

# **CSE 4/587 Data-Intensive Computing**

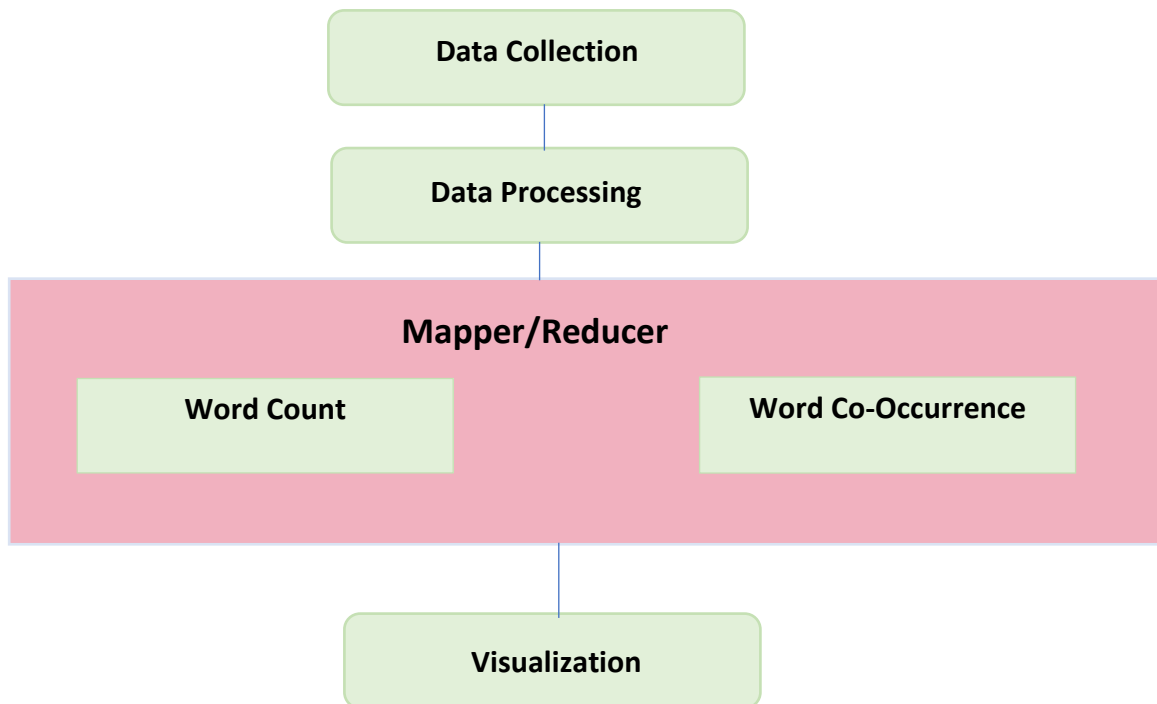
## **LAB2: DATA AGGREGATION, BIG DATA ANALYSIS AND VISUALIZATION**

Sonali Yeshwant Naidu (50299652)  
Pratibhaa Reddy Kandimalla (50299669)

## Introduction:

In this lab we worked on Data Aggregation, Big data analysis and Data Visualization. The topic chosen for this Data Analysis is 'Sports'.

Below is the block diagram showing the different steps involved.



## Part 1

### Data Collection:

- Collected data from three different sources such as "Twitter", "NY Times" and "Common Crawl".
- We have extracted the data related to topic "Sports" also using the subtopic/keywords for better data collection such as Basketball, Football, Golf, Badminton, Tennis, Cricket.
- The above data was collected using APIs for twitter and NY Times.

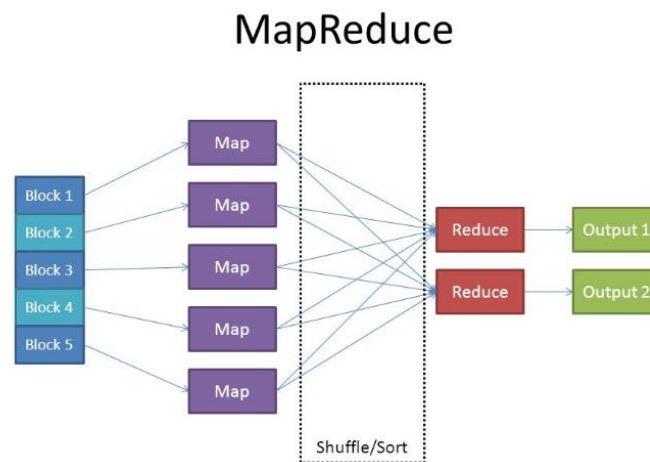
### Data Processing:

- The data extracted from the above sources i.e. Twitter and NY Times was cleaned for any duplicates, stopwords, unwanted charactes like punctuation marks, special characters etc.
- Downloaded and used NLTK in python to preprocess the same.

### Mapper/Reducer:

- Once HDFS and Virtual Machine were setup the WordCount and Word Co-Occurrence was run on the above processed data.
- On successful execution of Mapper/Reducer programs, the words along with their count was their output.

- We ran the word cooccurrence program on the same data collected.



#### Data Visualization:

- The output files of Mapper/Reducer for the WordCount and the Word Co-occurrence is loaded into Tableau for data visualization.
- Data processed is visualised in Tableau in the WordCloud format.

#### Project Extension:

The same project can be used on any other data sets for Data Collection, Analysis and Visualization. The data could be from any varied sources, which can be preprocessed and further analyzed. This data can be visualized using Tableau or any other data visualization tool.

## Part 2

We used the VM Box to setup the Hadoop Infrastructure for storing the Big Data collected from Twitter, NYTimes and Common Crawl. The mapper/reducer was run to clean and parse the data sets into words, remove stop words, stem the words (ex: running to run) and reduce will count the useful words.

**Word Count:** This Program is run on the data collected to parse the words and get the count of each word in the data which is later used for data visualization to understand the most commonly used words from Twitter, NYTimes and CommonCrawl.

**Word Co-Occurrence:** This Program was run to find the word co-occurrences in the collected data.

Below is the screenshot from VM Box with the hdfs directory and certain commands for running the Mapper/Reducer in it.

Few of the commands used were:

- `hdfs dfs -ls /` to view the list of directories

- `hadoop jar hadoop-3.1.2/share/hadoop/tools/lib/hadoop-streaming-3.1.2.jar -file /home/cse587/hduser/mapper.py -mapper mapper.py -file /home/cse587/hduser/reducer.py -reducer reducer.py -input /nytimes/nydata.txt -output /wc_nyt/`
- 

```

cse587@cse587:~$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
localhost: namenode is running as process 2504. Stop it first.
Starting datanodes
localhost: datanode is running as process 2602. Stop it first.
Starting secondary namenodes [cse587]
cse587: secondarynamenode is running as process 2884. Stop it first.
Starting resource manager
resource manager is running as process 3171. Stop it first.
Starting node managers
localhost: nodemanager is running as process 3334. Stop it first.
cse587@cse587:~$ hdfs dfs -ls /
Found 19 items
drwxr-xr-x - cse587 supergroup 0 2019-04-21 02:10 /a
drwxr-xr-x - cse587 supergroup 0 2019-04-21 02:23 /hoom
drwxr-xr-x - cse587 supergroup 0 2019-04-21 19:32 /nyt_ss
drwxr-xr-x - cse587 supergroup 0 2019-04-19 19:15 /nytimes
drwxr-xr-x - cse587 supergroup 0 2019-04-21 01:13 /p
drwxr-xr-x - cse587 supergroup 0 2019-04-21 01:32 /s
drwxr-xr-x - cse587 supergroup 0 2019-04-21 01:42 /small
drwxr-xr-x - cse587 supergroup 0 2019-04-21 02:21 /smallwordco
-rw-r--r-- 3 cse587 supergroup 80 2019-04-20 21:01 /test
-rw-r--r-- 3 cse587 supergroup 80 2019-04-19 18:09 /test1
drwxr-xr-x - cse587 supergroup 0 2019-04-21 19:27 /twitter_ss
drwxr-xr-x - cse587 supergroup 0 2019-04-21 03:16 /w
drwxr-xr-x - cse587 supergroup 0 2019-04-21 19:13 /wc_nyt_ss
drwxr-xr-x - cse587 supergroup 0 2019-04-21 15:53 /wc_nytimes
drwxr-xr-x - cse587 supergroup 0 2019-04-21 03:13 /wc_twitter
drwxr-xr-x - cse587 supergroup 0 2019-04-21 19:29 /wc_twitter_ss
drwxr-xr-x - cse587 supergroup 0 2019-04-21 18:09 /wco_nytimes
drwxr-xr-x - cse587 supergroup 0 2019-04-21 01:36 /y
drwxr-xr-x - cse587 supergroup 0 2019-04-21 01:44 /y
cse587@cse587:~$ hdfs dfs -rm -r /nyt_ss
Deleted /nyt_ss
cse587@cse587:~$ hdfs dfs -mkdir /nyt_ss
cse587@cse587:~$ hdfs dfs -put /home/cse587/examplehadoop/nydata_smallerSet_processed.txt /nyt_ss
cse587@cse587:~$ hdfs dfs -rm -r /wc_nyt_ss
Deleted /wc_nyt_ss
cse587@cse587:~$ hadoop jar hadoop-3.1.2/share/hadoop/tools/lib/hadoop-streaming-3.1.2.jar -file /home/cse587/examplehadoop/mapper.py -mapper mapper.py -file /home/cse587/examplehadoop/reducer.py -reducer
reducer.py -input /nyt_ss/nydata_smallerSet_processed.txt -output /wc_nyt_ss/
2019-04-21 19:47:21,062 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [/home/cse587/examplehadoop/mapper.py, /home/cse587/examplehadoop/reducer.py] [] /tmp/streamjob89021804618849675.jar tmpDir=null
2019-04-21 19:47:23,083 INFO Impl.MetricsConfig: loaded properties from hadoop-metrics2.properties
2019-04-21 19:47:23,335 INFO Impl.MetricsSystemImpl: Scheduled metric snapshot period at 10 second(s).
2019-04-21 19:47:23,336 INFO Impl.MetricsSystemImpl: JobTracker metrics system started
  
```

## Part 3

The mapper/reducer and word co occurrence code was run on all the three sets of data and visualized using Tableau.

Below are the screenshots of the code visualized for the word count and word co occurrences data.

**Word Clouds:**

**Word Count:**

Wordcount on Twitter:



Wordcount on NYTimes:



Word Count on smaller set of NYTimes data:



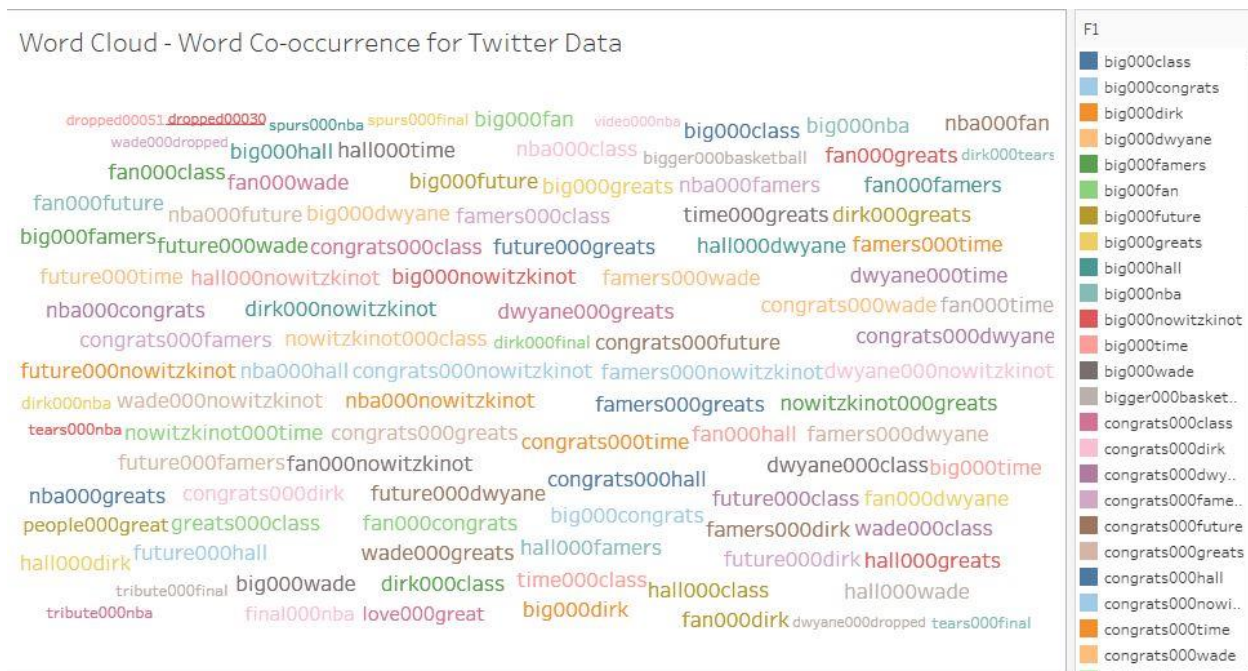


### Wordcount on Common Crawl:



### Word Co-Occurrence:

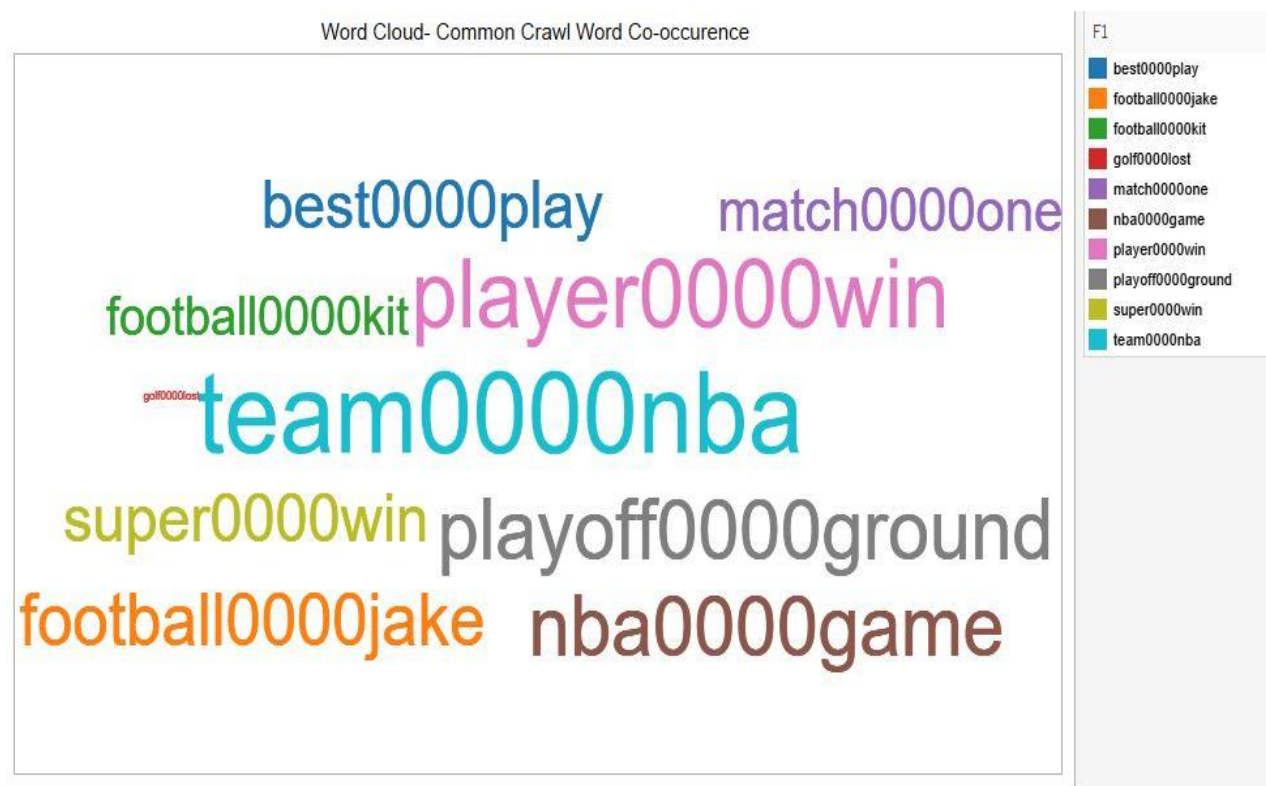
Word Co-occurrence on Twitter data:



Word Co-occurrence on top 30 NYTimes data:



Word Co-occurrence on top 10 Common Crawler data:





## Directory Structure:

### *sonaliyeLab2.zip*

- **Report.pdf**
- **Video.mp4**
- **part1 (folder)**
  - **Code** *included all the scripts related to data collection for Common Crawl, Twitter and NYT*
  - **Data**
    - Twitter -> Twitter Data, Twitter Data Processed, Twitter WC and Twitter Co-occurrence
    - NYT -> NYTimes Data, NYTimes Data Processed, NYTimes and NYTimes Co-occurrence
    - Commoncrawl -> Common Crawl Data, Common Crawl and Common Crawl Co-occurrence
- **part2 (folder)**
  - *mapper.py reducer.py, pairsmapper, dataprocessing*
  - *other files (screenshots/results)*
- **part3 (folder)**
  - **Twitter (folder)**
    - Code (folder) -> *mapper.py reducer.py, pairsmapper, dataprocessing*
    - Images (folder) -> *all visualizations .jpg/.png AND .js /.twbx files*
  - **NYT (folder)**
    - Code (folder) -> *mapper.py reducer.py, pairsmapper, dataprocessing*
    - Images (folder) -> *all visualizations .jpg/.png AND .js /.twbx files to create NYT images*
  - **Commoncrawl (folder)**
    - Code (folder) -> *mapper.py reducer.py, pairsmapper, dataprocessing*
    - Images (folder) -> *all visualizations .jpg/.png AND .js /.twbx files to create CC images*

## References:

NYTimes: [https://www.youtube.com/watch?v=zXif\\_9RVadI](https://www.youtube.com/watch?v=zXif_9RVadI)

CommonCrawl: <https://rushter.com/blog/python-fast-html-parser/>