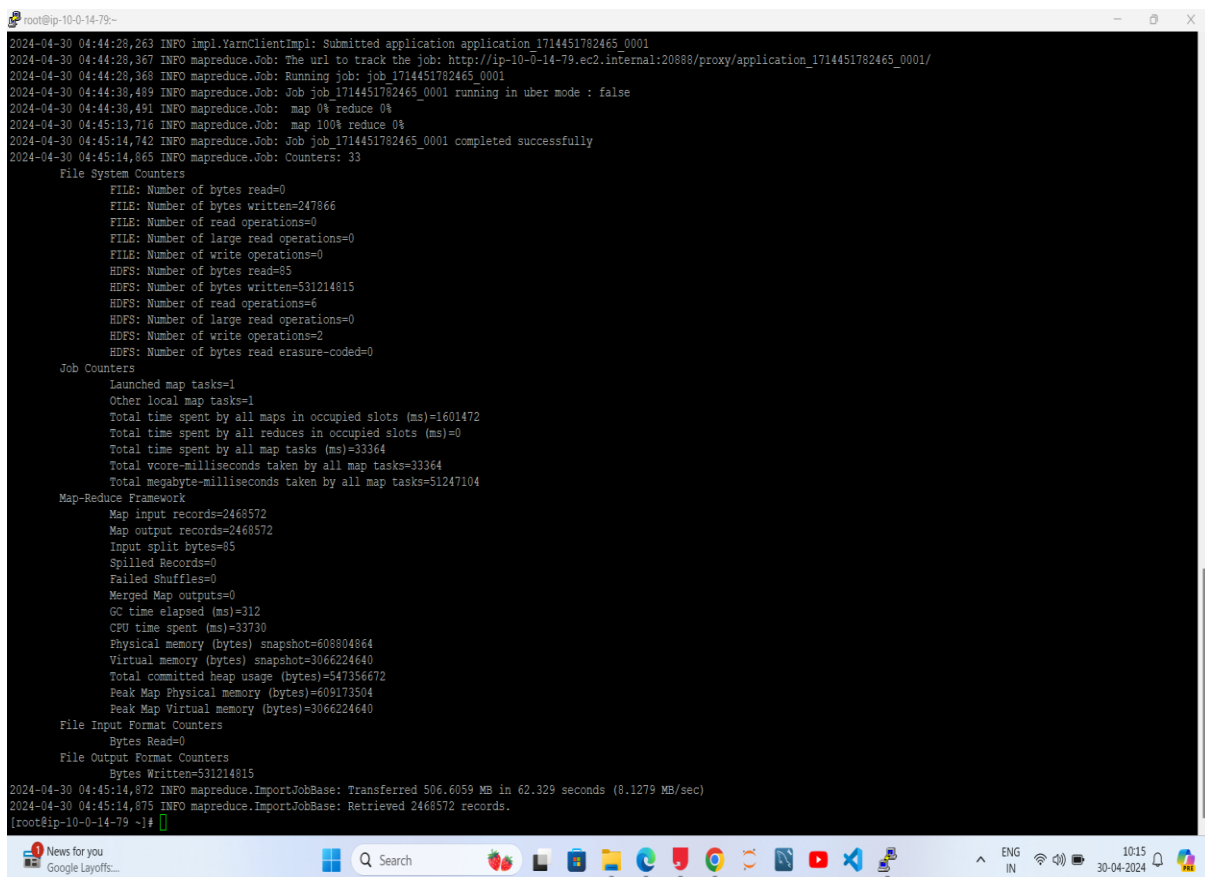# Data Ingestion from the RDS to HDFS using Sqoop

**Sqoop Import command used for importing table from RDS to HDFS**:

```
sqoop import \
--connect jdbc:mysql://upgraddetest.cyaielc9bmnf.us-east-1.rds.amazonaws.com/
testdatabase  \
--table SRC_ATM_TRANS \
--username student --password STUDENT123 \
--target-dir /user/root/spar_nord_bank_atm \
-m 1
```
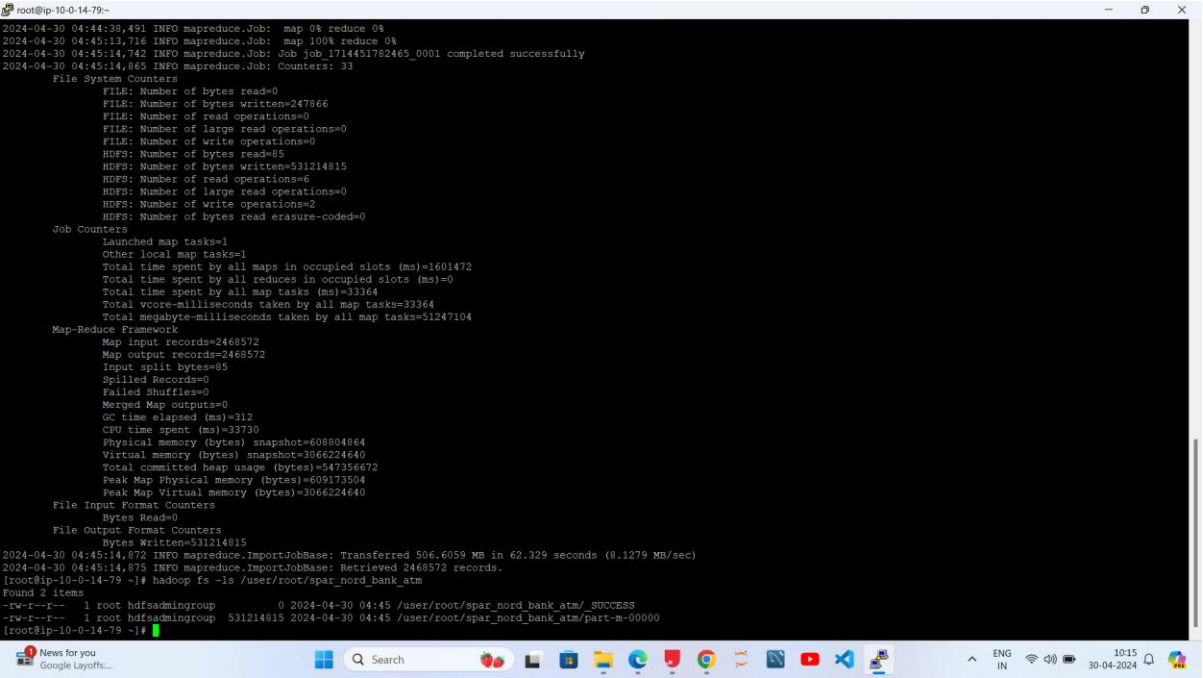


In the screenshot above we can see 2468572 rows have been retrieved.

# Command used to see the list of imported data in HDFS:

hadoop fs -ls /user/root/spar_nord_bank_atm



In the screenshot above we can see two items: - The first file is the success file, indicating that the MapReduce job was successful. - The second file 'part-m-00000' is the one that I imported. Since I used only one mapper in my import command, thus the data is in a single file.

Screenshot of the imported data:

 A portion of data read from part-m-00000 file