# Report: Cluster Setup
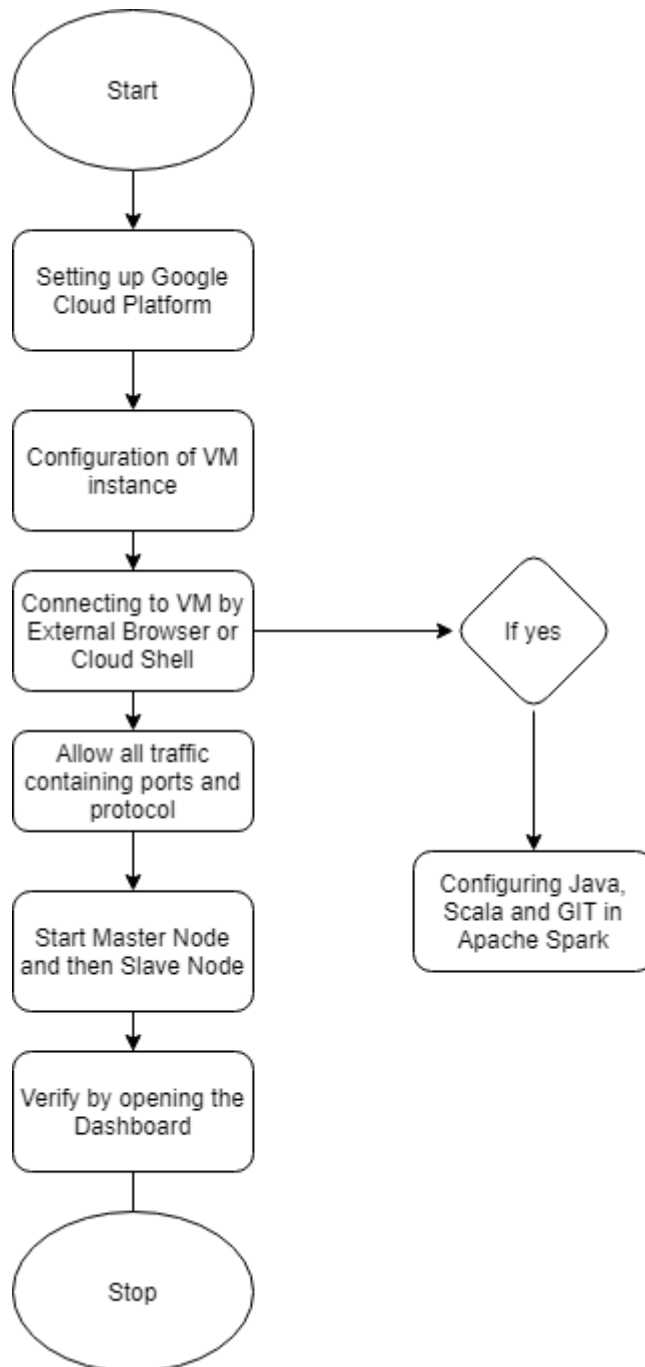
Flowchart:



Source: https:// app.diagram.net

Screenshots:

Configured & Initialized Apache Spark Cluster:





Started Master Node:

Report: Data Extraction

Algorithm:

1. Importing and installing pymongo, dnspython and tweepy.
2. Tweetcounter variable declared and initialised.
3. Loading the live tweets with the help of on_data method.
4. Appending all the incoming tweets inside a list and setting a condition after a certain number of tweets is reached.
5. On_error keeps track and handles the exception of the internal call.
6. In the main method connection to the MongoDB is set.
7. The database and the tables containing the tweets with particular attributes are inserted at the final step inside the DB.

Ouput when run on the VM instance:

Report: Pre-processing Engine

Algorithm:

1. Establishing connection to the MongoDB client
2. Storing RawDB, its collections and processedDB and its collection inside variables
3. Iterating through the collection of tables from the RawDB and passing each document of table in the JSON format through RegEx or Regular expression.
4. The regular expression removes the URLs, special symbols and white spaces.
5. The cleaned JSONs are appended to a new list which is finally inserted in ProcessedDB

Ouput when run on the VM instance:

Report: Reuter News Articles

Algorithm:

1. Connection to be established to the MongoDb by creating ReuterDb.
2. The .sgm files are read and stored in different variables
3. The findall method is used to find the specific tags
   <REUTERS></REUTERS>, <TEXT></TEXT>  and
   <TITLE></TITLE>.
4. The scanned texts between the tags are read appended in the list.
5. The final list is then inserted into the ReuterDb with the cleaned data.

ScreenShot:

Report: Data Visualization using Graph Database

Cypher:

CREATE (canada:country {title: 'Canada', continent: 'North America', climate: 'cold'}) - [:HAS_SEASON] -> (winter:weather {title: 'Winter', type: 'extreme', duration: '5 months'})

MATCH (canada:country {title: 'Canada', continent: 'North America', climate: 'cold'})

CREATE (canada) - [:GETS_FREQUENTLY_HITBY] -> (storm:Calamity {title: 'Storm', stormType: 'harsh', damage: 'extreme'})

RETURN canada, storm

MATCH (winter:weather {title: 'Winter', type: 'extreme', duration: '5 months'})

CREATE(flu:HealthCondition {title: 'Flu', type: 'pandemic', damage: 'medium'}) - [:MORE_CHANCES_TO_AFFECT_IN] -> (winter)

RETURN flu, winter

MATCH (winter:weather {title: 'Winter', type: 'extreme', duration: '5 months'})

CREATE (winter) - [:BREEZE_IS_VERY] -> (cold:Season {title: 'Cold', type: 'extreme', duration: '8 months'})

RETURN winter, cold

MATCH (cold:Season {title: 'Cold', type: 'extreme', duration: '8 months'})

CREATE (cold) - [:SO_RECOMMENDED_TO_STAY] -> (indoor:SafetyMeasure {title: 'Indoor', type: 'Outdoor/Indoor'})

MATCH(flu:HealthCondition {title: 'Flu', type: 'pandemic', damage: 'medium'})

CREATE (flu) - [:TAKE_PRECAUTIONS] -> (safety:Precaution {title: 'Safety', count: 4, priority: 'high'})

MATCH (safety:Precaution {title: 'Safety', count: 4, priority: 'high'})

MATCH (indoor:SafetyMeasure {title: 'Indoor', type: 'Outdoor/Indoor'})

CREATE (indoor) - [:A_PRECAUTION_FOR] -> (safety)


MATCH (storm:Calamity {title: 'Storm', stormType: 'harsh', damage: 'extreme'})

MATCH (safety:Precaution {title: 'Safety', count: 4, priority: 'high'})

CREATE (storm) - [:HIGHEST_PRIORITY_TO_TAKE] -> (safety)

MATCH (n) RETURN n


MATCH (storm:Calamity {title: 'Storm', stormType: 'harsh', damage: 'extreme'})

CREATE (storm) - [:LEADS_TO_HEAVY] -> (rain:SeasonAffect {title: 'Rain', type:'heavy', size:'508mm'})


MATCH (cold:Season {title: 'Cold', type: 'extreme', duration: '8 months'})

CREATE (cold) - [:MAKES_ROAD_TURN_INTO] -> (ice:consequence {title: 'Ice', damage:'medium', geography:'arctic', area: 'fourty-percent'})


MATCH (winter:weather {title: 'Winter', type: 'extreme', duration: '5 months'})

CREATE (winter) - [:VARIOUS_ROADS_ACCUMULATE] ->
(snow:SeasonAffect {title: 'Snow', type:'heavy', size:'400mm'})


MATCH (snow:SeasonAffect {title: 'Snow', type:'heavy', size:'400mm'})

MATCH (cold:Season {title: 'Cold', type: 'extreme', duration: '8 months'})

CREATE (snow) - [:MAKES_ENVIRONMENT] -> (cold)

MATCH (canada:country {title: 'Canada', continent: 'North America', climate: 'cold'})

CREATE (canada) - [:NEVER_EXPERIENCED_BREEZE_OF_TYPE] -> (hot:Season {title: 'Hot', type: 'low', duration: '4 months'})

Graph: