# AWS PROJECT

———

*Phase - 1*

# AWS PROJECT PHASE 1

- *Omkar Vartak –Year Selected 1987*

Year Selected 1987

# CREATING THE CLUSTER

# CREATING THE CLUSTER

# USING WINSCP TO TRANSFER FILES

# MODIFY HDFS



```
portmap                run a portmap service
secondarynamenode      run the DFS secondary namenode
sps                    run external storagepolicysatisfier
zkfc                   run the ZK Failover Controller daemon

SUBCOMMAND may print help when invoked w/o parameters or with -h.
[hadoop@ip-172-    ~]$ hdfs dfs -mkdir -p /user/hive/warehouse
[hadoop@ip-172-    0 ~]$ hdfs dfs -chmod g+w /user/hive/warehouse
```

After this we will extract the downloaded bz2 file using any extractor.

# HDFS COMMANDS

Below commands were used:
1. hdfs dfs -mkdir -p /user/hive/warehouse
2. hdfs dfs -chmod g+w /user/hive/warehouse

# CREATING DATABASE AND TABLES, LOADING THE DATA.

# DISPLAYING THE DATA USING THE LIMIT FUNCTION:-

```
1987    10    22    4       728     730     852     849     PS      1451    N
A       84    79    NULL    3       -2      SAN     SFO     447     NULL    N
ULL     0     NA    0       NULL    NULL    NULL    NULL    NULL
1987    10    23    5       731     730     902     849     PS      1451    N
A       91    79    NULL    13      1       SAN     SFO     447     NULL    N
ULL     0     NA    0       NULL    NULL    NULL    NULL    NULL
1987    10    24    6       744     730     908     849     PS      1451    N
A       84    79    NULL    19      14      SAN     SFO     447     NULL    N
ULL     0     NA    0       NULL    NULL    NULL    NULL    NULL
1987    10    25    7       729     730     851     849     PS      1451    N
A       82    79    NULL    2       -1      SAN     SFO     447     NULL    N
ULL     0     NA    0       NULL    NULL    NULL    NULL    NULL
1987    10    26    1       735     730     904     849     PS      1451    N
A       89    79    NULL    15      5       SAN     SFO     447     NULL    N
ULL     0     NA    0       NULL    NULL    NULL    NULL    NULL
1987    10    28    3       741     725     919     855     PS      1451    N
A       98    90    NULL    24      16      SAN     SFO     447     NULL    N
ULL     0     NA    0       NULL    NULL    NULL    NULL    NULL
1987    10    29    4       742     725     906     855     PS      1451    N
A       84    90    NULL    11      17      SAN     SFO     447     NULL    N
ULL     0     NA    0       NULL    NULL    NULL    NULL    NULL
1987    10    31    6       726     725     848     855     PS      1451    N
A       82    90    NULL    -7      1       SAN     SFO     447     NULL    N
ULL     0     NA    0       NULL    NULL    NULL    NULL    NULL
```

# QUERY DETERMINE THE THREE CARRIERS WITH THE HIGHEST DELAY TIME (IN HOURS)

```
hive>  with total as (
    > SELECT Year, UniqueCarrier , (round ((sum (ArrDelay)) /60, 2)) as sum_arrdelay, (round ((sum (DepDelay)) /60,2)) as sum_depdelay
    >  from Omkar1987
    >  group by Year, UniqueCarrier
    > )
    > select Year, UniqueCarrier, sum_arrdelay, sum_depdelay,(sum_arrdelay+sum_depdelay) as Total_delay from total
    >  order by Total_delay desc
    >  limit 6 ;
Query ID = hadoop_20221107060         31cc-46f2-a229-917f58d3b703
Total jobs = 1
Launching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1667799766282_0004)

----------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      8          8        0        0       0       0
Reducer 2 ...... container     SUCCEEDED      2          2        0        0       0       0
Reducer 3 ...... container     SUCCEEDED      1          1        0        0       0       0
----------------------------------------------------------------------------------------------
VERTICES: 03/03  [==========================>>] 100%  ELAPSED TIME: 12.80 s
----------------------------------------------------------------------------------------------
OK
1987    DL      36494.58        22571.73        59066.31
1987    CO      22340.97        25345.23        47686.2
1987    UA      21772.87        20607.1 42379.97
1987    PI      17991.15        14837.93        32829.08
1987    EA      16098.73        15651.95        31750.68
1987    NW      18948.92        12192.28        31141.199999999997
Time taken: 19.242 seconds, Fetched: 6 row(s)
hive>
```

# QUERY DETERMINE OVERALL WHICH TYPE OF DELAY (ARRIVALS OR DEPARTURES) IS THE LARGEST FOR AIRPORTS

```
hive> Select Sum(ArrDelay) AS  totArrival_delay From Omkar1987;
Query ID = hadoop_20221108051    _    -aed1-4476-b    -1ce14a37e0ca
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667883594833_0004)


----------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      6          6        0        0       0       0
Reducer 2 ...... container     SUCCEEDED      1          1        0        0       0       0
----------------------------------------------------------------------------
VERTICES: 02/02  [==========================>>] 100%  ELAPSED TIME: 9.90 s
----------------------------------------------------------------------------
OK
12170428
Time taken: 10.784 seconds, Fetched: 1 row(s)
hive> Select Sum(DepDelay) AS  totDep_delay From Omkar1987;
Query ID = hadoop_20221108051533_        d6-a7e5-d4fe4ebc7
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1667883594833_0004)


----------------------------------------------------------------------------
        VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED      6          6        0        0       0       0
Reducer 2 ...... container     SUCCEEDED      1          1        0        0       0       0
----------------------------------------------------------------------------
VERTICES: 02/02  [==========================>>] 100%  ELAPSED TIME: 10.57 s
----------------------------------------------------------------------------
OK
10419357
Time taken: 11.356 seconds, Fetched: 1 row(s)
hive>
```

# QUERY DETERMINES THE THREE **AIRPORTS** WITH THE HIGHEST DELAY TIME (IN HOURS)

- We will divide the result by 60, In order to get the results in hours we will then use the round function.

- After that we will get the result in hours which can be used in the visualization process.

- Arr_delay = 12170428/60=202840.46

- Dep_delay=10419357/60=173655.95

- In conclusion Arrival delay is more than that of Departure delay

# QUERY DETERMINES THE THREE **AIRPORTS** WITH THE HIGHEST DELAY TIME (IN HOURS)

```
hive> with total as (
    > select year, origin, dest, (round(sum(arrdelay)/60,2)) as sum_arrdelay, (round(sum(depdelay)/60,2)) as sum_depdelay
    > from Omkar1987
    > group by year, origin, dest
    > )
    > select year, origin, dest, sum_arrdelay, sum_depdelay, (sum_arrdelay+sum_depdelay) as total_delay from total
    > order by total_delay desc
    > limit 3;
Query ID = hadoop_2                          e0ab-4f49-9ea6-3e4aee303e25
Total jobs = 1
Launching Job 1 out of 1
Tez session was closed. Reopening...
Session re-established.
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1667937412624_0006)

--------------------------------------------------------------------------------
        VERTICES      MODE        STATUS    TOTAL   COMPLETED   RUNNING   PENDING   FAILED   KILLED
--------------------------------------------------------------------------------
Map 1 .......... container      SUCCEEDED      8         8         0         0         0         0
Reducer 2 ...... container      SUCCEEDED      2         2         0         0         0         0
Reducer 3 ...... container      SUCCEEDED      1         1         0         0         0         0
--------------------------------------------------------------------------------
VERTICES: 03/03  [==========================>>] 100%  ELAPSED TIME: 11.75 s
--------------------------------------------------------------------------------
OK
1987    LAX     SFO     2152.13 1502.68 3654.8100000000004
1987    SFO     LAX     1488.78 1225.87 2714.6499999999996
1987    PHX     LAX     1227.18 808.07  2035.25
Time taken: 17.897 seconds, Fetched: 3 row(s)
hive>
```
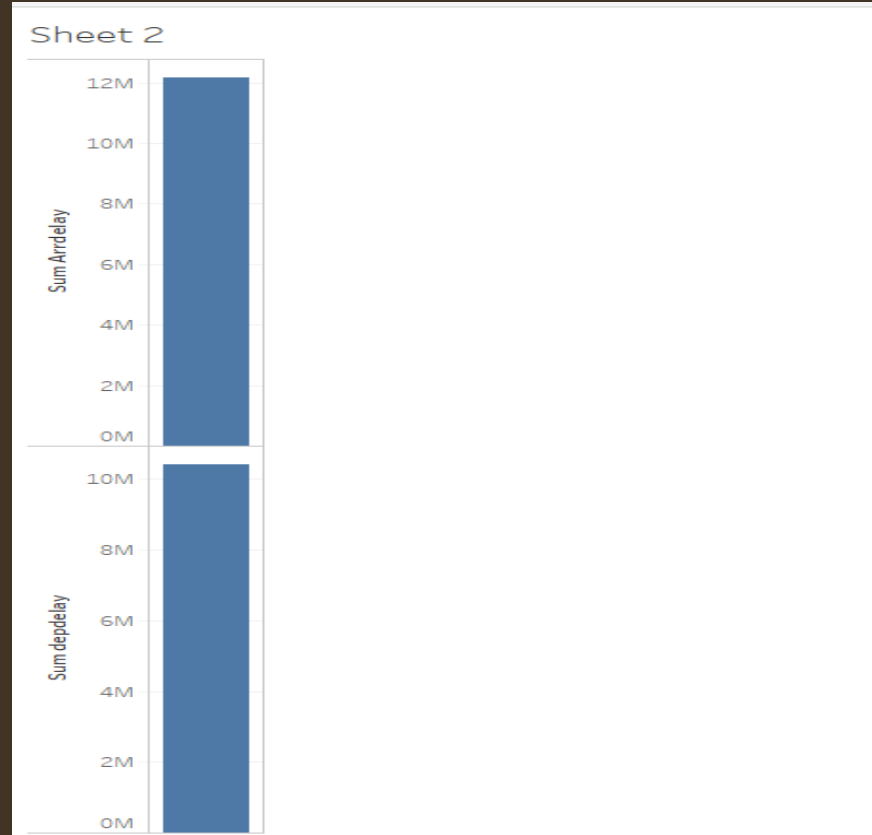
# VISUALIZATION- TOOL USED TABLEAU FOR YEAR 1987. DETERMINE THE THREE CARRIERS WITH THE HIGHEST DELAY TIME (IN HOURS)

# DETERMINE OVERALL WHICH TYPE OF DELAY (ARRIVALS OR DEPARTURES) IS THE LARGEST FOR AIRPORTS

# DETERMINE THE THREE AIRPORTS WITH THE HIGHEST DELAY TIME (IN HOURS