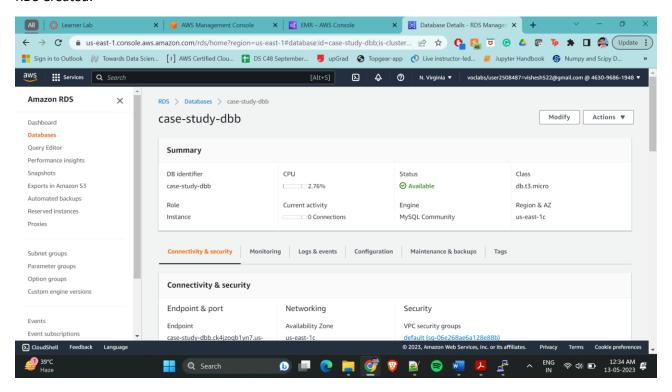# Task 1

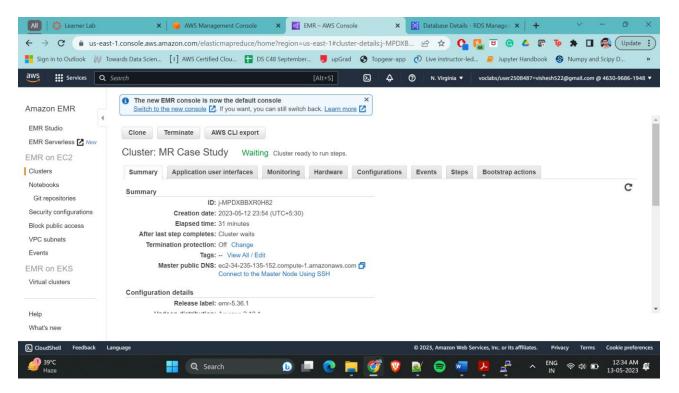## Create an RDS instance in your AWS account and upload the data to the RDS instance

1) We have created RDS instance from Learner Lab and Launched EMR cluster.

2) Downloaded the required files for this task - **yellow_tripdata_2017-01.csv** & **yellow_tripdata_2017-02.csv**

3) Connect RDS instance from EMR and create table and load records from csv files.
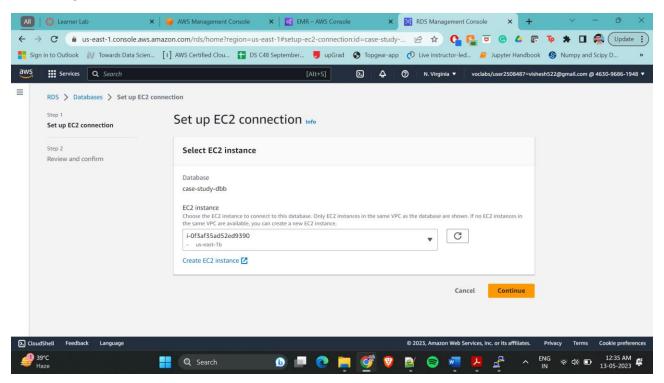
4) Loading the data to RDS

RDS Created:
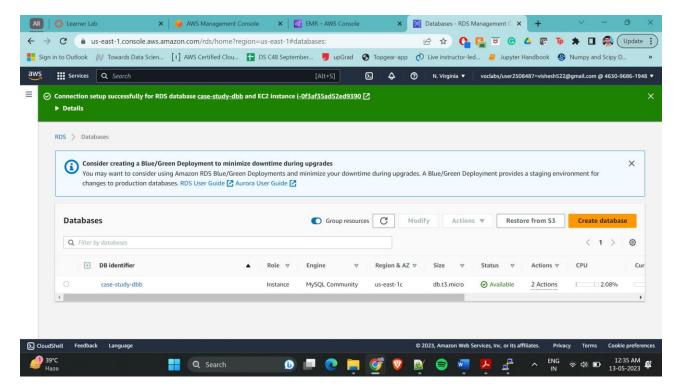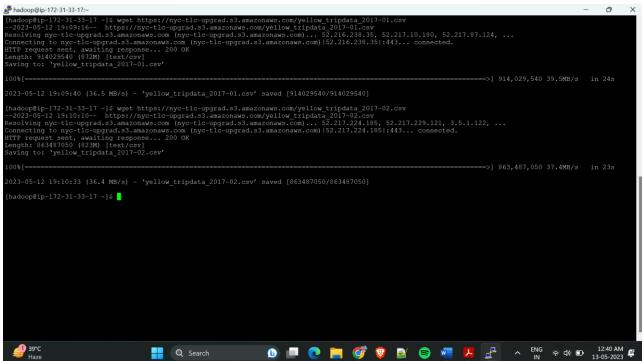


EMR Cluster:

Connecting RDS with EMR EC2 instance:

Connecting EMR instance with PuTTy and then downloading **yellow_tripdata_2017-01.csv** & **yellow_tripdata_2017-02.csv** :



**Connecting RDS with EMR Instance:**

**Hostname:** case-study-dbb.ck4jzoqb1yn7.us-east-1.rds.amazonaws.com

mysql -h case-study-dbb.ck4jzoqb1yn7.us-east-1.rds.amazonaws.com -P 3306 -u admin -p

```
root@ip-172-31-33-17:/home/hadoop
[root@ip-172-31-33-17 hadoop]# mysql -h case-study-dbb.ck4jzoqblyn7.us-east-1.rds.ama
Enter password:
Welcome to the MariaDB monitor.  Commands end with ; or \g.
Your MySQL connection id is 1297
Server version: 8.0.32 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)]> create database YellowTaxi;
Query OK, 1 row affected (0.04 sec)

MySQL [(none)]> use YellowTaxi;
Database changed
MySQL [YellowTaxi]> █
```

create database YellowTaxi;

user YellowTaxi;

Creating table--

CREATE TABLE taxi

(

vendorID INT,

tpep_pickup_datetime DATETIME,

tpep_dropoff_datetime DATETIME,

passenger_count INT,

trip_distance DOUBLE,

puLocationID INT,

doLocationID INT,

rateCodeID INT,

store_and_fwd_flag VARCHAR(255),

payment_type INT,

fare_amount DOUBLE,

extra DOUBLE,

mta_tax DOUBLE,

improvement_surcharge DOUBLE,

tip_amount DOUBLE,

tolls_amount DOUBLE,

total_amount DOUBLE,

congestion_Surcharge DOUBLE,

airport_fee DOUBLE

);

```
root@ip-172-31-33-17:/home/hadoop
MySQL [YellowTaxi]> CREATE TABLE taxi
    -> (
    -> vendorID INT,
    -> tpep_pickup_datetime DATETIME,
    -> tpep_dropoff_datetime DATETIME,
    -> passenger_count INT,
    -> trip_distance DOUBLE,
    -> puLocationID INT,
    -> doLocationID INT,
    -> rateCodeID INT,
    -> store_and_fwd_flag VARCHAR(255),
    -> payment_type INT,
    -> fare_amount DOUBLE,
    -> extra DOUBLE,
    -> mta_tax DOUBLE,
    -> improvement_surcharge DOUBLE,
    -> tip_amount DOUBLE,
    -> tolls_amount DOUBLE,
    -> total_amount DOUBLE,
    -> congestion_Surcharge DOUBLE,
    -> airport_fee DOUBLE
    -> );
Query OK, 0 rows affected (0.38 sec)

MySQL [YellowTaxi]>
```

**Load Data into Above table**

LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-01.csv'

INTO TABLE taxi

FIELDS TERMINATED BY ','

LINES TERMINATED BY '\n'

IGNORE 1 LINES;

LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-02.csv'

INTO TABLE taxi

FIELDS TERMINATED BY ','

LINES TERMINATED BY '\n'

IGNORE 1 LINES;

```
MySQL [YellowTaxi]> LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-01.csv'
    -> INTO TABLE taxi
    -> FIELDS TERMINATED BY ','
    -> LINES TERMINATED BY '\n'
    -> IGNORE 1 LINES;
Query OK, 9710820 rows affected, 65535 warnings (2 min 19.95 sec)
Records: 9710820  Deleted: 0  Skipped: 0  Warnings: 29132460

MySQL [YellowTaxi]> LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-02.csv'
    -> INTO TABLE taxi
    -> FIELDS TERMINATED BY ','
    -> LINES TERMINATED BY '\n'
    -> IGNORE 1 LINES;
Query OK, 9169775 rows affected, 65535 warnings (2 min 2.74 sec)
Records: 9169775  Deleted: 0  Skipped: 0  Warnings: 27509325
```