In this document the top 200000 post by viewcount data set is extracted from StackExchange Posts table and further cleaned and processed to obtain query results.

**Data Extraction:**

Link:  https://data.stackexchange.com/meta.codereview/query/new

Queries ran on the online terminal of StackExchange:

select top 50000 * from Posts where ViewCount > 100000 order by ViewCount desc;

select top 50000 * from Posts where ViewCount >= 65000 and ViewCount <= 100000 order by ViewCount desc;

select top 50000 * from Posts where ViewCount >= 48200 and ViewCount < 65000 order by ViewCount desc;

select top 50000 * from Posts where ViewCount >= 38194 and ViewCount < 48200 order by ViewCount desc;

select top 50000 * from Posts where ViewCount >= 38059 and ViewCount < 38194 order by ViewCount desc;

Few extra records were captured; hence those were deleted manually from the last dataset collected.

The files were merged into a single dataset Final_PostData.csv using Google Colab JupyterNotebook:
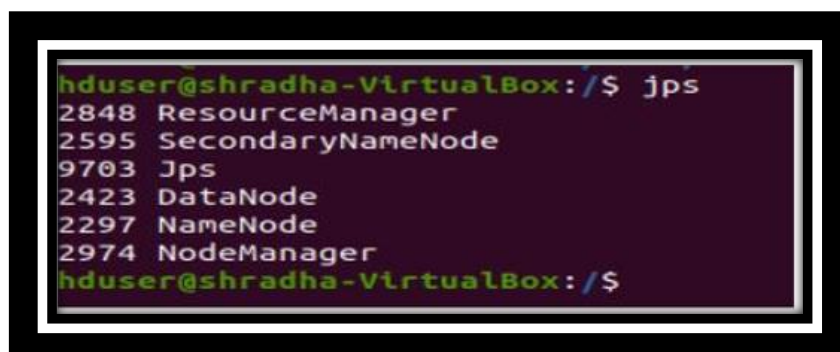


The merged file is then downloaded.

*Local Ubuntu Server: Once Hadoop Cluster (Namenode, Datanode)(v3.3.1), HIVE(v3.1.2) and PIG(v0.17.0) are installed, further steps can be followed.*
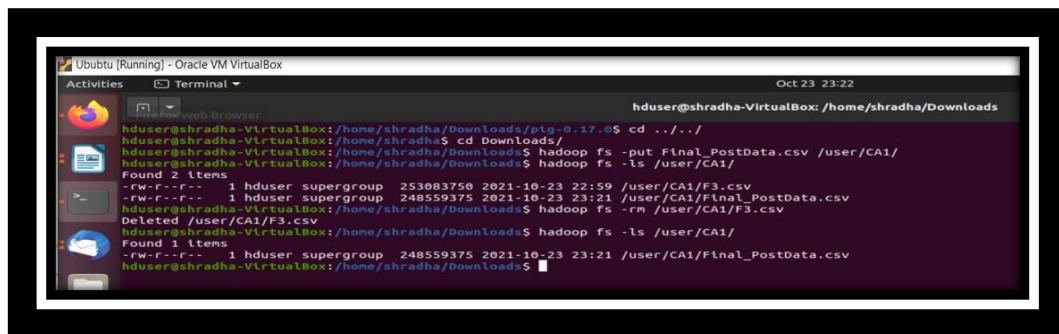*Cloud GCP: Since DataProc utility is utilized, all components are already installed.*

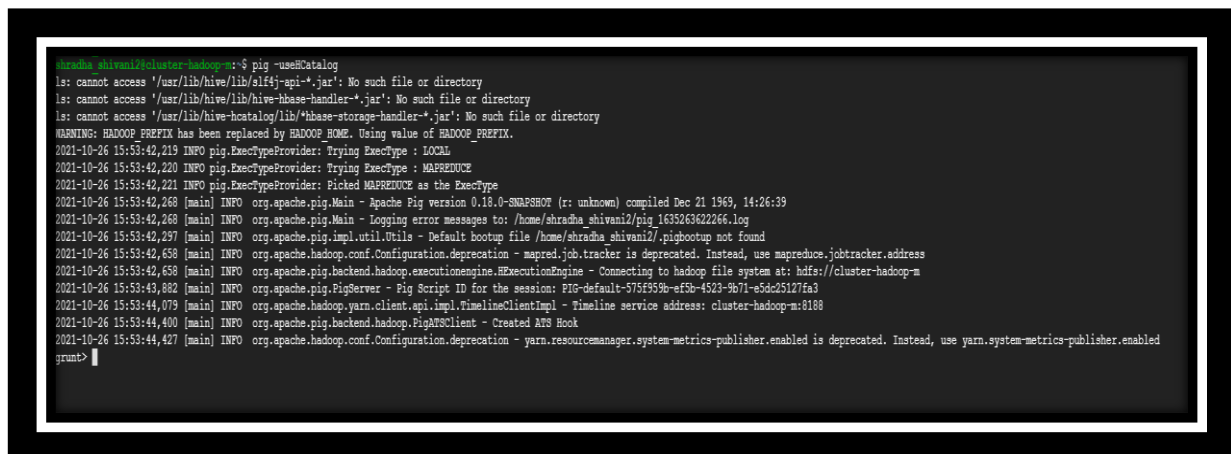1.  The status of Hadoop Cluster is checked.
    **Command**: jps

2. Dataset Final_PostData.csv is placed in HDFS location '/user/CA1'
   **Command**: hadoop fs -put Final_PostData.csv /user/CA1



3. Logged in to pig terminal using HCatalog. The Command used is pig -useHCatalog.



4. Using Pig, the data present in Final_Posts.csv is further cleaned. New lines, tab or carriage return characters, single characters, special symbols etc are replaced with space in 'Body' column. Additionally, the dataset is filtered for NOT NULL OWNERUSERID tuples.

**Command:**

grunt>loadposts = load 'hdfs:///user/CA1/Final_PostData.csv' using org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'YES_MULTILINE','NOCHANGE','SKIP_INPUT_HEADER') as(id:int,posttypeid:int,acceptedanswerid:int, parentid:int,creationdate:DATETIME,deletiondate:DATETIME,score:int,viewcount:int,body:chararray,owneruserid:int,ow nerdisplayname:chararray,lasteditoruserid:int,lasteditordisplayname:chararray,lasteditdate:DATETIME,lastactivitydate:DA TETIME,title:chararray,tags:chararray,answercount:int,commentcount:int,favoritecount:int,closeddate:DATETIME,commu nityowneddate:DATETIME,contentlicense:chararray);
grunt>posts = foreach loadposts generate id,posttypeid,acceptedanswerid,parentid,creationdate,deletiondate,score,viewcount,REPLACE(body,'\\n','') as body,owneruserid,ownerdisplayname,lasteditoruserid,lasteditordisplayname,lasteditdate,lastactivitydate,title,tags,answerco unt,commentcount,favoritecount,closeddate,communityowneddate,contentlicense;
grunt> posts = foreach posts generate id,posttypeid,acceptedanswerid,parentid,creationdate,deletiondate,score,viewcount,REPLACE(body,'\\t','') as body,owneruserid,ownerdisplayname,lasteditoruserid,lasteditordisplayname,lasteditdate,lastactivitydate,title,tags,answerco unt,commentcount,favoritecount,closeddate,communityowneddate,contentlicense;
grunt> posts = foreach posts generate id,posttypeid,acceptedanswerid,parentid,creationdate,deletiondate,score,viewcount,REPLACE(body,'\\r','') as body,owneruserid,ownerdisplayname,lasteditoruserid,lasteditordisplayname,lasteditdate,lastactivitydate,title,tags,answerco unt,commentcount,favoritecount,closeddate,communityowneddate,contentlicense;
grunt> formatted_posts = FOREACH posts GENERATE  id AS id, score AS score, REPLACE(body,'*','') AS body, owneruserid AS owneruserid, REPLACE(title,'*','') AS title, REPLACE(tags,'*','') AS tags;

grunt>
grunt> valid_posts = FILTER formatted_posts BY (owneruserid IS NOT NULL) AND (score IS NOT NULL);
grunt>



5.  The cleaned file is then stored in HDFS at location /user/CA1/Output. (Note: while implementing in Google Cloud Platform, the cleaned file was stored in HDFS at location /user/CA1/combined')
    **Command:**

    grunt> store valid_posts into 'hdfs:///user/CA1/Output';



6.  After successful log generation when the store command was run, the cleaned files can been seen at '/user/CA1/Output'
    **Command:** hadoop fs -ls /user/CA1/Output



7.  The _SUCCESS log was removed.
    **Command:**
    hduser@shradha-VirtualBox:/home/shradha/Downloads$ hadoop fs -rm /user/CA1/Output/_SUCCESS
    Deleted /user/CA1/Output/_SUCCESS
    hduser@shradha-VirtualBox:/home/shradha/Downloads$

8.  HIVE_POSTS table was created in HIVE.
    **Command:**

    > hduser@shradha-VirtualBox:/home/shradha/Downloads$ hive
    SLF4J: Class path contains multiple SLF4J bindings.
    SLF4J: Found binding in [jar:file:/home/shradha/Downloads/apache-hive-3.1.2-bin/lib/log4j-slf4j-impl-
    2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
    SLF4J: Found binding in [jar:file:/home/shradha/Downloads/hadoop/share/hadoop/common/lib/slf4j-log4j12-
    1.7.30.jar!/org/slf4j/impl/StaticLoggerBinder.class]
    SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
    SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
    Hive Session ID = 81dfa9d9-3907-4af3-aa07-dc64f299f967

    Logging initialized using configuration in jar:file:/home/shradha/Downloads/apache-hive-3.1.2-bin/lib/hive-common-
    3.1.2.jar!/hive-log4j2.properties Async: true
    Hive Session ID = b339dd7b-291a-492e-b659-3e16121185b8
    Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution
    engine (i.e. spark, tez) or using Hive 1.X releases.
    hive> CREATE TABLE HIVE_POSTS (id int, score int, body String, owneruserid Int, title String, tags String) ROW
    FORMAT DELIMITED
        > FIELDS TERMINATED BY ','

    OK
    Time taken: 1.393 seconds
    hive> show tables;
    OK
    hive_posts
    Time taken: 0.335 seconds, Fetched: 1 row(s)
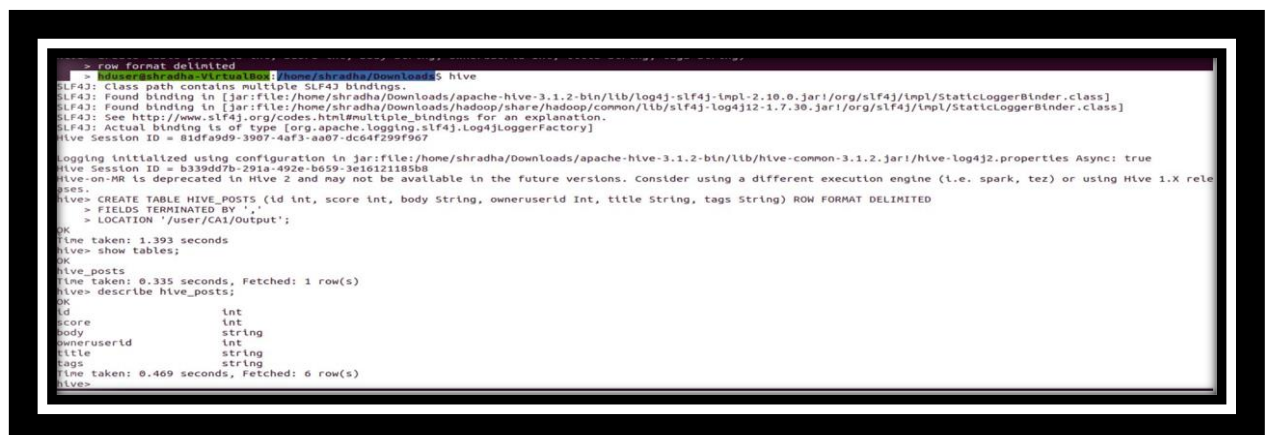    hive> describe hive_posts;
    OK
    id                      int
    score                   int
    body                    string
    owneruserid             int
    title                   string
    tags                    string
    Time taken: 0.469 seconds, Fetched: 6 row(s)
    hive>



9.  The cleaned file part-m-00001 at HDFS location '/user/CA1/Output/' was loaded to the created HIVE_POSTS table using
    HIVE LOAD command. The cleaned file is selected depending on the size and content of the file. For this implementation
    part-m-00001 was used for further processing.
    **Command:**
    hive> LOAD DATA INPATH 'hdfs:///user/CA1/Output/part-m-00001' INTO TABLE hive_posts;

### Queries Ran:

### Query1:
Use Pig/Hive/MapReduce - Extract, Transform and Load the data as applicable to get the top 10 posts by score

### Command:
hive> select id, score, title from hive_posts order by score desc limit 10;



### Query2:
Use Pig/Hive/MapReduce - Extract, Transform and Load the data as applicable to get the top 10 users by post score

### Command:
hive> select owneruserid,sum(score) as Total_Score from hive_posts group by owneruserid order by Total_Score desc limit 10;
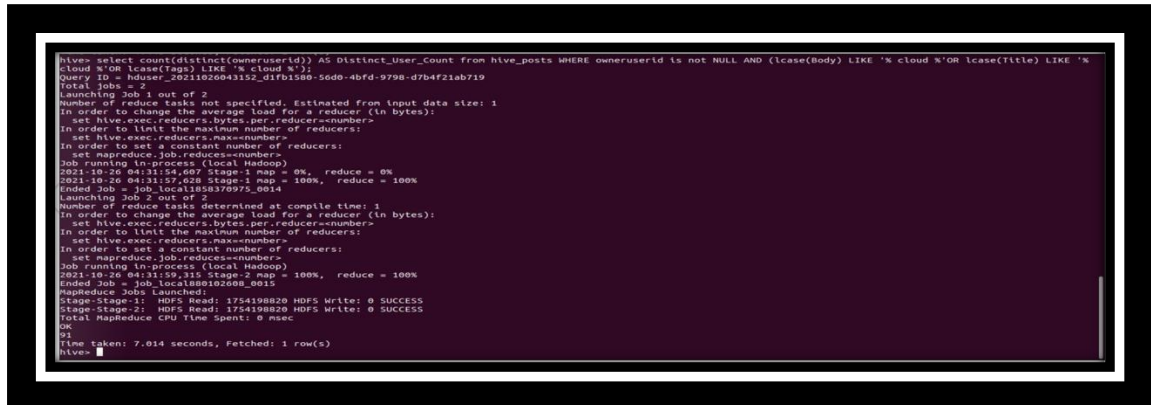
**Query3:**

Use Pig/Hive/MapReduce - Extract, Transform and Load the data as applicable to get The number of distinct users, who used the word "cloud" in one of their posts.

**Command**: Note( Anyone of these commands can be used for this query)
hive> select count(distinct(owneruserid)) AS Distinct_User_Count from hive_posts WHERE locate (" cloud ",concat(Body,Title,Tags))>0;

<div align="center">OR</div>

hive> select count(distinct(owneruserid)) AS Distinct_User_Count from hive_posts where (Body like '% cloud %' OR title like '% cloud %' OR tags like '% cloud %');



We can used hive -e directly form the terminal to run the above 3 Hive queries and capture the output in a text file.

**Commands:**

**Query1:**

hduser@shradha-VirtualBox:/home/shradha/Downloads/hadoop/sbin$ hive -e "SELECT ID, SCORE, TITLE FROM HIVE_POSTS ORDER BY SCORE DESC LIMIT 10;" > /home/shradha/Downloads/Query1_Output.txt

**Query2:**

hduser@shradha-VirtualBox:/home/shradha/Downloads/hadoop/sbin$ hive -e "SELECT OWNERUSERID, SUM(SCORE) as TOTAL_SCORE FROM HIVE_POSTS GROUP BY OWNERUSERID ORDER BY TOTAL_SCORE DESC LIMIT 10;" > /home/shradha/Downloads/Query2_Output.txt

**Query3:**

hduser@shradha-VirtualBox:/home/shradha/Downloads$ hive -e "SELECT COUNT(DISTINCT(OWNERUSERID)) AS DISTINCT_OWNERS_COUNT FROM HIVE_POSTS WHERE LOCATE(' cloud ',concat(BODY,TITLE,TAGS))>0;" > /home/shradha/Downloads/Query3_Output.txt