# LAB2: DATA AGGREGATION, BIG DATA ANALYSIS AND VISUALIZATION:
# CSE- 587

**SUBMITTED BY:**

**ABHISHEK MUNI**
**UB ID: 50291410**

**SHWETASREE CHOWDHURY**
**UB ID: 50296995**

**NBA**

Data aggregation from Twitter, New York Times and Common Crawl and their analysis in the following sequence-

1. Data aggregation from more than one source using the APIs (Application programming interface) exposed by data sources.
2. Applying classical big data analytic method of MapReduce to the unstructured data collected.
3. Store the data collected on WORM infrastructure Hadoop using S3.
4. Building a visualization data product

**IMPLEMENTATION:**

**1. Data aggregation:**

**a. Data aggregation from tweets:**

For collecting data from twitter,  API search  was used-

- A twitter developer account was created  and an APP was created within it from which the consumer_ key,consumer_token, acess_key and access token were taken-

**Keys and tokens**

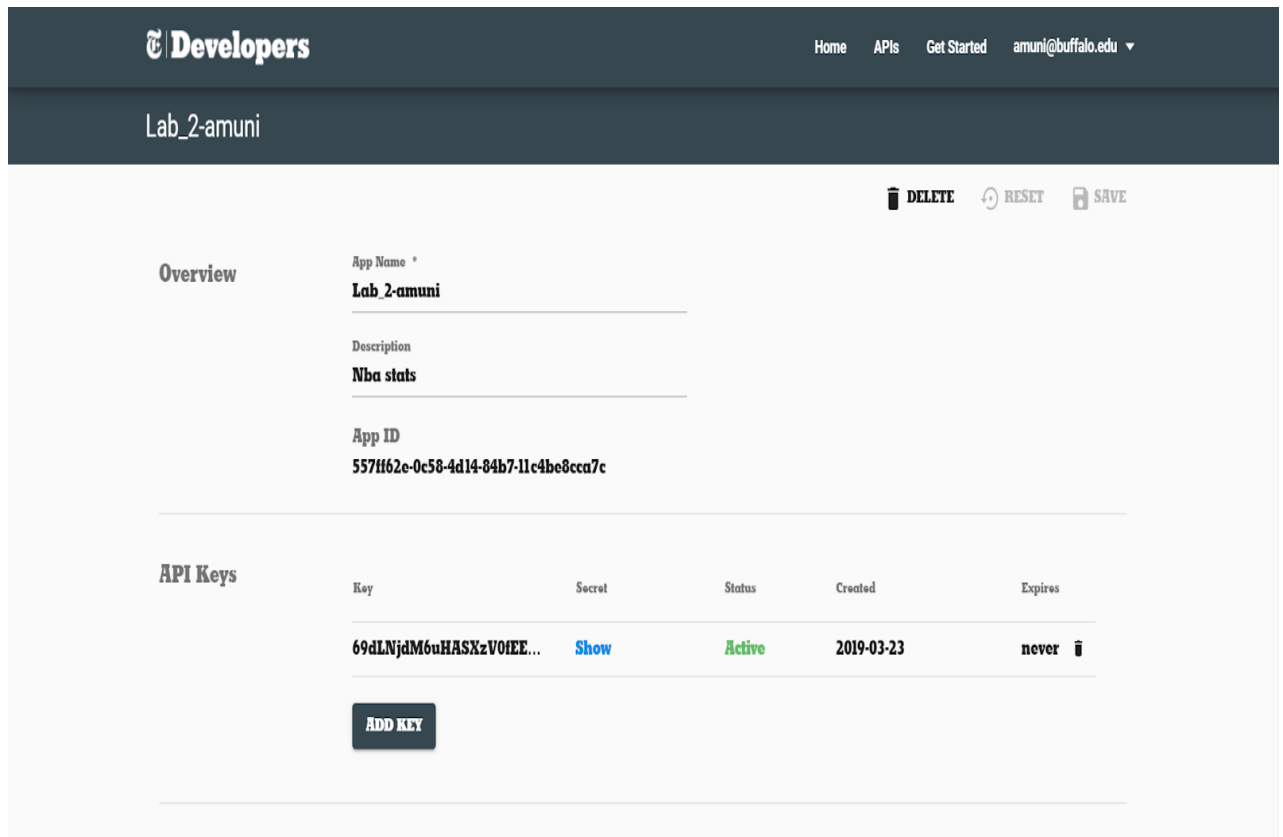Keys, secret keys and access tokens management.

**Consumer API keys**

Guxhi58a04HKqNbiOvLyTfZcL (API key)

e31biYv1TwyVIeDmQ59FtgkIPSocUSmrOlIa6wBWfApcTnD6m1 (API secret key)

Regenerate

**Access token & access token secret**

1097193550196944896-p5aweBkg79oxWL6FxmEPhItiOYCd2c (Access token) NEW

XtHHZnh8N398m2DxPbxItHpqE6RK6aPwudWuf5D9K5G68 (Access token secret) NEW

Read and write (Access level)

Revoke     Regenerate

- These were used in Rtweet package and the tweets were collected using the keyword 'NBA' -
- The duplicate tweets were removed and were written into a csv file.

**b) Data aggregation from NewYork Times**:

- For collection of data from New york times, we used the article search API in the NYT app.
- The access token was taken from above and the URLs related to the keyword 'NBA' was obtained using the NYT Article search package in python.
- Then using Beautiful Soup package the meta-data was obtained from the URLs .
- Also using the above package,HTML parsing was done and the output was written in a text file.Below given is the screenshot of the NYT app from which the article search API was obtained.



**c) Data aggregration from common crawl :**

- We have used python to mine common crawl data.
- Common Crawl is a gigantic dataset that is created by crawling the web. They provide the data in both downloadable format (gigantic) or we can query against their indices and only retrieve back the information you are after.
- We access the compressed archive files stored on Amazon S3 and pull out the actual content.

**d)** After data was obtained from three different sources some amount of pre-processing was done before sending it for Big data analysis. The following operations were performed-

- **Stemming operation**: Stemming operation is performed to produce different morphological variants of the stem word. This is used to put together the words that are different in grammatical perspective, but are same for the purpose of data analysis.

- **Stop words:** These are a set of words like grammatical articles('the','is','in') which are redundant for data analysis. These were removed before the file was sent for big-data processing.

- The above implementations were used in a python script to filter out the same and provide the results which were then fed into the MR framework.

**2) Application of classical big data analytic method of MapReduce to the unstructured data collected-**

We have used the AWS Elastic Map Reduce framework to achieve the implementation.

AWS EMR is a flexible and scalable approach and provides an immense enterprise level flexibility when it comes to working on Big Data.

**Setting up AWS:**

We have largely followed the traditional AWS documentation to create an EMR cluster, create a key pair, add steps and so on.

https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-gs-launch-sample-cluster.html

**AWS Architecture:**

**Cluster Dashboard:**



**Properties of a cluster**



**Properties of an MR Streamline step:**

*We use custom jar option if we are to run the MR job on the basis of a jar file.

*We use streamlined operation if we are to use .py approach to separately use python files as mapper, reducer, locations of input and output folder.

## Setting up the S3 bucket:

Amazon S3 is cloud storage for the internet. To upload your data (photos, videos, documents etc.), you first create a bucket in one of the AWS Regions. You can then upload any number of objects to the bucket.

In terms of implementation, buckets and objects are resources, and Amazon S3 provides APIs for you to manage them. For example, you can create a bucket and upload objects using the Amazon S3 API. You can also use the Amazon S3 console to perform these operations. The console uses the Amazon S3 APIs to send requests to Amazon S3.

Files can be uploaded on a drag and drop basis.

**4) Data visualization using Tableau-**

Finally the outputs obtained after the big data analysis was visualized using Tableau. The outputs obtained after running word-count on three different data-bases were as follows-

1. **Output for running word-count on twitter dataset:**



**Interpretations:**

- As we were doing data twitter data analysis on the key-word nba we could see highest number of "**nba**" words in the related tweets.
- As the **playoff** season of nba-2019 was onset in mid-march, we also find considerable amount of playoff related tweets during this time of the year.
- As "James harden" is one of the leading scores in this season, we get some considerable number of tweets related to him.
- Also we could see many tweets related to the team that are table-toppers in both the conferences like Boston Celtics, mlb, Toronto Raptors, Golden states warriors
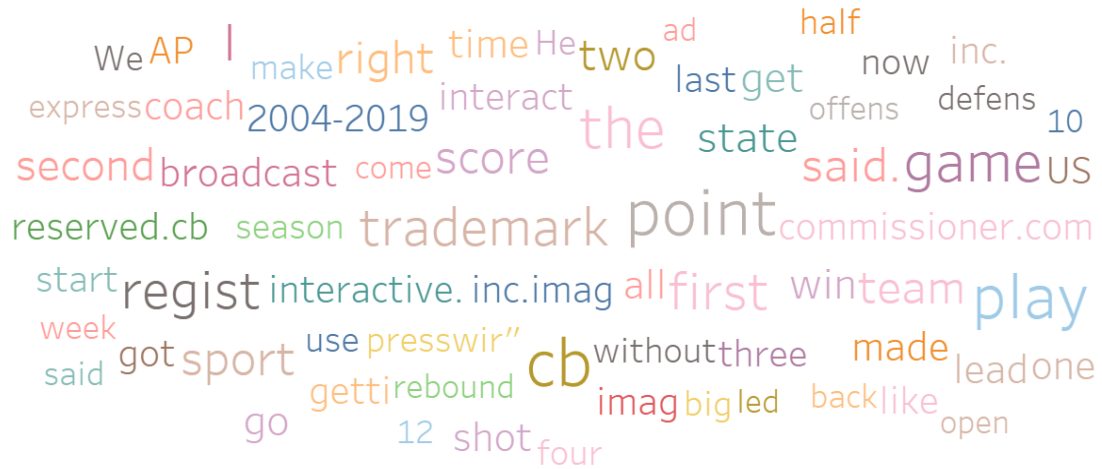
**2)Output for running word-count on NewYork Times data:**



**Interpretation**:

- As we searched for nba keyword for article search  we could find considerable **nba**  words inside the articles.
- As **golden state warrior** has been the league champion for past couple of seasons we could find many analyses on the same team.
- As nba is a basketball league we could find many **basketball** words in the article analysis.
- As nba involves a lot of matches we could find many analyses related to winning of team.
- As recently we had summer **draft** for the nba we could find some of the related articles as it is crucial for the upcoming season.

**Output generated for running word-count on common-crawl data**:



**Data visualization using Tableau:**

https://public.tableau.com/views/lab2_15559316767240/nyt_worsdcouccurence?:embed=y&:display
_count=yes&publish=yes