# 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 --query 'select * from testdatabase.SRC_ATM_TRANS where $CONDITIONS' --username "student" --password "STUDENT123" --target-dir "/user/root/bank_atm" --m 1
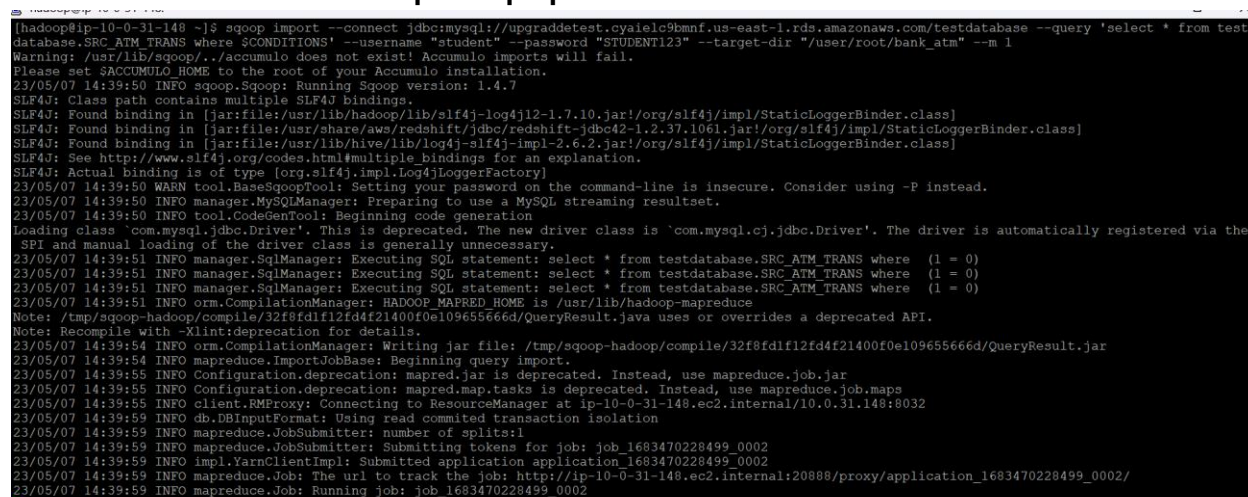

**Command used to see the list of imported data in HDFS:**
hadoop fs -mv /user/root/bank_atm/part-m-00000 /user/livy

hadoop fs -ls /user/livy/part-m-00000

**Screenshot of the imported data:**

1. **Screenshot of data import sqoop command**

2. Screenshot of the data imported



```
        File System Counters
                FILE: Number of bytes read=0
                FILE: Number of bytes written=189185
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=87
                HDFS: Number of bytes written=531214815
                HDFS: Number of read operations=4
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
        Job Counters
                Launched map tasks=1
                Other local map tasks=1
                Total time spent by all maps in occupied slots (ms)=1575984
                Total time spent by all reduces in occupied slots (ms)=0
                Total time spent by all map tasks (ms)=32833
                Total vcore-milliseconds taken by all map tasks=32833
                Total megabyte-milliseconds taken by all map tasks=50431488
        Map-Reduce Framework
                Map input records=2468572
                Map output records=2468572
                Input split bytes=87
                Spilled Records=0
                Failed Shuffles=0
                Merged Map outputs=0
                GC time elapsed (ms)=342
                CPU time spent (ms)=29670
                Physical memory (bytes) snapshot=603389952
                Virtual memory (bytes) snapshot=3289088000
                Total committed heap usage (bytes)=538443776
        File Input Format Counters
                Bytes Read=0
        File Output Format Counters
                Bytes Written=531214815
23/05/07 14:40:46 INFO mapreduce.ImportJobBase: Transferred 506.6059 MB in 50.607 seconds (10.0106 MB/sec)
23/05/07 14:40:46 INFO mapreduce.ImportJobBase: Retrieved 2468572 records.
[hadoop@ip-10-0-31-148 ~]$
```

3. Screen shot of the data moved to livy location



```
                Bytes Written=531214815
23/05/07 14:40:46 INFO mapreduce.ImportJobBase: Transferred 506.6059 MB in 50.607 seconds (10.0106 MB/sec)
23/05/07 14:40:46 INFO mapreduce.ImportJobBase: Retrieved 2468572 records.
[hadoop@ip-10-0-31-148 ~]$ hadoop fs -mv /user/root/bank_atm/part-m-00000 /user/livy
[hadoop@ip-10-0-31-148 ~]$ ^C
[hadoop@ip-10-0-31-148 ~]$ hadoop fs -ls /user/livy/part-m-00000
-rw-r--r--   1 hadoop hadoop  531214815 2023-05-07 14:40 /user/livy/part-m-00000
[hadoop@ip-10-0-31-148 ~]$ ^C
```