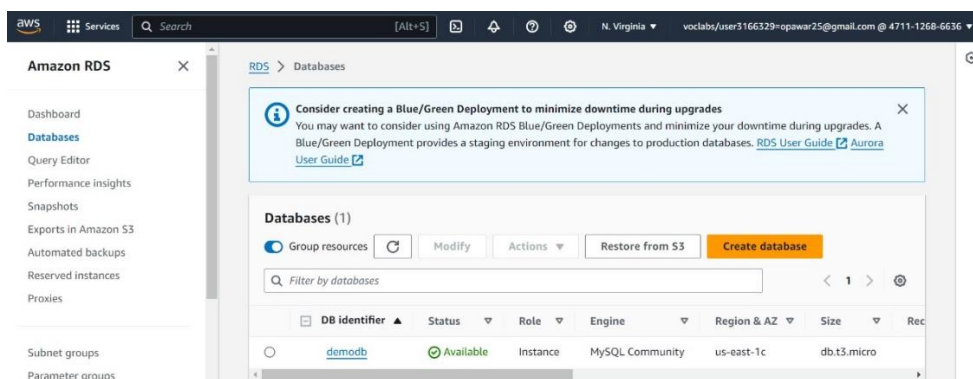# TASK 1

**Create an RDS instance in your AWS account and upload the data to the RDS instance. Since the dataset is huge, you need to upload the data from only two files**
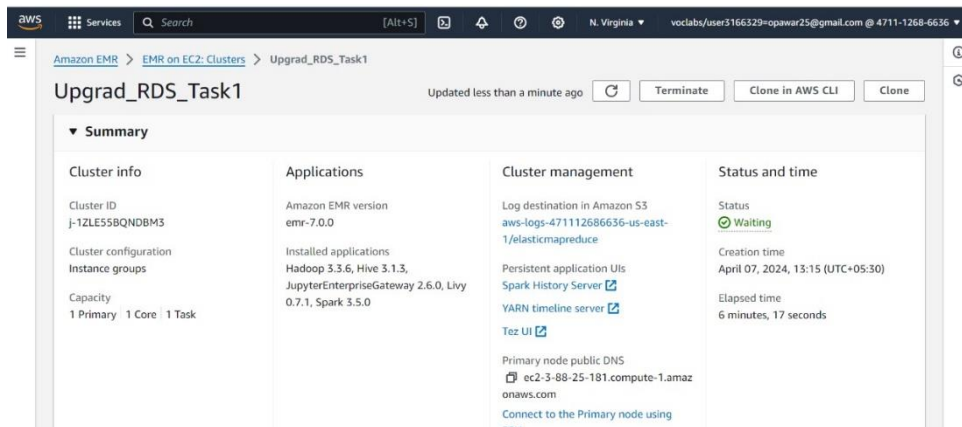
**Step 1.** Creating an RDS instance

- Search for RDS on Amazon Web Server.
- Click on create database under database menu.
- Further choose "Standard create" from the database creation method.
- Next click on MySQL from Engine Options to run a MySQL engine.
- Click on Free tier under the Templates option to use the engine for free.
- Now, give a name to your database under DB cluster identifier.
- Next, give a public access to your database by clicking on yes under the Public Access option.
- Ensure to make the database password authentication under database authentication for security.
- At last, click on Create Database to create the RDS Database.
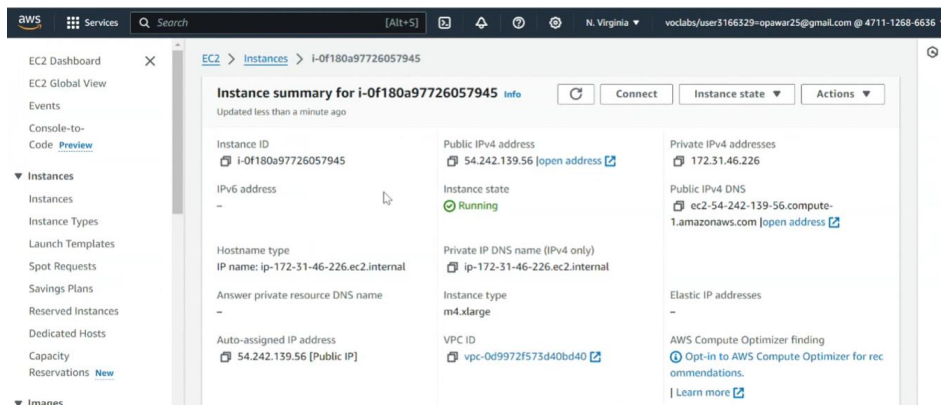- The database will take 5-10 minutes to get active.



**Step 2.** Creating an EMR Cluster

- Search for EMR on Amazon Web Server.
- Click on "create cluster" to create a EMR cluster
- Now, give name to your cluster under the name column.
- Next, tick the required applications under the application bundle for example Hadoop, HBase, Sqoop etc.
- Further choose an EC2 instance type for Primary node.
- Next enter the Size required (30 GiB) under EBS root volume.
- Now, select the EC2 key pair for SSH to the cluster by creating a new key pair.
- Now at last, select the service role and instance profile as default role
- Finally, click on Create Cluster to create an EMR Cluster
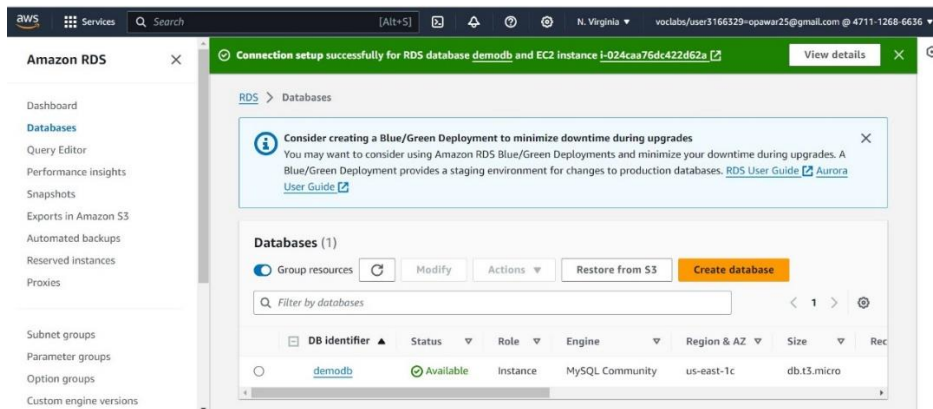- The cluster will take 5-10 minutes to get active

**Step 3**. Creating EC2 Instance

- Search for EC2 on Amazon Web Server
- Click on Launch instances to create a EC2 instance
- Now, give a name to your EC2 instance under the Name and Tags option
- Now, Select a desired Amazon Machine Image for example, Amazon Linux, macOS etc.
- Now select the instance type which is required for example t2.micro with 1 vCPU and 1GiB memory
- Now, select a key pair name from the drop-down menu or create a new key pair
- At last, click on Launch Instance to create a EC2 Instance
- The creation of the instance may take around 5-10 minutes to get active
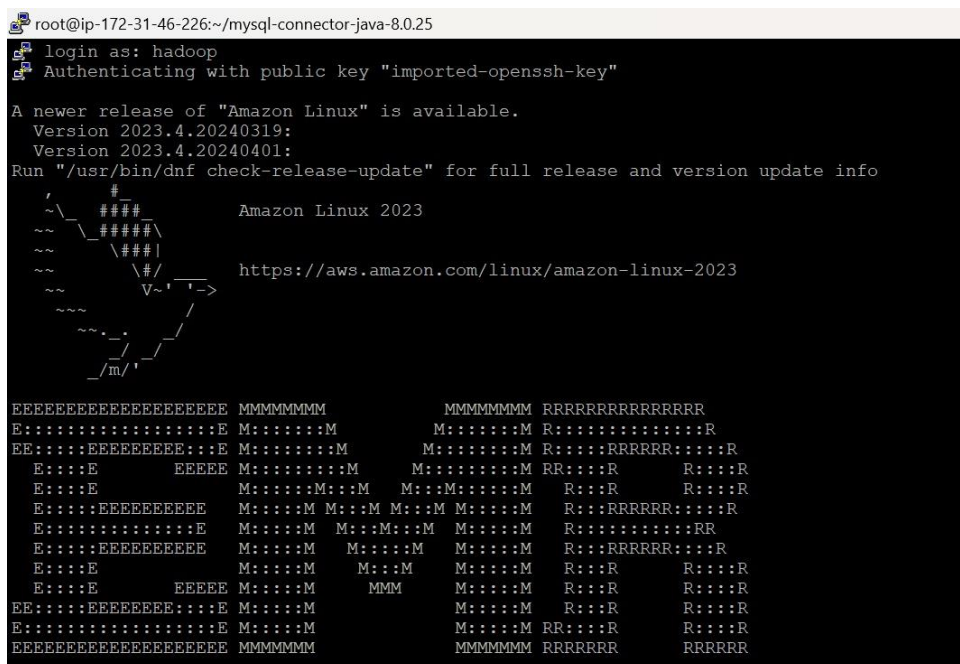
**Step 4**. Connecting a setup between RDS database and EC2 instance



**Step 5**. Connecting EMR instance with the local computer using Putty.

- Firstly, copy the Primary node public DNS from the EMR cluster and paste it as a host name or IP Address in Putty
- Now, click on credentials under the Auth menu from SSH and browse for private key file for authentication.
- After inserting the private key, click open and the connection between EMR instance with the local computer will be executed.

**Step 6**: Connecting to MySQL database remotely through putty using following commands

1. wget https://de-mysql-connector.s3.amazonaws.com/mysql-connector-java-8.0.25.tar.gz
2. tar -xvf mysql-connector-java-8.0.25.tar.gz
3. cd mysql-connector-java-8.0.25/
4. sudo cp mysql-connector-java-8.0.25.jar /usr/lib/sqoop/lib/

```
[hadoop@ip-172-31-46-226 ~]$ sudo -i

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::::EEEEEEEEE:::E  M:::::M              M:::::M  R:::R       R::::R
E::::::::::::::::::E M:::::M              M:::::M RR::::R     R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMM              MMMMMM RRRRRRR       RRRRRR

[root@ip-172-31-46-226 ~]# ^[[200~wget https://de-mysql-connector.s3.amazonaws.c
om/mysql-connector-java-8.0.25.tar.gz~
-bash: $'\E[200~wget': command not found
[root@ip-172-31-46-226 ~]# wget https://de-mysql-connector.s3.amazonaws.com/mysq
l-connector-java-8.0.25.tar.gz
--2024-04-07 18:53:30--  https://de-mysql-connector.s3.amazonaws.com/mysql-conne
ctor-java-8.0.25.tar.gz
Resolving de-mysql-connector.s3.amazonaws.com (de-mysql-connector.s3.amazonaws.c
om)... 52.216.217.233, 52.217.230.145, 52.217.131.201, ...
Connecting to de-mysql-connector.s3.amazonaws.com (de-mysql-connector.s3.amazona
ws.com)|52.216.217.233|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 4079310 (3.9M) [application/x-gzip]
Saving to: 'mysql-connector-java-8.0.25.tar.gz'

mysql-connector-jav 100%[===================>]   3.89M  --.-KB/s    in 0.03s

2024-04-07 18:53:30 (137 MB/s) - 'mysql-connector-java-8.0.25.tar.gz' saved [407
9310/4079310]
```

```
2024-04-07 19:14:18 (40.5 MB/s) - 'yellow_tripdata_2017-01.csv' saved [914029540
/914029540]

[hadoop@ip-172-31-46-226 ~]$ ls
yellow_tripdata_2017-01.csv
[hadoop@ip-172-31-46-226 ~]$ wget https://de-mysql-connector.s3.amazonaws.com/my
sql-connector-java-8.0.25.tar.gz
--2024-04-07 19:15:21--  https://de-mysql-connector.s3.amazonaws.com/mysql-conne
ctor-java-8.0.25.tar.gz
Resolving de-mysql-connector.s3.amazonaws.com (de-mysql-connector.s3.amazonaws.c
om)... 52.217.172.137, 52.216.39.105, 52.216.246.100, ...
Connecting to de-mysql-connector.s3.amazonaws.com (de-mysql-connector.s3.amazona
ws.com)|52.217.172.137|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 4079310 (3.9M) [application/x-gzip]
Saving to: 'mysql-connector-java-8.0.25.tar.gz'

mysql-connector-jav 100%[===================>]   3.89M  --.-KB/s    in 0.05s

2024-04-07 19:15:21 (83.6 MB/s) - 'mysql-connector-java-8.0.25.tar.gz' saved [40
79310/4079310]

[hadoop@ip-172-31-46-226 ~]$ tar -xvf mysql-connector-java-8.0.25.tar.gz
```

**Step 7:** Uploading the data from 2 given files



**Step 8:** Creating schema using MySQL connection and also creating table trip_log

.

Step 9: Finally checking the count of the given task 1 files.



```
/863487050]

[hadoop@ip-172-31-14-221 ~]$ ls
mysql-connector-java-8.0.25        yellow_tripdata_2017-01.csv
mysql-connector-java-8.0.25.tar.gz  yellow_tripdata_2017-02.csv
[hadoop@ip-172-31-14-221 ~]$ mysql -h demodb.crog2ckcif6b.us-east-1.rds.amazonaw
s.com -P 3306 -u admin -p
Enter password:
Welcome to the MariaDB monitor.  Commands end with ; or \g.
Your MySQL connection id is 28
Server version: 8.0.35 Source distribution

Copyright (c) 2000, 2018, Oracle, MariaDB Corporation Ab and others.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

MySQL [(none)]> use taxi_records;
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
MySQL [taxi_records]> show tables;
+----------------------+
| Tables_in_taxi_records |
+----------------------+
| trip_log             |
+----------------------+
1 row in set (0.002 sec)

MySQL [taxi_records]> LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-
02.csv' INTO TABLE trip_log
    -> FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' IGNORE 1 LINES;
Query OK, 9169775 rows affected, 65535 warnings (1 min 59.713 sec)
Records: 9169775  Deleted: 0  Skipped: 0  Warnings: 9169775

MySQL [taxi_records]> SELECT COUNT(*) FROM taxi_records.trip_log;
+----------+
| COUNT(*) |
+----------+
| 18880595 |
+----------+
1 row in set (53.893 sec)

MySQL [taxi_records]>
```