# Big Data Video Presentation

**Question d:**

**Uploading CSV in S3**



**Spark Initialization**



**Spark Code**

```
from pyspark.sql import SparkSession

# Create a Spark session

spark = SparkSession.builder.appName("DelayFlightsAnalysis").getOrCreate()


# Upload CSV and create DataFrame

csv_path = "path/to/DelayedFlights-updated.csv"

delay_flights_df = spark.read.csv(csv_path, header=True, inferSchema=True)


# Create a temporary table named delay_flights

delay_flights_df.createOrReplaceTempView("delay_flights")
```

```
# Run the first query using HiveQL

result_hiveql = spark.sql("""

    SELECT Year, AVG((CarrierDelay / ArrDelay) * 100) AS avg_delay_percentage

    FROM delay_flights

    GROUP BY Year""").show()


# Run the second query using Spark SQL

result_sparksql = spark.sql("""

    SELECT Year, AVG((CarrierDelay / ArrDelay) * 100) AS avg_delay_percentage

    FROM delay_flights

    GROUP BY Year""").show()
```

## Question e:

### EMR Cluster Initialization



### Adding Shell Script and table in Hadoop directory



### Grant permission to Shell Script and execute

```
[hadoop@ip-172-31-55-229 ~]$ chmod +x assignment.sh
[hadoop@ip-172-31-55-229 ~]$ ./assignment.sh
Hive Session ID = 6a8f9b0d-62d1-4063-b0fe-e8df4f23f963
```

## Map Reducer Worker distribution

```
Map 1: -/-      Reducer 2: 0/1
Map 1: 0/1      Reducer 2: 0/1
Map 1: 0/1      Reducer 2: 0/1
Map 1: 0(+1)/1  Reducer 2: 0/1
Map 1: 0/1      Reducer 2: 0/1
Map 1: 1/1      Reducer 2: 0(+1)/1
Map 1: 1/1      Reducer 2: 1/1
OK
```
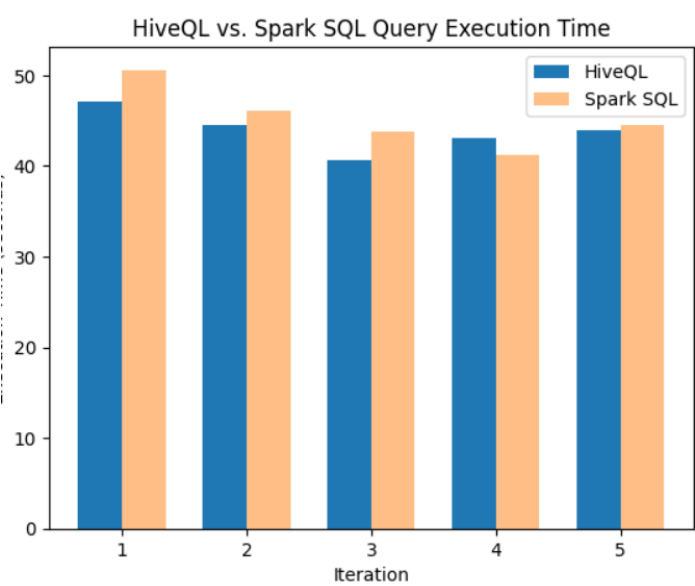
## Run Time Calculation

```
Time taken: 23.629 seconds, Fetched: 1 row(s)
Iteration 1
HiveQL Execution Time: 0m43.761s seconds
Spark SQL Execution Time: 0m44.568s seconds
Spark SQL Execution Time: 0m44.568s seconds
Iteration 2
HiveQL Execution Time: 0m45.451s seconds
Spark SQL Execution Time: 0m44.386s seconds
Iteration 3
HiveQL Execution Time: 0m41.443s seconds
Spark SQL Execution Time: 0m44.888s seconds
Iteration 4
HiveQL Execution Time: 0m44.945s seconds
Spark SQL Execution Time: 0m44.805s seconds
Iteration 5
HiveQL Execution Time: 0m39.552s seconds
Spark SQL Execution Time: 0m49.065s seconds
```

## Download into local path

```
C:\Users\tuanm\OneDrive\Desktop\DS\bigdata\new>scp -i Bigdata.pem hadoop@ec2-52-91-134-150.compute-1.amazonaws.com:execution_times.csv .
execution_times.csv
```

## Results

| | Iteration | HiveQL_Time | SparkSQL_Time |
|---|---|---|---|
| 0 | 1 | 47.165 | 50.591 |
| 1 | 2 | 44.515 | 46.060 |
| 2 | 3 | 40.606 | 43.760 |
| 3 | 4 | 43.132 | 41.166 |
| 4 | 5 | 43.893 | 44.501 |



## Shell Script file is **assignment_for_5_iteration.sh**

**Question f:**

Same command approach as above. But Shell script is different. I have attached it in the same folder
assignment_for_5_queries.sh

**Results**

| | Mapping | HiveQL_Time | SparkSQL_Time |
|---|---|---|---|
| 0 | Career delay query | 58.340 | 43.824 |
| 1 | Nas delay query | 42.040 | 41.719 |
| 2 | Weather delay query | 43.682 | 44.415 |
| 3 | Late aircraft delay query | 42.455 | 43.383 |
| 4 | Security delay query | 34.011 | 44.908 |

**Question g:**

**Results**