# NJIT

## New Jersey Institute of Technology

# DS 644101 - COURSE PROJECT
## FLIGHT ANALYSIS REPORT

**Instructor:**

YAJUAN LI

**Report by:**

Ankush Ranapure (ar2653@njit.edu)

Dheeraj Ananth Kumar (da559njit.edu)

Jaya Sathwika Gadi (jg899njit.edu)

Nivedita Reddy Addulla (na566njit.edu)

# Table of context

# Table of Figures

# 1. STRUCTURE OF OOZIE WORKFLOW:
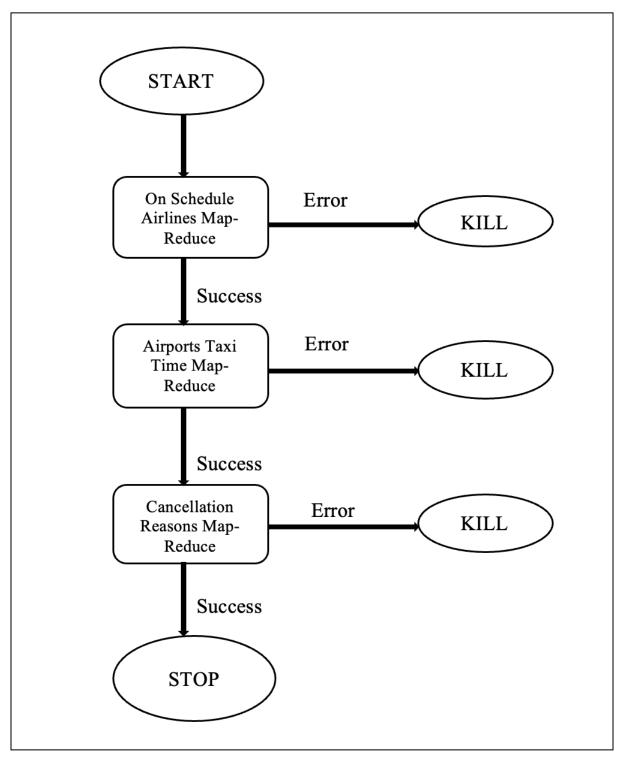


Figure 1.1 - Oozie Workflow

## 2. ALGORITHM:

**First Map-Reduce: On Schedule Airlines**

1.  Mapper Function: <Key,Value> : <UniqueCarrier, 1 or 0>
    For each input line (ignoring the first line and NA data):
    ○  Read the UniqueCarrier and ArrDelay columns.
    ○  If ArrDelay is less than or equal to 10 minutes, output: `<UniqueCarrier, 1>`.
    ○  Otherwise, output: `<UniqueCarrier, 0>`.

2.  Reducer Function: <Key,Value> : <UniqueCarrier, Probability>
    Initialize variables: `sum_1 = 0`, `sum_0_1 = 0`, `sum_0 = 0`.
    For each input from the mapper:
    ○  Sum the values based on the UniqueCarrier key.
    ○  If the value is 1, increment `sum_1`.
    ○  Increment `sum_0_1` for each value (both 1 and 0).
    ○  Calculate the probability for each UniqueCarrier: `Probability = sum_1 / sum_0_1`.
    Sort the probabilities using a comparator function.

3.  Output Handling:
    Output the top 3 and bottom 3 UniqueCarriers with their probabilities.
    If the data is NULL, output: "There is no value used, so no output."

**Second Map-Reduce: On Airport Taxi Time**

1.  Mapper Algorithm Read Data: <key,value>:<IATA airport code, TaxiTime>:<Origin,TaxiOut>or<Dest,Taxi>
    ○  Read the dataset line by line, skipping the first line (header).
    ○  Check Taxi Time: For each line, check if the TaxiIn or TaxiOut value is not 'N/A'.
    ○  Output Key-Value Pair: If the taxi time is valid, output a key-value pair: <IATA airport code, TaxiTime>.

2.  Reducer Algorithm: <key,value>:<IATA airport code, Average TaxiTime>
    ○  Sum Taxi Times: For each unique key (IATA airport code), sum the taxi times from all its associated values.
    ○  Count Occurrences: Count the total number of times each key appears.

○ Calculate Average: For each key, calculate the average taxi time using the formula: average = total taxi time / number of occurrences.
○ Sort Results: Sort the airports by their average taxi times. Output Top and Bottom

3. Output Handling:
Output the three airports with the longest and three with the shortest avg t.
If the data is NULL, then output: There is no value used, so no output.

**Third Map-Reduce: On Cancellation Reasons**

1. Mapper Function: <Key,Value> : <CancellationCode, 1>
   ○ Input: Key-Value pairs (CancellationCode, 1)
   ○ Read the data line by line.
   ○ Ignore the first line.
   ○ If the value of Canceled is 1 and CancellationCode is not "N/A", output a new key-value pair: (CancellationCode, 1)

2. Reducer Function: <key,value>:<CancellationCode, sum of 1s>
   ○ Input: Key-Value pairs (CancellationCode, list of 1s)
   ○ Sum the values (count of 1s) for each CancellationCode.
   ○ Use a comparator function to sort the CancellationCode based on the sum of 1s.
   ○ Output: The most common reason for flight cancellations.

3. Output Handling:
   ○ If the data is NULL, output: "There is no common reason for flight cancellation."

**3. A performance measurement plot that compares the workflow execution time in response to an increasing number of VMs used for processing the entire data set (22 years) and an in-depth discussion on the observed performance comparison results**
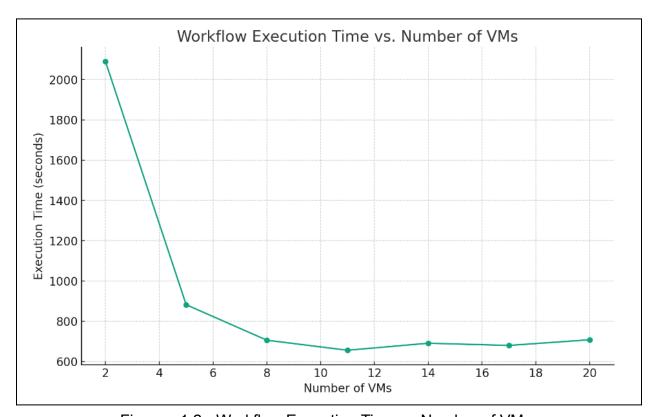


Figure - 1.2 - Workflow Execution Time vs Number of VMs

In Figure 1.2, an increase in the number of virtual machines (VMs) leads to a reduction in workflow execution time. This enhancement is attributed to the heightened processing capacity of the Hadoop cluster, facilitated by the parallel handling of data across multiple datanodes.

Consequently, the execution time for each map-reduce job diminishes, resulting in a shorter overall duration for the Oozie workflow. It's important to note, however, that the execution time for dealing with the same data size may not consistently decrease with a rise in the number of VMs.

Once the execution time reaches a certain threshold, further attempts to augment the number of VMs may not yield additional time savings. This occurs due to the increased information interaction time among datanodes within the Hadoop cluster as the number of VMs escalates.

**4. A performance measurement plot that compares the workflow execution time in response to an increasing data size (from 1 year to 22 years) and an in-depth discussion on the observed performance comparison results**
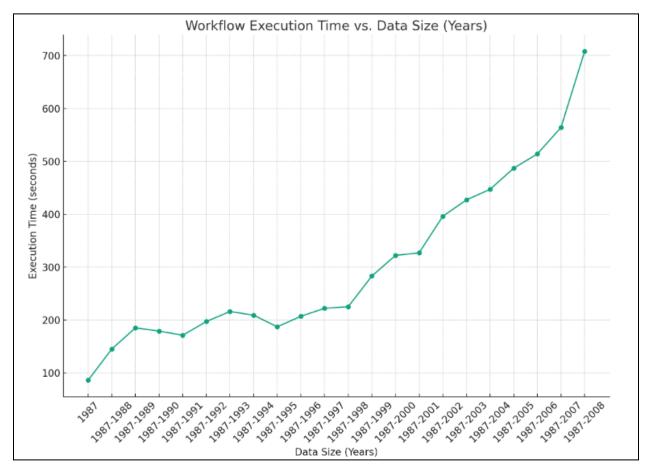


Figure 1.3 - Workflow Execution Time vs Data Size (in Years)

In Figure 1.3, as the data size grows, the execution time of the oozie workflow consistently rises. Initially, the increase in time consumption is gradual due to the relatively moderate growth in data during the early years.

However, after 1998, the time consumption experienced a rapid surge, indicating a steeper slope compared to the earlier years. This increase is due to a faster growth in flight data from 1998 to 2008, reflecting a rising trend in the number of individuals wanting for air travel.
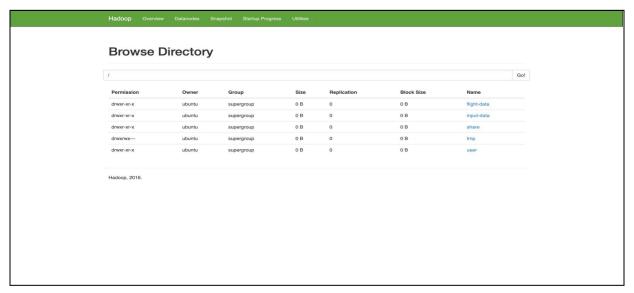
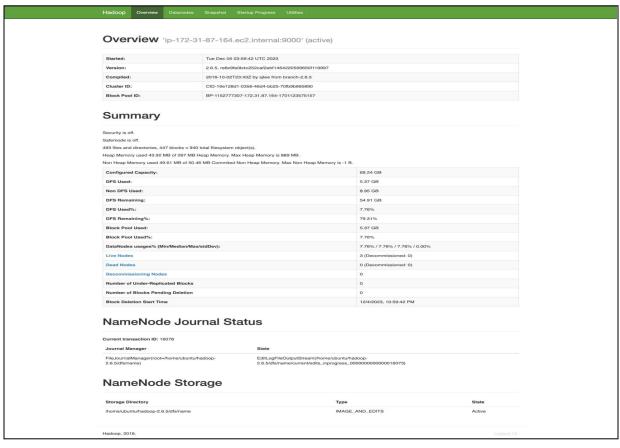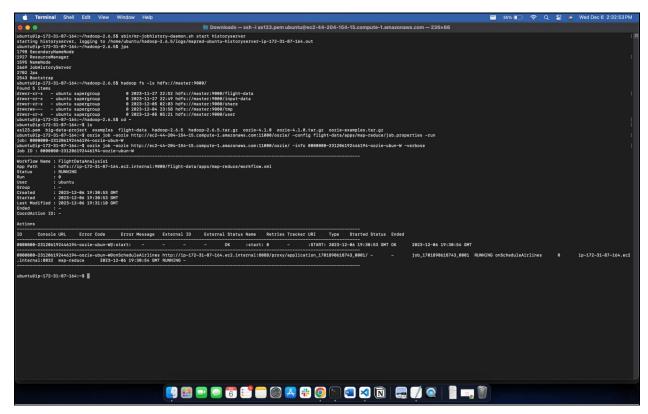## Miscellaneous Screenshots:



Figure 2.1 - Hadoop Directory screenshot



Figure 2.2 - HDFS Overview

Figure 2.3 - Oozie Job running screenshot