

# Project Report - Part A

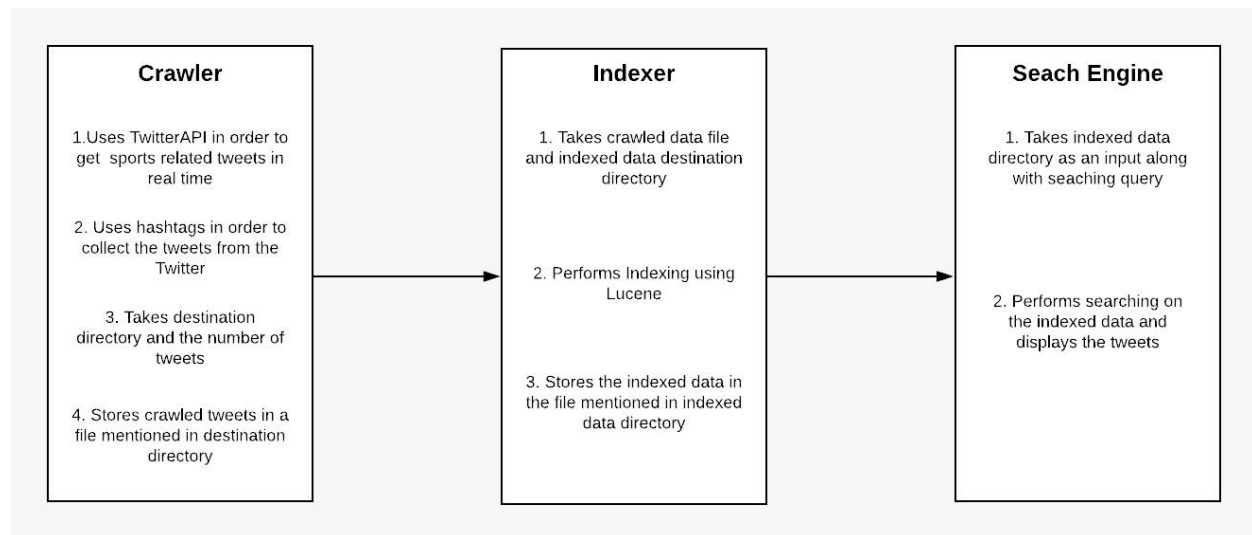
## CS 242 : Information Retrieval and Web Search

Abhishek Ayachit, Gowtham Tumati, Lovepreet Singh Dhaliwal, Sudip Bala, Teja Vemparala  
{aayac001, gtuma002, ldhal001, sbala027, tvemp001}@ucr.edu  
Department of Computer Science, University of California - Riverside

### Project Overview:

Information is key in today's world and the Internet plays a major role in providing adequate amounts of information about anything. Data is available in abundance over the Internet, which can be extracted easily by anyone. People rely on the information available on a large scale. Information could be of any form, articles, pictures, videos, and other well-known forms. Learning from the knowledge extracted on the internet, it can be used to get profitable productivity thereby prospering the exchange.

It is clear enough that we have adequate amounts of information about sports, latest news of events, leagues and all related stuff. Taking the data from Twitter into consideration, we make a search engine for the user, wherein the user can search for any sport-related term and have all the information about it, at hand. Described as a flowchart, is our project procedure(Fig. below)



(Figure: Project Overview)

The project is divided into two phases, which includes Crawling in the first one and then, Indexing.

## **Crawling:**

We worked on this project with an idea of building a sports-based search. In order to implement this, a lot of tweets related to the latest leagues, players and broadcast corporations are required. A question that turned out soon was the method to extract the tweets. Soon enough, we had two options to do that, one of which was the usage of Twitter's REST API, and the other was the usage of Twitter's Streaming API.

The "Search Engine" that we are building requires the latest information regarding everything, the tweets that we extracted must be recent ones. Therefore, we have decided to go forward with Twitter's Streaming API, since Twitter's Rest API is more useful if we wanted to do some analysis or search on historical data, which is not the current case. Streaming API gives us high volumes of live tweets data as per our request until we want to stop receiving.

We have chosen Python to code for tweet crawling procedure. For fetching the live tweets, we have made use of the "tweepy" library. Tweepy library handles connection, authentication and sessions, which makes it easier in using the Streaming API.

We seeded a lot of hashtags related to sports, channels that broadcast content, teams that participate, latest and implemented the crawler. Some of the hashtags included are '#Soccer', '#Formula1', '#MotoGP', '#Cricket', '#ESPN', '#SuperBowl', '#49ers', '#NBA', '#Lakers', '#CR7', '#Badminton'. The biggest challenge till now was to get those tweets faster. Then, we made the use of multiprocessing in our program, as it utilizes all the processors in the given machine and tries to get the data much faster than earlier. The performance has been improving gradually. A size of 50MB was decided in a single file, wherein the incoming tweets are stored. In a similar way, all the tweets are gathered and processed as a batch of files.

## **Indexing:**

After the crawling process is completed, we will need to perform indexing on all the tweets. On learning about the various indexing architectures, we have learnt that the current architecture for a search engine uses the inverted indexes concept. The procedure inside it is that all the indexes are stored as pairs in which the word is key and document ID is the value. In addition to these, there are ranking parameters such as term frequency. Lucene is a text search engine that helps search over huge amounts of data along with the provision of the necessary API's to provide high performance during the process.

The standard analyzer of Lucene is used to prevent stopwords while index generation. Stopwords are not considered as they have a chance to affect the overall ranking of the documents. We have not performed a case sensitive search, where the same word with a different case can be taken as the same word. Indexing has been done only on two attributes as they are the key. The rest of the attributes have been stored for use in the next part.

List of tweet attributes:

| <b>Tweet Attribute</b> | <b>Description</b>         |
|------------------------|----------------------------|
| HashTag                | Keyword for search         |
| Tweet                  | Content of the tweet       |
| Title                  | Title of the URL           |
| createdAt              | The timestamp of the tweet |
| CoOrdinates            | Tweet Location             |
| URL                    | Complete URL of the tweet  |

### **Implementation Procedure:**

Deployment of the entire project:

```
$ bash execute_partA.sh <index-file-directory> <crawled-data-directory>  
<number-of-tweets-to-crawl>  
#eg bash execute_partA.sh /Users/lavidhaliwal/Desktop/cs242/indexed/  
/Users/lavidhaliwal/Desktop/cs242/crawled/ 1000
```

Deployment of Crawler:

```
$ python <file-name> <crawled-data-directory> <number-of-tweets-to-crawl>  
#eg: python tweets_crawl.py /Users/lavidhaliwal/Desktop/cs242/crawled/ 1000
```

Building Lucene Index:

```
$ java -jar <jar-file> <index-file-directory> <crawled-data-directory> 1  
#eg: java -jar tse-0.0.1-SNAPSHOT-jar-with-dependencies.jar  
/Users/lavidhaliwal/Desktop/cs242/indexed/ /Users/lavidhaliwal/Desktop/cs242/crawled/ 1
```

Using the Search Engine:

```
$ java -jar <jar-file> <index-file-directory> 2  
#eg: java -jar tse-0.0.1-SNAPSHOT-jar-with-dependencies.jar  
/Users/lavidhaliwal/Desktop/cs242/indexed/ 2
```

More details regarding this are mentioned in the README file.

## Results:

### 1. Crawling

```
Lovepreets-MacBook-Pro:abc lavidhaliwal$ bash execute_partA.sh indexed/ crawled/ 10
.... Connected to twitter streaming API ....
('Tweet', 1, 'F.size = ', 3534, ' on file:', 0)
Hashtags:[] Tweet:Fuji Xerox Super Cup 2020 Yokohama FM(J-league winner) X Vissel Kobe(Cup winner) Splayer missed penalty kick in a 1 https://t.co/oddne5YubN Coordinates: Date:Sat Feb 08 06:59:24 +0000 20
20 RetweetCount:0 ReplyCount:0 FavoriteCount:0 URL:https://twitter.com/1/web/status/1226037799331057665 Title:Tatau88 Tokyo ENG SalesWeb engineer on Twitter Fuji Xerox Super Cup 2020 Yokohama FM J league W
inner X Vissel Kobe Cup winner Splayer missed penalty kick in a low lol Congrats Vissel Kobe football Iniesta legend record https t co oddne5YubN

('Tweet', 2, 'F.size = ', 4090, ' on file:', 0)
Hashtags:Sports Tweet:RT @PressTV: 7-year-old Iranian #football sensation bends it like Mo Salah #Iran #Sports https://t.co/WTgZhyx3Ip Coordinates: Date:Sat Feb 08 06:59:26 +0000 2020 RetweetCount:0 Repl
yCount:0 FavoriteCount:0 URL:None Title:None

('Tweet', 3, 'F.size = ', 4340, ' on file:', 0)
Hashtags:peace Tweet:You are my peace of mind, in this crazy world. #read #Reading #Kindle #chai #team #funtimes #bed #peace https://t.co/OWfj7hk92C Coordinates:[14.4014883, 50.09109229] Date:Sat Feb 08 0
6:59:26 +0000 2020 RetweetCount:0 ReplyCount:0 FavoriteCount:0 URL:https://twitter.com/i/web/status/1226037805115002880 Title:Pushpendra Pandya on Twitter You are my peace of mind in this crazy world read
#reading kindle chai team funtimes bed peace peaceandlove Prague Castle https t co 12E25Xsyk0

('Tweet', 4, 'F.size = ', 4844, ' on file:', 0)
Hashtags:NBA Tweet:Honestly, they should all be suspended for a month. You can let that incompetence stand. #RipCity #NBA #NBA https://t.co/27JByEGjct Coordinates: Date:Sat Feb 08 06:59:26 +0000 2020 Retw
eetCount:0 ReplyCount:0 FavoriteCount:0 URL:https://twitter.com/i/web/status/1226037806796922881 Title:Chris Barron on Twitter Honestly they should all be suspended for a month You can let that incompeten
ce stand RipCity NBA NBA NBAofficial https t co 27JByEGjct

('Tweet', 5, 'F.size = ', 5313, ' on file:', 0)
Hashtags:StepOnEarth Tweet:RT @TurkishAirlines: Earth is here, waiting to be reached! #StepOnEarth with the airline that flies to the most countries on the planet. Coordinates: Date:Sat Feb 08 06:59:27 +
0000 2020 RetweetCount:0 ReplyCount:0 FavoriteCount:0 URL:None Title:None

('Tweet', 6, 'F.size = ', 5592, ' on file:', 0)
Hashtags:Superstar Tweet:RT @DCblackpink: BLACKPINK x Adidas Superstar 50th Anniversary #BLACKPINK #Adidas #Superstar #adidasoriginals #Originals_kr #KillThisLov Coordinates: Date:Sat Feb 08 06:59:28 +
0000 2020 RetweetCount:0 ReplyCount:0 FavoriteCount:0 URL:None Title:None

('Tweet', 7, 'F.size = ', 5871, ' on file:', 0)
Hashtags:[] Tweet:What is this. Step up your game refs Coordinates: Date:Sat Feb 08 06:59:29 +0000 2020 RetweetCount:0 ReplyCount:0 FavoriteCount:0 URL:None Title:None

('Tweet', 8, 'F.size = ', 6040, ' on file:', 0)
Hashtags:[] Tweet: Coordinates: Date:Sat Feb 08 06:59:29 +0000 2020 RetweetCount:0 ReplyCount:0 FavoriteCount:0 URL:None Title:None

('Tweet', 9, 'F.size = ', 6173, ' on file:', 0)
Hashtags:[] Tweet:RT @daniel86cricket: Good to see many International cricket matches being played in Pakistan. for years people of Pakistan were deprived of Coordinates: Date:Sat Feb 08 06:59:30 +0000 20
20 RetweetCount:0 ReplyCount:0 FavoriteCount:0 URL:None Title:None

('Tweet', 10, 'F.size = ', 6445, ' on file:', 0)
Hashtags:NBA Tweet:RT @SportsPM: You know what time it is... Damian Lillard ties the game! #NBA https://t.co/vSQouJ2AXh Coordinates: Date:Sat Feb 08 06:59:31 +0000 2020 RetweetCount:0 ReplyCount:0 Favo
riteCount:0 URL:None Title:None
```

### 2. Indexing and Searching

```
=====Indexing Started=====
=====Indexing Ended=====
Time Taken = 244ms

===== Twitter Search Engine =====
Search Query =
Sports

===== Query Results =====
Title =
HashTag = #
Tweet = RT @SportsHoochi: 14 https://t.co/UA4r2sQgq4 # #sports #
createdAt = at Feb 08 06:57:46 +0000 2020
URL = https://hoochi.news/articles/20200208-OWTlTS0041.html?utm_source=divr.it&utm_medium=twitter
<== Score and Rank Info ==>
Rank = 1
Score = 1.1305887
=====
Title =
HashTag = #
Tweet = RT @SportsHoochi: 14 https://t.co/UA4r2sQgq4 # #sports #
createdAt = at Feb 08 06:57:46 +0000 2020
URL = https://hoochi.news/articles/20200208-OWTlTS0041.html?utm_source=divr.it&utm_medium=twitter
<== Score and Rank Info ==>
Rank = 2
Score = 1.1305887
=====
Title = None
HashTag = #Sports
Tweet = RT @PressTV: 7-year-old Iranian #football sensation bends it like Mo Salah #Iran #Sports https://t.co/WTgZhyx3Ip
createdAt = at Feb 08 06:59:26 +0000 2020
URL = None
<== Score and Rank Info ==>
Rank = 3
Score = 0.65951
=====
Lovepreets-MacBook-Pro:abc lavidhaliwal$
```

## **Collaboration Details:**

Abhishek Ayachit

- Contributed for Tweet Indexing using Lucene
- Implemented Indexing of tweets using Lucene in Java
- Learnt how Lucene works
- Helped in designing crawler strategies

Gowtham Tumati

- Designed text analyzer choices
- Learnt the working of Twitter Streaming API
- Helped in visualizing the indexed file
- Worked on the project report

Lovepreet Singh Dhaliwal

- Coding for Tweet crawling from the Streaming API
- Researched the basic layout of the search engine
- Implemented Query Search Logic to retrieve top results
- Checked out the optimization procedures for the crawling code

Sudip Bala

- Worked on visualizing the indexed file
- Implemented Multiprocessing
- Helped in Tweet Indexing using Lucene
- Worked on the project report

Teja Vemparala

- Performed necessary analysis and design for Tweet Indexing
- Learnt how Lucene works
- Designed crawler strategies
- Helped in deciding the hashtags to seed

## **References:**

1. <http://lucene.apache.org>
2. <http://www.getopt.org/luke>
3. <https://docs.python.org/2/library/multiprocessing.html>
4. [http://docs.tweepy.org/en/3.7.0/streaming\\_how\\_to.html](http://docs.tweepy.org/en/3.7.0/streaming_how_to.html)
5. <https://developer.twitter.com/>