# PRINCIPLES OF BIG DATA MANAGEMENT

## PHASE #1

**TEAM SIZE:3**

**TEAM MEMBERS:** Sai Tejaswi Koppuravuri ([sk6zb@mail.umkc.edu](mailto:sk6zb@mail.umkc.edu))

Pallavi Desai ([pd2qd@mail.umkc.edu](mailto:pd2qd@mail.umkc.edu))
Anusha Palla ([apgmc@mail.umkc.edu](mailto:apgmc@mail.umkc.edu))

## Github Link of the Project:

[https://github.com/SaitejaswiK/Principles-of-BigData-Management](https://github.com/SaitejaswiK/Principles-of-BigData-Management)

## Objective:

- The principle point of this stage is to build up a framework to store, break down, and envision a social network's.
- Tasks:
    1. Collect social network's information (e.g. tweets) in any format preferred JSON.

    2. Store the content substance (e.g. tweet's content) from the information into a document in HDFS.

    3. Run a Word Count program in Apache Spark and Hadoop on the content document and store the yield and log records locally

## Applications/Software's Used:

Twitter Developer Account, Apache Spark, Python Hadoop.

## Collecting tweets from Twitter:

- Firstly, we have made an developer account in Twitter utilizing beneath connect.
[https://apps.twitter.com/](https://apps.twitter.com/)

- Below are the factors that contains the client certifications to get to Twitter API
  - ACCESS_TOKEN = " 779311765163171844-RCUoOhu2R53ugDk3O8xTX50rgi2zj4o"
  - ACCESS_SECRET = " y9Evdnwz1tfI43fIyun18OQOxgt6HQjWh6g3Gb99ExwOI"
  - CONSUMER_KEY = " xMJiyum9ZLKuGeZDPI1uL3qeU"
  - CONSUMER_SECRET = "6df8h8k2O7AwBJgYREWwTfwB1MFXVBuUm4PttByrGiRKDj6bI5"

- We have composed python program that is utilized to bring tweets in JSON design. (tweet_ data.py)

  Link:https://github.com/SaitejaswiK/Principles-of-BigDataManagement/blob/master/Source/Python%20Programs/tweet_data.py



**Fig1: Tweets collection**

- The extricated record in JSON arrange contains all the tweet points of interest, for example, id, created at, text, profile_background_image_url and so forth.
- From JSON tweets record just the content substance is extricated utilizing Python program. The got content points of interest are put away in a record. (twittertextconvert.py)

  Link:https://github.com/SaitejaswiK/Principles-of-BigData-Management/blob/master/Source/Python%20Programs/twittertextconvert.py

```
hadoop@pallavidesai-VirtualBox:~$ python twittertextconvert.py
hadoop@pallavidesai-VirtualBox:~$ python tweet_hash.py
```

**Fig 2: Creating a python file for Hashtags extraction**

## Store the text content (e.g. tweet's text) from the data into a file in HDFS.

- The twitter tweets content substance record is moved from local to HFDS.
- First a folder is made in HDFS and the content document is moved from local to HDFS utilizing underneath order.

  **Make directory in local**: hadoop fs - mkdir pbproject/input

  **Move content record from local to HDFS**: hadoop fs - copyFromLocal FileOutput.txt pbproject/input

  **To list the records under a registry**: hadoop fs - ls pbproject/input

```
hadoop@pallavidesai-VirtualBox:~$
hadoop@pallavidesai-VirtualBox:~$
hadoop@pallavidesai-VirtualBox:~$ hdfs dfs -ls /pbproject
Found 2 items
drwxr-xr-x   - hadoop supergroup          0 2018-09-25 13:00 /pbproject/input
drwxr-xr-x   - hadoop supergroup          0 2018-09-25 13:06 /pbproject/output
hadoop@pallavidesai-VirtualBox:~$ hdfs dfs -ls /pbproject/input
Found 2 items
-rw-r--r--   1 hadoop supergroup     132286 2018-09-24 10:39 /pbproject/input/FileOutput_hash.txt
-rw-r--r--   1 hadoop supergroup   16198591 2018-09-25 13:00 /pbproject/input/tweets_out.json
hadoop@pallavidesai-VirtualBox:~$
```

**Fig 3: HDFS Commands**

```
drwxr-xr-x   2 hadoop hadoop          4096 Sep 11 09:34 Videos
hadoop@pallavidesai-VirtualBox:~$
hadoop@pallavidesai-VirtualBox:~$
hadoop@pallavidesai-VirtualBox:~$
hadoop@pallavidesai-VirtualBox:~$ hdfs dfs -copyFromLocal extractedOutput.txt /p
binput/input/
```

**Fig 4: Creating Extracted Text in Hadoop**

- The directory created and the files moved to HDFS can be viewed as shown below.
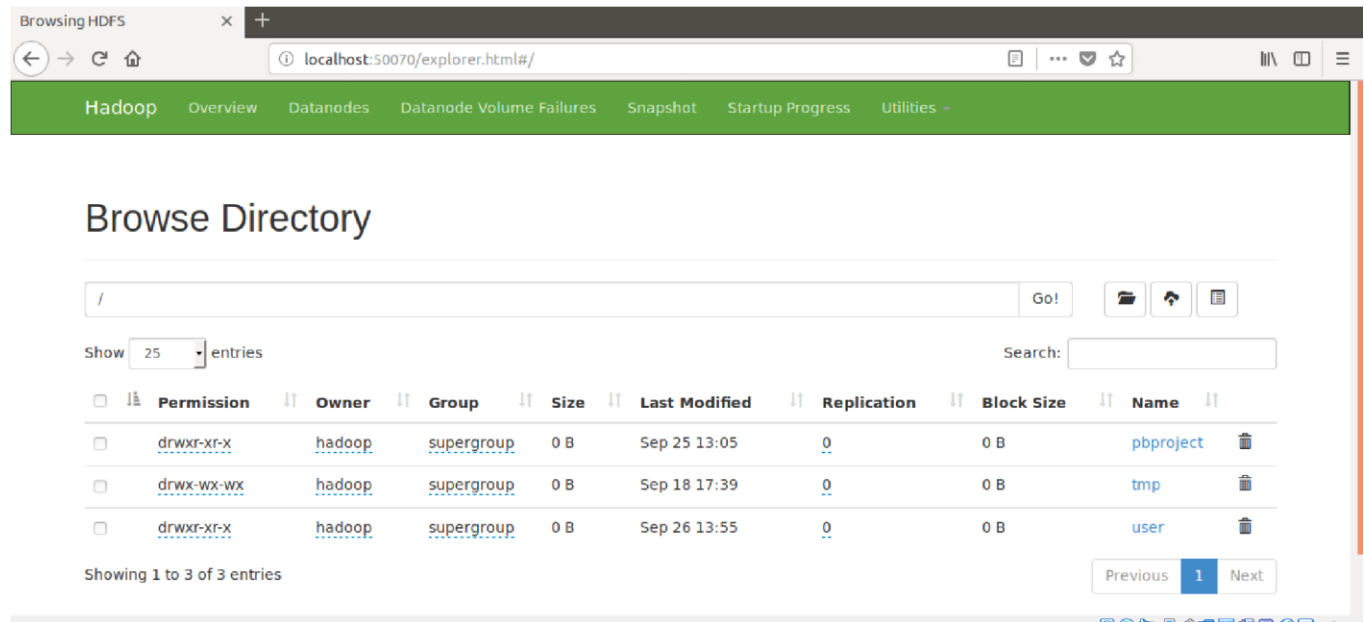
**Fig 5: Directory in HDFS**

Link of the tweets output: https://github.com/SaitejaswiK/Principles-of-BigDataManagement/blob/master/Source/Twitter%20Tweets/FileOutput.txt

**<u>Extracting Hashtags from the obtained output:</u>**

- A python program has been written for the extraction of URL's and hashtags from the obtained output. (tweets_hash.py)

   Link of the code: https://github.com/SaitejaswiK/Principles-of-BigData-Management/blob/master/Source/Python%20Programs/tweet_hash.py

Link of the extracted Hashtags output: https://github.com/SaitejaswiK/Principles-of-BigDataManagement/blob/master/Source/Twitter%20Tweets/FileOutput_hash.txt

## Run a Word Count program in Apache Hadoop on the text file and store the output and log files locally.

- First of all, to run word count program on set of data we require.
- Using Hadoop, run the word count example for the obtained tweets file.



**Fig 6: Running Wordcount in Hadoop**

**Output in Hadoop Browser**

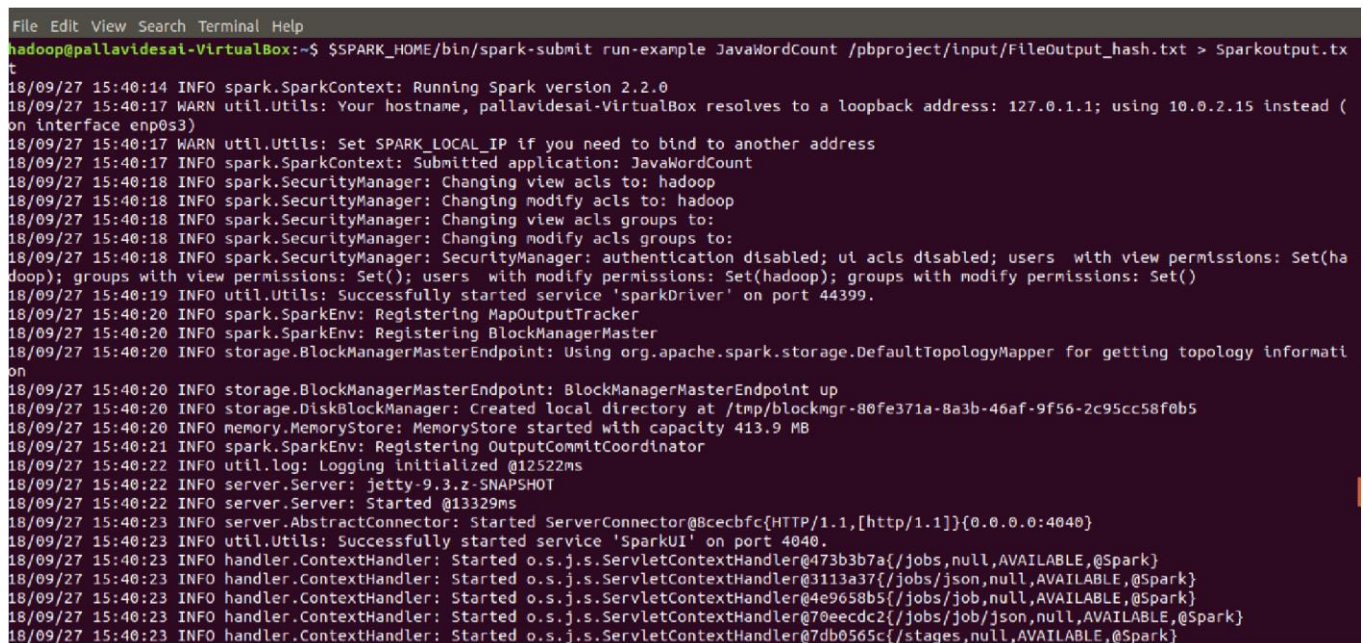

**Fig 7: Files in HDFS**



**Fig 8: Directory Information**

## Run a Word Count program in Apache Spark on the text file and store the output and log files locally.

- Then, after running the word count example on Hadoop, now it's time to run the same word count example using Apache Spark.

- The output obtained from the word count running on Apache Hadoop is almost similar to the output obtained from Apache Spark except the minor differences.



**Fig 9: Spark Commands**



**Fig 10: Sample Word Count Output**

Word count output Link:

https://github.com/SaitejaswiK/Principles-of-BigData-Management/blob/master/Source/WordCount%20Output/part-r-00000