# 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 --username student --password STUDENT123 --table SRC_ATM_TRANS --target-dir /user/root/etlassignment -m 1

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

hdfs dfs -ls

**Screenshot of the imported data:**

```
[root@ip-10-0-0-196 ~]# hdfs dfs -ls
Found 2 items
drwx------   - root supergroup          0 2021-07-03 11:30 .staging
drwxr-xr-x   - root supergroup          0 2021-07-03 11:30 etlassignment
[root@ip-10-0-0-196 ~]#
```