screen snapshot of names of the files or directories that are listed and include it in your assignment submission:

Ans. Command:  **hadoop fs -ls /**

```
MINGW64:/c/Users/gashm/OneDrive/Desktop/Big Data
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -ls /
Found 4 items
drwxr-xr-x   - hdfs hdfsadmingroup          0 2023-09-11 00:40 /apps
drwxrwxrwt   - hdfs hdfsadmingroup          0 2023-09-11 00:42 /tmp
drwxr-xr-x   - hdfs hdfsadmingroup          0 2023-09-11 01:40 /user
drwxr-xr-x   - hdfs hdfsadmingroup          0 2023-09-11 00:40 /var
```

**10. (2 points) Execute a command (you needed to figure out which one) to list the files and directories under the hdfs directory listed below:**

**/user**

Write down the command you executed and also take a screen snapshot of names of the files or directories that are listed and include it in your assignment submission.

Ans. **hadoop fs -ls /user**

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -ls /user
Found 6 items
drwxrwxrwx   - hadoop hdfsadmingroup          0 2023-09-11 00:40 /user/hadoop
drwxr-xr-x   - mapred mapred                  0 2023-09-11 00:40 /user/history
drwxrwxrwx   - hdfs   hdfsadmingroup          0 2023-09-11 00:40 /user/hive
drwxrwxrwx   - hue    hue                     0 2023-09-11 00:40 /user/hue
drwxrwxrwx   - oozie  oozie                   0 2023-09-11 00:42 /user/oozie
drwxrwxrwx   - root   hdfsadmingroup          0 2023-09-11 00:40 /user/root
```

**11. (2 points) Execute a command to create the following HDFS directory:**

**/user/csp554**

Record the command you executed and include it in your assignment submission.

Ans. Command: **hadoop fs -mkdir /user/csp554**

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -ls /user
Found 6 items
drwxrwxrwx   - hadoop hdfsadmingroup          0 2023-09-11 00:40 /user/hadoop
drwxr-xr-x   - mapred mapred                  0 2023-09-11 00:40 /user/history
drwxrwxrwx   - hdfs   hdfsadmingroup          0 2023-09-11 00:40 /user/hive
drwxrwxrwx   - hue    hue                     0 2023-09-11 00:40 /user/hue
drwxrwxrwx   - oozie  oozie                   0 2023-09-11 00:42 /user/oozie
drwxrwxrwx   - root   hdfsadmingroup          0 2023-09-11 00:40 /user/root
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -mkdir /user/csp554
```

**12. (2 points) Execute a command to create the following HDFS directory:**

**/user/csp554-2**

Ans. Command: **hadoop fs -mkdir /user/csp554-2**

Both csp554 and csp554-2 HDFS directories created as shown below:

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -mkdir /user/csp554-2
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -ls /user
Found 8 items
drwxr-xr-x   - hadoop hdfsadmingroup          0 2023-09-11 01:41 /user/csp554
drwxr-xr-x   - hadoop hdfsadmingroup          0 2023-09-11 01:42 /user/csp554-2
drwxrwxrwx   - hadoop hdfsadmingroup          0 2023-09-11 00:40 /user/hadoop
drwxr-xr-x   - mapred mapred                  0 2023-09-11 00:40 /user/history
drwxrwxrwx   - hdfs   hdfsadmingroup          0 2023-09-11 00:40 /user/hive
drwxrwxrwx   - hue    hue                     0 2023-09-11 00:40 /user/hue
drwxrwxrwx   - oozie  oozie                   0 2023-09-11 00:42 /user/oozie
drwxrwxrwx   - root   hdfsadmingroup          0 2023-09-11 00:40 /user/root
```

Record the command you executed and include it in your assignment submission.

**13. (2 points) Execute a command that copies a given local file to the given hdfs directory :**

**Source local file: /home/hadoop/myname.txt   (where the actual name is your name as described above)**

Destination HDFS directory: /user/csp554

Ans. Command: **hadoop fs -put /home/hadoop/ashmitagupta.txt /user/csp554**

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -put /home/hadoop/ashmitagupta.txt /user/csp554
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -ls /user/csp554
Found 1 items
-rw-r--r--   1 hadoop hdfsadmingroup         21 2023-09-11 01:43 /user/csp554/ashmitagupta.txt
```

**14. (2 points) Copy a file from one hdfs directory to another hdfs directory and write down the command**

Source hdfs file: /user/csp554/myname.txt (where the actual name is your name as described above)

Destination HDFS directory: /user/csp554-2

Ans. Command: **hadoop fs -cp /user/csp554/ashmitagupta.txt /user/csp554-2**

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -cp /user/csp554/ashmitagupta.txt /user/csp554-2
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -ls /user/csp554-2
Found 1 items
-rw-r--r--   1 hadoop hdfsadmingroup         21 2023-09-11 01:44 /user/csp554-2/ashmitagupta.txt
```

**15.  (2 points) Copy the object myid.txt you uploaded to an S3 bucket into the Hadoop master node Linux file system. The actual object includes your student id as above.**

Ans. Command: **aws s3 cp s3://a20512498-mybucket/A20512498.txt /home/hadoop/A20512498.txt**

```
[hadoop@ip-172-31-37-54 ~]$ aws s3 cp s3://a20512498-mybucket/A20512498.txt /home/hadoop/A20512498.txt
download: s3://a20512498-mybucket/A20512498.txt to ./A20512498.txt
[hadoop@ip-172-31-37-54 ~]$ ls
A20512498.txt  ashmitagupta.txt
```

**16. (2 points) Copy the same object myid.txt you created in an S3 bucket into HDFS into the directory /users/csp554**

> hadoop fs -cp s3://mybucket/myid.txt    hdfs:///user/csp554-2

After you executed the above command, execute another command (you needed to figure out which one) to list the files and directories under the hdfs directory listed below:

/user/csp554-2

Write down the command you executed and also take a screen snapshot of names of the files or directories that are listed and include it in your assignment submission.

Ans. Commands:

**hadoop fs -cp s3://a20512498-mybucket/A20512498.txt hdfs:///user/csp554-2**

**hadoop fs -ls /user/csp554-2**

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -cp s3://a20512498-mybucket/A20512498.txt    hdfs:///user/csp554-2
2023-09-11 01:53:30,768 INFO s3n.S3NativeFileSystem: Opening 's3://a20512498-mybucket/A20512498.txt' for reading
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -ls /user/csp554-2
Found 2 items
-rw-r--r--   1 hadoop hdfsadmingroup         19 2023-09-11 01:53 /user/csp554-2/A20512498.txt
-rw-r--r--   1 hadoop hdfsadmingroup         21 2023-09-11 01:44 /user/csp554-2/ashmitagupta.txt
```

**17. (2 points) Execute a command to show the contents of the myid.txt file in the hdfs directory /user/csp554-2**

**Clue: look up about how to use the "cat" command in the file system shell document.**

Write down the command you executed and also take a screen snapshot of the listed content of the file and include it in your assignment submission.

Ans. Command: **hadoop fs -cat /user/csp554-2/A20512498.txt**

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -cat /user/csp554-2/A20512498.txt
this is the id file
```

**18. (2 points) Execute a command to remove the myid.txt file in the hdfs directory /user/csp554-2**

Clue: look up about how to use the "rm" command in the file system shell document.

Write down the command you executed, then list the content of the /user/csp554-2 HDFS directory and take a screen snapshot of the listed content of the directory and include it in your assignment submission.

Ans. Command: **hadoop fs -rm /user/csp554-2/A20512498.txt**

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -rm /user/csp554-2/A20512498.txt
Deleted /user/csp554-2/A20512498.txt
```

**Screenshot of contents in directory:**

Command: **hadoop fs -ls /user/CS554-2**

```
[hadoop@ip-172-31-37-54 ~]$ hadoop fs -ls /user/csp554-2
Found 1 items
-rw-r--r--   1 hadoop hdfsadmingroup         21 2023-09-11 01:44 /user/csp554-2/ashmitagupta.txt
```