# Data Ingestion Tasks
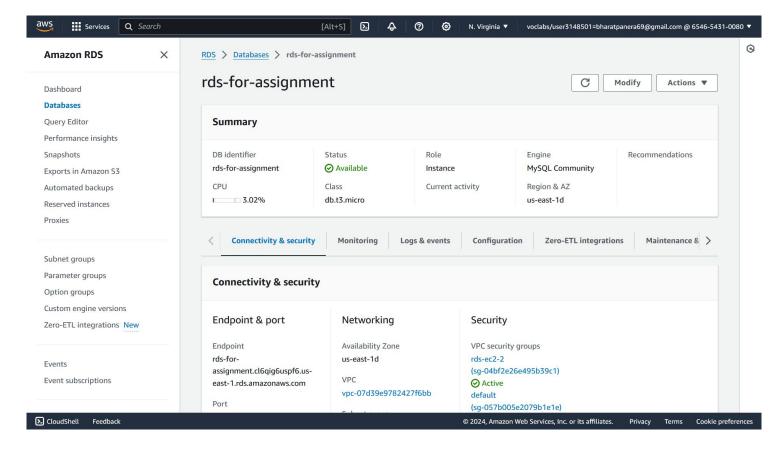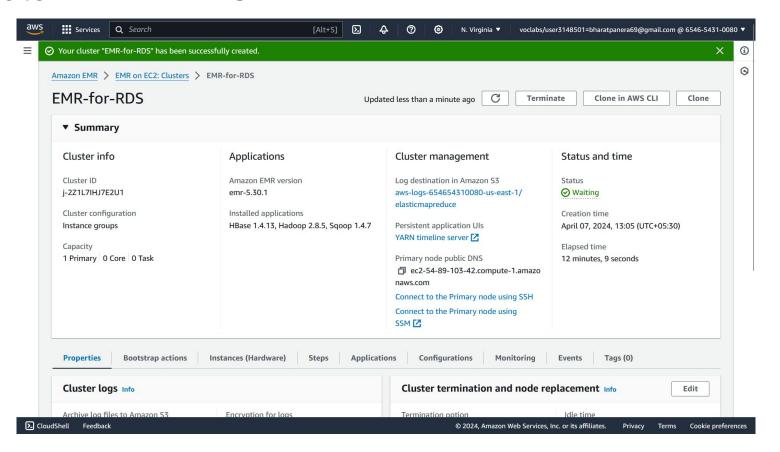
Task 1: Create RDS and import the data into it

# Create RDS instance in AWS
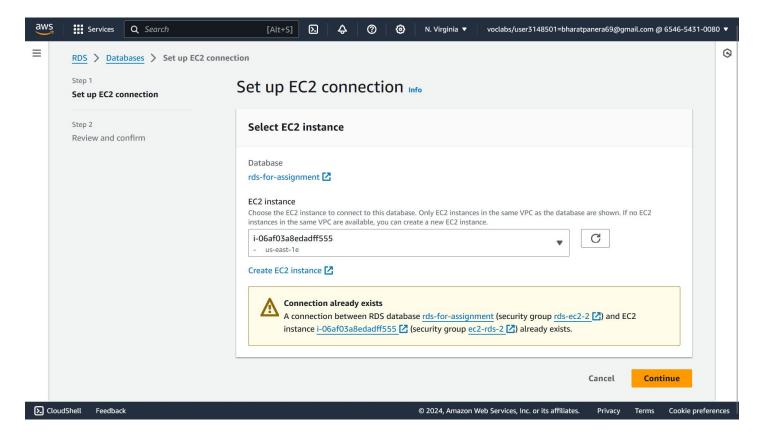
# Create EMR in AWS

# Connect RDS to EMR

- To connect RDS to EMR select the RDS and click the "Actions" button
- Then click "set up EC2 connection"
- Choose an appropriate EC2 instance and click "continue"
- Then review and confirm
- It will connect the RDS to EMR cluster
- Then login to EMR from local machine and try to access mysql

      **Command:** mysql -h rds-for-assignment.cl6qig6uspf6.us-east-1.rds.amazonaws.com -P 3306 -u admin -p

# Connect RDS to EMR

# Login to EMR

# Access the MySQL from EMR

```
[ec2-user@ip-172-31-51-243 ~]$ mysql -h rds-for-assignment.cl6qig6uspf6.us-east-1.rds.amazonaws.com -P 3306 -u admin -p12345678
Welcome to the MariaDB monitor.  Commands end with ; or \g.
Your MySQL connection id is 30
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)]> show databases;
+--------------------+
| Database           |
+--------------------+
| information_schema |
| mysql              |
| performance_schema |
| sys                |
+--------------------+
4 rows in set (0.00 sec)

MySQL [(none)]> 
```

# Create database "taxi_data" & table "taxi_tripdata"

- We will create a database "taxi_data"
  - **Syntax:** CREATE DATABASE taxi_data;
- Now we will create a table "taxi_tripdata"
  - **Syntax:** CREATE TABLE taxi_tripdata (vendorid INT, tpep_pickup_datetime TIMESTAMP, tpep_dropoff_datetime TIMESTAMP, passenger_count INT, trip_distance DECIMAL(10,2), ratecodeid INT, store_and_fwd_flag BOOLEAN, puocationId INT, doLocationid INT, payment_type INT, fare_amount DOUBLE, extra DOUBLE, mta_tax DOUBLE, tip_amount DOUBLE, tolls_amount DOUBLE, improvement_surcharge DOUBLE, total_amount DOUBLE, congestion_surcharge DOUBLE, airport_fee DOUBLE);

# Create database "taxi_data" & table "taxi_tripdata"

```
MySQL [taxi_data]> CREATE TABLE taxi_tripdata (
    ->     vendorid INT,
    ->     tpep_pickup_datetime TIMESTAMP,
    ->     tpep_dropoff_datetime TIMESTAMP,
    ->     passenger_count INT,
    ->     trip_distance DECIMAL(10,2) ,
    ->     ratecodeid INT,
    ->     store_and_fwd_flag BOOLEAN,
    ->     puocationId INT,
    ->     doLocationid INT,
    ->     payment_type INT,
    ->     fare_amount DOUBLE,
    ->     extra DOUBLE,
    ->     mta_tax DOUBLE,
    ->     tip_amount DOUBLE,
    ->     tolls_amount DOUBLE,
    ->     improvement_surcharge DOUBLE,
    ->     total_amount DOUBLE,
    ->     congestion_surcharge DOUBLE,
    ->     airport_fee DOUBLE
    -> );
Query OK, 0 rows affected (0.03 sec)

MySQL [taxi_data]> SHOW TABLES;
+--------------------+
| Tables_in_taxi_data |
+--------------------+
| taxi_tripdata      |
+--------------------+
1 row in set (0.00 sec)

MySQL [taxi_data]> DESC taxi_tripdata;
+-----------------------+--------------+------+-----+---------+-------+
| Field                 | Type         | Null | Key | Default | Extra |
+-----------------------+--------------+------+-----+---------+-------+
| vendorid              | int          | YES  |     | NULL    |       |
| tpep_pickup_datetime  | timestamp    | YES  |     | NULL    |       |
| tpep_dropoff_datetime | timestamp    | YES  |     | NULL    |       |
| passenger_count       | int          | YES  |     | NULL    |       |
| trip_distance         | decimal(10,2)| YES  |     | NULL    |       |
| ratecodeid            | int          | YES  |     | NULL    |       |
| store_and_fwd_flag    | tinyint(1)   | YES  |     | NULL    |       |
| puocationId           | int          | YES  |     | NULL    |       |
| doLocationid          | int          | YES  |     | NULL    |       |
| payment_type          | int          | YES  |     | NULL    |       |
| fare_amount           | double       | YES  |     | NULL    |       |
| extra                 | double       | YES  |     | NULL    |       |
| mta_tax               | double       | YES  |     | NULL    |       |
| tip_amount            | double       | YES  |     | NULL    |       |
| tolls_amount          | double       | YES  |     | NULL    |       |
| improvement_surcharge | double       | YES  |     | NULL    |       |
| total_amount          | double       | YES  |     | NULL    |       |
| congestion_surcharge  | double       | YES  |     | NULL    |       |
| airport_fee           | double       | YES  |     | NULL    |       |
+-----------------------+--------------+------+-----+---------+-------+
19 rows in set (0.00 sec)
```

# Download the required files into EMR

- Download the yellow_tripdata_2017-01.csv & yellow_tripdata_2017-02.csv files into EMR.
    - **Commands:**
        - wget https://nyc-tlc-upgrad.s3.amazonaws.com/yellow_tripdata_2017-01.csv
        - wget https://nyc-tlc-upgrad.s3.amazonaws.com/yellow_tripdata_2017-02.csv

```
[ec2-user@ip-172-31-51-243 ~]$ du -sch *
1008M    yellow_tripdata_2017-01.csv
824M     yellow_tripdata_2017-02.csv
1.8G     total
```

# Load the taxi data into MySQL table

- Take absolute path of the CSV files using command "readlink -f *csv"
  - /home/ec2-user/yellow_tripdata_2017-01.csv
  - /home/ec2-user/yellow_tripdata_2017-02.csv
- Load the .csv files into MySQL table using below query.
  - LOAD DATA LOCAL INFILE '/home/ec2-user/yellow_tripdata_2017-01.csv' INTO TABLE taxi_data.taxi_tripdata FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' IGNORE 1 LINES;
  - LOAD DATA LOCAL INFILE '/home/ec2-user/yellow_tripdata_2017-02.csv' INTO TABLE taxi_data.taxi_tripdata FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' IGNORE 1 LINES;
- Then we can cross check with counting the number of imported records with number of lines present in csv files.
  - To check number of lines in csv files: "wc -l *csv"
  - To check imported records in MySQL table:
    - SELECT COUNT(*) FROM taxi_data.taxi_tripdata;
  - It should be totale of both the csv files.

# Load the taxi data into MySQL table

```
[ec2-user@ip-172-31-51-243 ~]$ wc -l *csv
  9710821 yellow_tripdata_2017-01.csv
  9169776 yellow_tripdata_2017-02.csv
 18880597 total
```

```
MySQL [taxi_data]> SELECT COUNT(*) FROM taxi_tripdata;
+----------+
| COUNT(*) |
+----------+
| 18880595 |
+----------+
1 row in set (55.25 sec)
```