# PRINCIPLES OF BIG DATA MANAGEMENT

## PHASE #1

**TEAM SIZE: 4**

**TEAM MEMBERS:** Sampath Gattu (sgvhb@mail.umkc.edu)
Moulika Chadalavada (mc7d8@mail.umkc.edu)
Prathyusha (ppf43@mail.umkc.edu)
Lakshmi Kona Nikitha (lnkwd5@mail.umkc.edu)

## Objective:

- The main aim of this phase is to develop a system to store, analyze, and visualize a social network's.
- Tasks:
    1. Collect social network's data (e.g. tweets) in JSON format.
    2. Store the text content (e.g. tweet's text) from the data into a file in HDFS.
    3. Run a Word Count program in Apache Spark on the text file and store the output and log files locally.

**Applications/Software's Used:** Apache Spark, Hadoop, Scala, Cloudera, Twitter Developer Account, Python.

## Collecting tweets from Twitter:
- Firstly, we have created a developer account in Twitter using below link.
    https://apps.twitter.com/

- Below are the variables that contains the user credentials to access Twitter API
    - ACCESS_TOKEN = "1974127951-N1QRXNCsVCywXl67BU1Wx7VlJ1fw2TlScZDY07s"
    - ACCESS_SECRET = "zLLfZ4ZF9BxvgzjXAkjJFAVFOL4P1i0fLSaZzc6LrnqZJ"
    - CONSUMER_KEY = "RfqrQLUtEAvwqSmNpGjZydmOn"
    - CONSUMER_SECRET = "1XdUvxMDfcL5eE89oAo7ZO8UESv6H7HacNa9B0Y7cGFa5MiPjo"

- We have written python program that is used to fetch tweets in JSON format. (Tweets-Collect.py)

- The extracted file in JSON format contains all the tweet details such as id, created at, text, profile_background_image_url etc.

- From JSON tweets file only the *text* content is extracted using Python program. The fetched text details are stored in a file. (TwitterTextConvert.py)
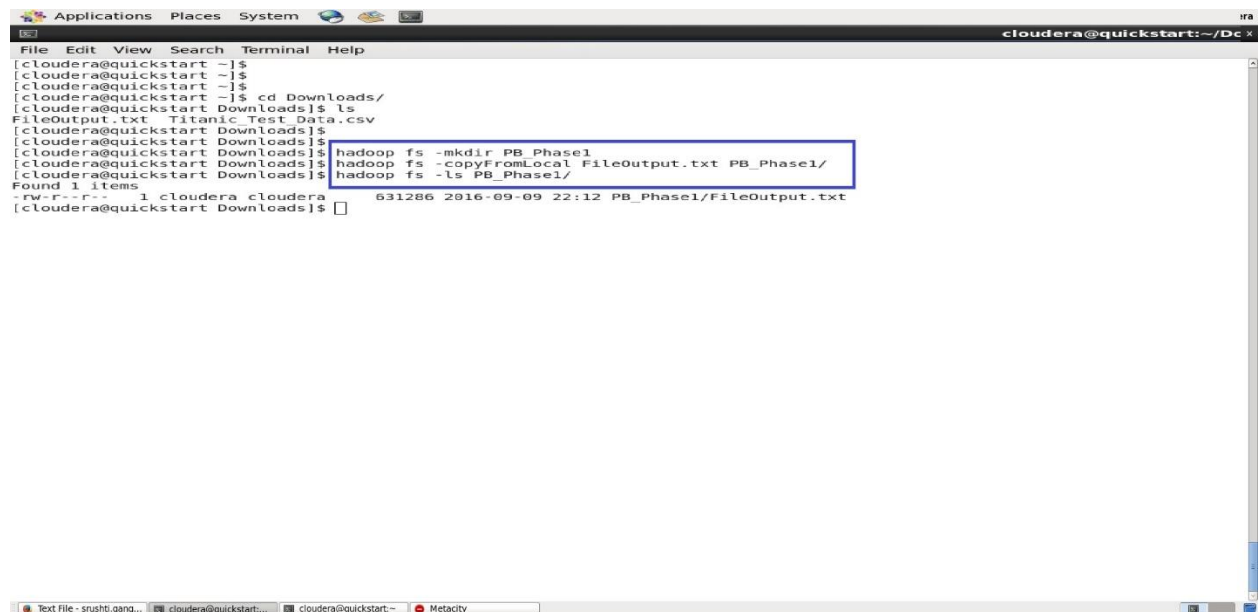
## Store the text content (e.g. tweet's text) from the data into a file in HDFS.

- The twitter tweets text content file is moved from local to HFDS.

- First a folder is created in HDFS and the text file is moved from local to HDFS using below command.

   **Create directory in local:** hadoop fs -mkdir PB_Phase1

   **Move text file from local to HDFS:** hadoop fs -copyFromLocal FileOutput.txt PB_Phase1/

   **To list the files under a directory:** hadoop fs -ls PB_Phase/



**Fig 1: HDFS Commands**

- The directory created and the files moved to HDFS can be viewed as shown below.



**Fig 2: Directory in HDFS**



**Fig 3: Files in HDFS**

**Fig 4: Directory Information**

## Run a Word Count program in Apache Spark on the text file and store the output and log files locally.
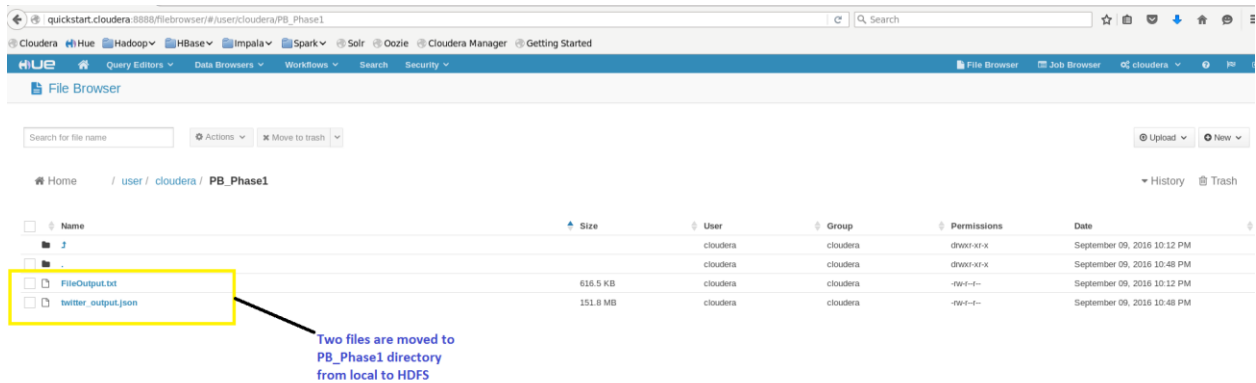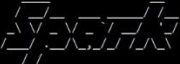
- First of all, to run word count program on set of data we require Apache Spark and scala.
- Once the Spark and Scala is installed successfully, in cmd prompt once the spark directory with **spark-shell** is given the Spark is integrated on Scala.
  Eg: C:/Users/Nik/Desktop/Spark>bin\spark-shell



**Fig 5: Spark**

- Once spark is opened successfully then word count program is written and executed in *scala.*
  1. **val textfile = sc.textFile("/Users/nikky/Desktop/PB Phase-1/FileOutput.txt")**
     - sc.textFile in Spark Shell, creates a RDD with each line as an element. For example, if there are 10 files in folder, 10 partitions will be created.
     - RDD means *Resilient Distributed Datasets* which are immutable and partitioned collection of records, which can only be created by coarse grained operations such as map, filter, group by etc.
     - RDDs can only be created by reading data from a stable storage such as HDFS or by transformations on existing RDDs.

  2. **val outcount = textFile.flatMap(line => line.split(" "))map(word => (word, 1))reduceByKey(_ + _)**
     - Once the file is read and stored in a variable textfile MapReduce function is used to execute word count program
     - MapReduce works by breaking the processing into 2 phases *Map Phase and Reduce Phase.*
     - Each phase has key-value pairs as input and output, by specifying two functions *Map Function and Reduce Function.*

  3. **outcount.saveAsTextFile("/Users/nikky/Desktop/PB Phase-1/WordCount")**
     - Once the reduce function group the words with respective count, the data is stored in WordCount folder.



**Fig 6: Word Count Commands**

```
part-00001
 1  (Installing,6)
 2  (https://t.co/zzYKV2ahXh,2)
 3  (https://t.co/PqHhQYEA8t,1)
 4  (Sulawesi.18KT,1)
 5  (#Fashion,1)
 6  (https://t.co/Im2dvNjzKZRT,4)
 7  (your...Tweets,1)
 8  (BALL,1)
 9  (fridge,1)
10  (AC~NEVADA,2)
11  (proponiendo,1)
12  (@mlp_RubyRay)WOW!!!,1)
13  (Trinity,1)
14  (https://t.co/PDjuTlyMFvMESMERIZING,1)
15  (Wonderlocke,2)
16  (https://t.co/HSp4OCkTVgRT,1)
17  (rialzo,1)
18  (end,7)
19  ((Nintendo,1)
20  (#046RT,4)
21  ("Red,3)
22  (https://t.co/tOrOskPC46,1)
23  (#trends,1)
24  (todo?,1)
25  (really,,2)
26  (https://t.co/kG08xQ0pqwReport:,1)
27  (9.,3)
28  (less,2)
29  (execute,1)
30  (https://t.co/7yZyGYfGSr,1)
31  (JD,5)
32  (04:25,1)
33  (Hiring,63)
34  (paso,3)
35  (shape,1)
36  (nosotros,2)
37  (different,12)
```

**Fig 7: Sample Word Count Output**