# Data Intensive Computing

## LAB – 2
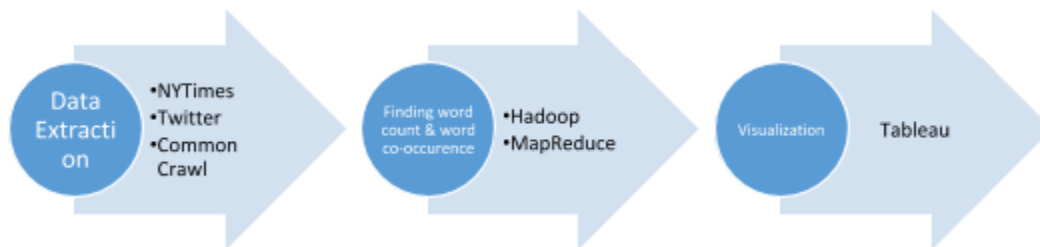
### Vineel Kurma, vineelku, 50291799

### Meghana Vadlapally, mvadlapa, 50298635

**Project Exploration:**

- In this project, we collected sports data from NYTimes, Twitter, Common crawl using APIs.
- Python has been used as a tool to collect data from the sources.
- We used VM to set up HDFS for applying big data analytical methods on the collected data.
- The obtained data has been given as an input to the MapReduce to find word count and word co-occurrence.
- The extracted word count and word co-occurrence files are visualized by creating word cloud using Tableau.

**Flow chart:**

1.  **Data Extraction:** Contains sports data from NYTimes, Twitter, Common Crawl.
    **Tool:** Used python to extract data
    - **NYTimes code**: nytimescode.ipynb
    - **Twitter code**:  used LAB-1's code
    - **Common Crawl**: CommonCrawlGetPy2.py, CommonCrawlEditPy3.py

**SubTopics:**
Used the FootBall, BasketBall, Cricket, NFL, GOLF as subtopics to extract the data.

**Twitter Data Collection:**
For collecting the tweets used Tweepy library and used the twitter API credentials created from the previous lab, referenced the code from the GitHub.
link:https://github.com/shuzhanfan/Geo-tagged-streaming-tweets-collect.
After collecting the data, preprocessed by removing the duplicate tweets and used a combination of tweets cleaning library and regular expression to filter out emojis and non-ASCII characters. Overall we have collected around 60,000 total tweets in which  25K tweets are unique.

**NewYork Times Data Collection:**
For collecting the NewYork Times data, created an developer API and extracted the URLs and put them into a file. Then for each URL, used request module and extracted the raw HTML data. After fetching the HTML data, used the beautiful soup library and extracted the articles. In total, we have collected around 500 articles

**Common Crawl Data Collection:**
For collecting the Common Crawl Data, used the example project code as reference https://rushter.com/blog/python-fast-html-parser/. But since the example is in Python 2 and the library WARC is not supported in Python 3, created an python 2.7 environment in Anaconda and used the code to collect the data. To download the data used the https://index.commoncrawl.org/ link. From the link, we get a unique WARC file path and then filtered out the relevant sports links. Used the major sports website links and extracted the WARC files. Did a lot of hits and trails with various links as most of them have irrelevant data. Then used the Python 3 beautiful soup code used for NYTimes and extracted the data. We have collected around 450 links and articles.

2.  **Word Count:** Word count of the obtained data has been done using the MapReduce model (Hadoop).  In the programming used, we imported re library to extract only alphanumeric characters. Also, we are filtering the stop words we got from the NLTK library stop words and also we have manually added additional ~15 stop words.

**Tool:** We used a mapper.py, reducer.py files provided in cse587.



In the HDFS, a directory was made using

Hdfs dfs –mkdir /filename

The obtained data has been sent to the HDFS directory using

Hdfs dfs –put /home/cse587/examplehadoop/dataContainingFolder.txt /filename

The MapReduce tool was used on the uploaded file using

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
/filename/dataContainingFolder.txt -output /outputfolder

Once the data has been processed with MapReduce the created file (A part-00000 file has been
created containing word count) has been taken from HDFS to disk using

hdfs dfs -get /outputfolder/part-00000  /home/cse587/examplehadoop/results.txt

3. **Word co-occurrence:** To find the word co-occurrences of the collected data we used MapReduce

**Tool:** We used Python (Pairsmapper.py) programming language as a tool for the MapReduce implementation.

**To Visualise the word co-occurrence using Tableau, we faced issues with separators. So we used the number 98765 as the separator in visualisation.**

4. **Visualization:** To visualize the obtained word count and word co-occurrence, we used Tableau forming word cloud.

**Link to our Tableau server:**
https://public.tableau.com/profile/vineel.kurma#!/vizhome/TwitterWordCoOccuranceAndWordCount/NyTimesWordCount

Below are the screenshots of the word clouds of NYTimes, Twitter, Common Crawl.

## NYTimes:

**Word count:**



NyTimesWordCount

left
back since two player made world teams
final game would last years scored first
team games season said points
coach players one year league
soccer time nba also city win new like

**Word co-occurrence:**



NYTimesCoccurance

one98765last
scored98765first   league98765team   world98765cup
last98765one   last98765first   one98765two
first98765points premier98765league   game98765one team98765first
points98765points, one98765players   first98765game
one98765team playoff98765points last98765league points,98765points
rebounds98765points two98765first points98765rebounds last98765team
points98765first   scored98765points points9876510 game98765points
first98765team   manchester98765city   game98765first points98765game
first98765time points98765scored   team98765one   league98765one
one98765first   champions98765league   points98765playoff
first98765last one98765league league98765last   team98765players
first98765league   first98765one team98765league two98765one
one98765game team98765last   one98765player

# Twitter:

**Word count:**

TwitterWordCount

player one
greats class great hall life fan people dirk
act dwyane famers win one last dance
congrats nfl basketball bigger final nba
love wade game future spurs nowitzkinot last
time big watch tribute really tears
video

**Word co-occurrence:**

TwitterCoOccurance

fan98765time
big98765hall fan98765dirk nba98765hall nba98765time hall98765class fan98765future
fan98765class dirk98765class fan98765greats big98765greats wade98765class
future98765hall nba98765greats fan98765dwyane hall98765famers
time98765greats future98765class congrats98765hall greats98765class big98765dwyane
big98765nba famers98765class dwyane98765greats future98765wade big98765congrats
fan98765hall congrats98765class dwyane98765dirk dwyane98765class dirk98765time nba98765congrats
future98765greats nowitzkinot98765time congrats98765wade congrats98765future
dirk98765nowitzkinot wade98765nowitzkinot wade98765dirk big98765dirk famers98765nowitzkinot
nowitzkinot98765greats congrats98765nowitzkinot dwyane98765nowitzkinot dwyane98765wade
future98765nowitzkinot congrats98765famers congrats98765greats nowitzkinot98765class
congrats98765dwyane hall98765nowitzkinot future98765dwyane big98765nowitzkinot big98765wade
nba98765wade nba98765nowitzkinot fan98765nowitzkinot congrats98765dirk congrats98765time nba98765fan
dirk98765greats future98765famers famers98765dwyane dwyane98765time famers98765time nba98765dwyane
fan98765famers wade98765greats famers98765greats nba98765future hall98765dwyane nba98765famers
hall98765greats future98765dirk famers98765wade wade98765time future98765time big98765time
fan98765wade time98765class famers98765dirk hall98765wade big98765famers
hall98765time nba98765class big98765future big98765class hall98765dirk
big98765fan nba98765dirk

**Common Crawl:**

**Word count:**



**Word co-occurrence:**

CommonCrawlWordCoOccurance

op98765verkoop th98765ng
op98765voorraad padding98765px
voorraad98765door
door98765spartoo
nh98765ng passport98765agency
voorraad98765verkoop ng98765th
verkoop98765door important98765padding
px98765importantverkoop98765spartoo
spartoo98765nl ng98765nh
spartoo98765op ch98765ng
nl98765op

**Observations:**

- The main issue observed in this project is there are a lot of irrelevant words which are not part of mainstream stop words from NLP libraries, hence the word cloud may not exactly reflect the keywords of our topic.
- Stopped doing stemming and lemmatization as we were getting undesired output with many invalid words.
- The problem with co-occurrence is that the pairs are duplicated I.E pair (a,b) was again observed as (b, a), to avoid that we need to store the pair so that we don't have to repeat the words.
- While working with common crawl, we have observed a lot of junk data, even after filtering with keywords. Due to lack of time, we could not proceed further with cleaning the data.