

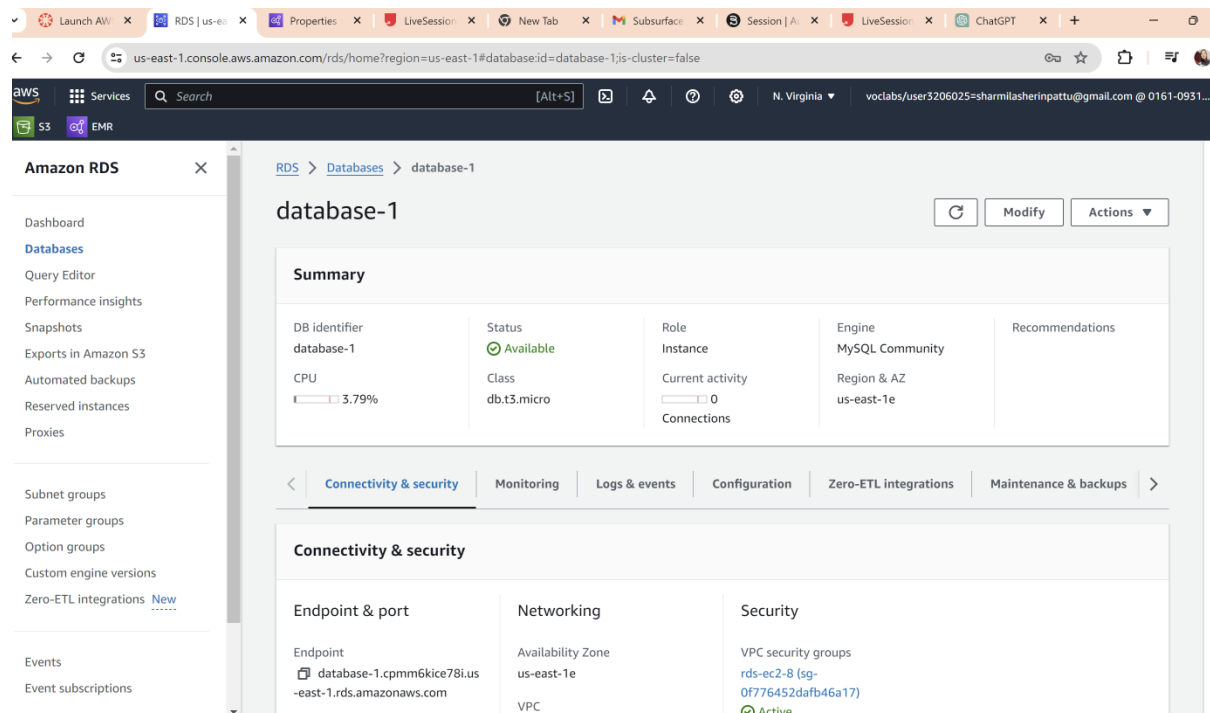
# Task 1: AWS Environment Setup and Data Upload

**Objective:** Creating an RDS instance in AWS and uploading specific data files.

## Instructions:

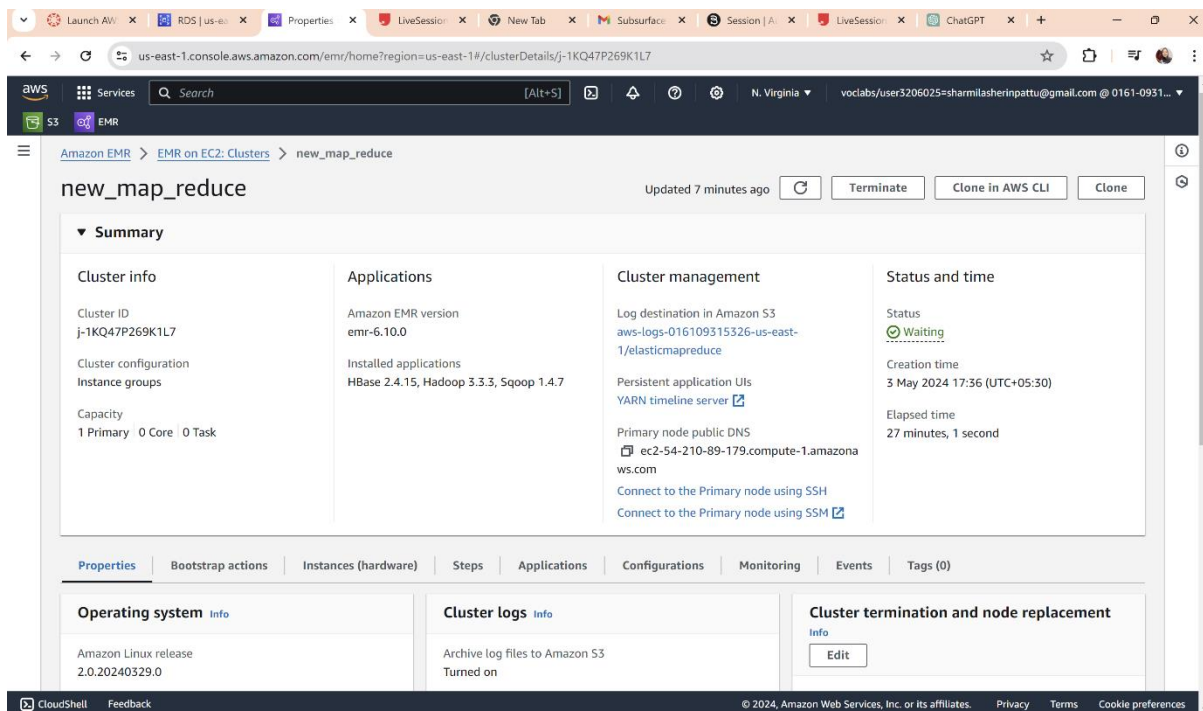
1. Begin by setting up an AWS environment and configuring an RDS instance.
2. Utilize the provided AWS account credentials.
3. Upload only two files, yellow\_tripdata\_2017-01.csv and yellow\_tripdata\_2017-02.csv, from the dataset due to its size.
4. Ensured a suitable schema is created for the datasets to facilitate their upload to the RDS instance.

## 1.RDS instance creation in AWS



## 2. EMR creation

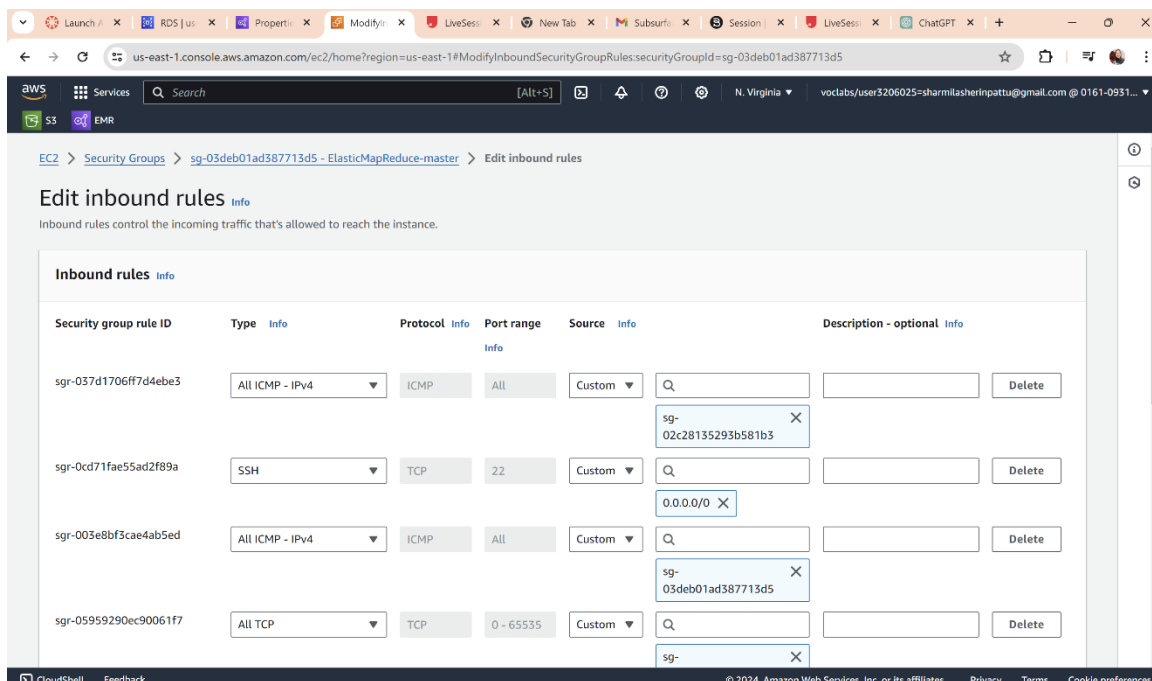
- Including Apache Sqoop, Apache Hbase, Hadoop



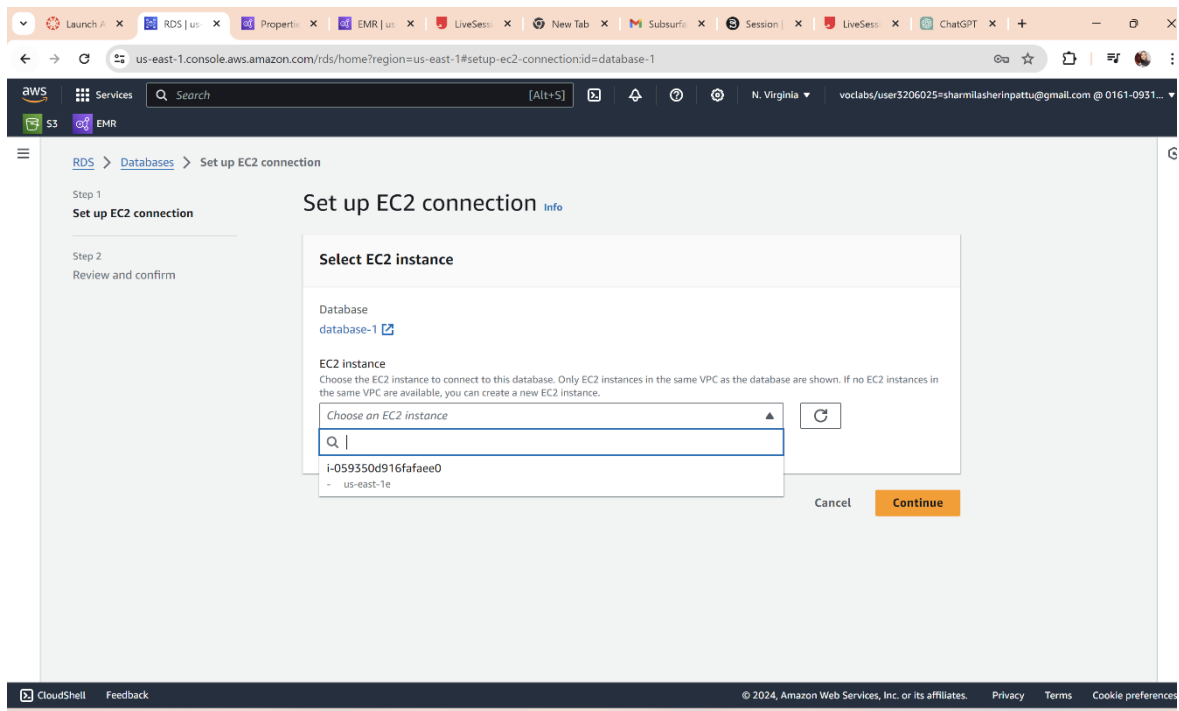
## 3: Connecting the RDS instance with the EMR instance

- To connect the RDS instance with the EMR instance, I adjusted the security group settings.
- Here's what I did: Accessed the AWS Management Console. Navigated to the EC2 section and selected "Security Groups".
- Identified the security group associated with the RDS instance.
- Edited the inbound rules to permit traffic from the EMR instance.

- This entailed specifying either the EMR instance's security group ID or its IP address range, along with the relevant port for database connectivity (such as 3306 for MySQL). Saved the modifications to the security group.
- Through these adjustments, I ensured secure connectivity between the EMR and RDS instances, facilitating data processing and analysis tasks.



- Then we click on 'Action' button on RDS menu and then 'Set up EC2 connection'.



- To access the RDS instance through the EMR instance, we used the following command: **“mysql -h database-1.cpmm6kice78i.us-east-1.rds.amazonaws.com -P 3306 -u admin -p “**
- Upon executing the command, we were prompted to enter the password. After providing the password, the login process was completed successfully.



```

MySQL [(none)]> USE yellow_taxi;
Database changed
MySQL [yellow_taxi]> CREATE TABLE trips (
-> VendorID VARCHAR(255),
-> tpep_pickup_datetime TIMESTAMP NOT NULL DEFAULT '0000-00-00 00:00:00',
-> tpep_dropoff_datetime TIMESTAMP NOT NULL DEFAULT '0000-00-00 00:00:00',
-> passenger_count INT,
-> trip_distance DOUBLE,
-> RatecodeID VARCHAR(255),
-> store_and_fwd_flag VARCHAR(255),
-> PULocationID VARCHAR(255),
-> DOLocationID VARCHAR(255),
-> payment_type VARCHAR(255),
-> 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 [yellow_taxi]>

```

To download the necessary CSV files, I executed the following 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"

```

These commands fetched the specified CSV files from the provided URLs.

Just to showcase that table data was 0 before importing data to database:

```

Database changed
MySQL [yellow_taxi]> CREATE TABLE trips (
  -- VendorID VARCHAR(255),
  -- tpep_pickup_datetime TIMESTAMP NOT NULL DEFAULT '0000-00-00 00:00:00',
  -- tpep_dropoff_datetime TIMESTAMP NOT NULL DEFAULT '0000-00-00 00:00:00',
  -- passenger_count INT,
  -- trip_distance DOUBLE,
  -- RatecodeID VARCHAR(255),
  -- store_and_fwd_flag VARCHAR(255),
  -- pickup_location VARCHAR(255),
  -- dropoff_location VARCHAR(255),
  -- payment_type VARCHAR(255),
  -- 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 [yellow_taxi]> exit
bye
[hadoop@ip-172-31-58-114 ~]$ wget "https://nyc-tlc-upgrad.s3.amazonaws.com/yellow_tripdata_2017-01.csv"
--2024-05-06 00:54:21-- https://nyc-tlc-upgrad.s3.amazonaws.com/yellow_tripdata_2017-01.csv
Resolving nyc-tlc-upgrad.s3.amazonaws.com (nyc-tlc-upgrad.s3.amazonaws.com)... 52.217.161.25, 52.217.224.209, 3.5.25.102, ...
Connecting to nyc-tlc-upgrad.s3.amazonaws.com (nyc-tlc-upgrad.s3.amazonaws.com) [52.217.161.25]:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 914029540 (872M) [text/csv]
Saving to: 'yellow_tripdata_2017-01.csv'

100%[=====>] 914,029,540 37.0MB/s in 27s

2024-05-06 00:54:49 (31.8 MB/s) = 'yellow_tripdata_2017-01.csv' saved [914029540/914029540]

[hadoop@ip-172-31-58-114 ~]$ wget "https://nyc-tlc-upgrad.s3.amazonaws.com/yellow_tripdata_2017-02.csv"
--2024-05-06 00:55:07-- https://nyc-tlc-upgrad.s3.amazonaws.com/yellow_tripdata_2017-02.csv
Resolving nyc-tlc-upgrad.s3.amazonaws.com (nyc-tlc-upgrad.s3.amazonaws.com)... 3.5.10.233, 3.5.25.62, 16.182.106.97, ...
Connecting to nyc-tlc-upgrad.s3.amazonaws.com (nyc-tlc-upgrad.s3.amazonaws.com) [3.5.10.233]:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 863487050 (823M) [text/csv]
Saving to: 'yellow_tripdata_2017-02.csv'

100%[=====>] 863,407,050 29.1MB/s in 26s

2024-05-06 00:55:33 (31.4 MB/s) = 'yellow_tripdata_2017-02.csv' saved [863487050/863487050]

[hadoop@ip-172-31-58-114 ~]$ pwd
/home/hadoop
[hadoop@ip-172-31-58-114 ~]$ ls
yellow_tripdata_2017-01.csv yellow_tripdata_2017-02.csv
[hadoop@ip-172-31-58-114 ~]$

```

**To load data into the MySQL table, I logged in and executed the following SQL commands:**

**LOAD DATA LOCAL INFILE '/home/hadoop/yellow\_tripdata\_2017-01.csv'  
 INTO TABLE trips  
 FIELDS TERMINATED BY ','  
 LINES TERMINATED BY '\n'  
 IGNORE 1 LINES;**

**LOAD DATA LOCAL INFILE '/home/hadoop/yellow\_tripdata\_2017-02.csv'  
 INTO TABLE trips  
 FIELDS TERMINATED BY ','  
 LINES TERMINATED BY '\n'  
 IGNORE 1 LINES;**

These commands imported the data from the specified CSV files into the MySQL table “trips”

```
hadoop@ip-172-31-50-15:~$ mysql -h database-1.cpm6k1ce781.us-east-1.rds.amazonaws.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)]> LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-01.csv'
-> INTO TABLE trips
-> FIELDS TERMINATED BY ','
-> LINES TERMINATED BY '\n'
-> IGNORE 1 LINES;
ERROR 1046 (3D000): No database selected
MySQL [(none)]> use yellow_taxi;
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 [yellow_taxi]> LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-01.csv'
-> INTO TABLE trips
-> FIELDS TERMINATED BY ','
-> LINES TERMINATED BY '\n'
-> IGNORE 1 LINES;

Query OK, 9710820 rows affected, 65535 warnings (2 min 23.75 sec)
Records: 9710820 Deleted: 0 Skipped: 0 Warnings: 19421640

MySQL [yellow_taxi]>
MySQL [yellow_taxi]>
MySQL [yellow_taxi]> LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-02.csv'
-> INTO TABLE trips
-> FIELDS TERMINATED BY ','
-> LINES TERMINATED BY '\n'
-> IGNORE 1 LINES;

Query OK, 9169775 rows affected, 65535 warnings (2 min 14.43 sec)
Records: 9169775 Deleted: 0 Skipped: 0 Warnings: 18339550

MySQL [yellow_taxi]>
```

Confirming that data is loaded: to do this, we run simple SQL queries:

- > select count (\*) from trips;
- > select \* from trips limit 5;

```
MySQL [yellow_taxi]>
MySQL [yellow_taxi]>
MySQL [yellow_taxi]> LOAD DATA LOCAL INFILE '/home/hadoop/yellow_tripdata_2017-02.csv'
-> INTO TABLE trips
-> FIELDS TERMINATED BY ','
-> LINES TERMINATED BY '\n'
-> IGNORE 1 LINES;

Query OK, 9169775 rows affected, 65535 warnings (2 min 14.43 sec)
Records: 9169775 Deleted: 0 Skipped: 0 Warnings: 18339550

MySQL [yellow_taxi]> select count (*) from trips;
ERROR 1064 (42000): You have an error in your SQL syntax; check the manual that corresponds to your MySQL server version for the right syntax to use near '*) from trips' at line 1
MySQL [yellow_taxi]> SELECT COUNT(*) FROM trips;

+-----+
| COUNT(*) |
+-----+
| 18880595 |
+-----+
1 row in set (52.88 sec)

MySQL [yellow_taxi]> select * from trips limit 5;

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| VendorID | tpep_pickup_datetime | tpep_dropoff_datetime | passenger_count | trip_distance | RatecodeID | store_and_fwd_flag | PULocationID | DOLocationID | payment_type | fare_amount | extra | mta_tax | tip_amount | tolls_amount | improvement_surcharge | total_amount | congestion_surcharge | airport_fee |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 1 | 2017-01-01 00:32:05 | 2017-01-01 00:37:48 | 1 | 1.2 | 1 | N | 140 | 236 | 2 | 6.5 | 0.5 | 0.5 | 0 | 0 | 0 | 7.8 | 0 | 0 |
| 1 | 2017-01-01 00:43:25 | 2017-01-01 00:47:42 | 2 | 0.7 | 1 | N | 140 | 140 | 2 | 5 | 0.5 | 0.5 | 0 | 0 | 0 | 6.3 | 0 | 0 |
| 1 | 2017-01-01 00:49:10 | 2017-01-01 00:53:53 | 2 | 0.8 | 1 | N | 140 | 237 | 2 | 5.5 | 0.5 | 0.5 | 0 | 0 | 0 | 6.8 | 0 | 0 |
| 1 | 2017-01-01 00:36:42 | 2017-01-01 00:41:09 | 1 | 1.1 | 1 | N | 41 | 42 | 2 | 6 | 0.5 | 0.5 | 0 | 0 | 0 | 7.3 | 0 | 0 |
| 1 | 2017-01-01 00:07:41 | 2017-01-01 00:18:16 | 1 | 3 | 1 | N | 263 | 12 | 2 | 11 | 0.5 | 0.5 | 0 | 0 | 0 | 12.3 | 0 | 0 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
5 rows in set (0.01 sec)
```

After Importing the final values in dataset is around : 18880595