# Final Project: Twitter Geolocations

Team Members: Edward Zabrensky, Siddharth Menon, Joshua Beto

## Description

Our final project uses the Twitter Stream API to collect Geolocated Tweets. We create an index on these Tweets using Lucene to generate results for search queries. We plot these results on a map using the Google Maps API.

## Collaboration Details

We decided to split up the 3 parts per person. Our collaboration details are as follows:

**Part 1: Twitter Stream** - Joshua Beto
**Part 2: Lucene Index / JSON Parsing** - Siddharth Menon
**Part 3: Geolocation Extension** - Edward Zabrensky

## Part 1: Twitter Stream

Assigned: Joshua Beto

### Architecture:

The architecture is split into a few basic modules. Their functionalities are listed below:

streams module: Handles any logic that deals with the Twitter Stream API (Connection, Reconnecting, Tweets, etc.). This module also adds the 'url_title' field to the Tweet JSON through the map_urls function

listener module: Callback that the streams module calls on when the data is collected. This module allows us to record the number of Tweets streamed and writes the data to its corresponding file.

config module: Handles any configuration data including the Twitter API keys, endpoints, and output file names.

oauth module: Handles any oauth authentication. This module is used by the streams module to connect to the Twitter endpoints.

app module: Handles our main logic and calls our streams module. This module uses user input to determine when to stop streaming.

# Data Collection:

**Threads / Processes:**

The call to the Twitter Streaming API uses a single thread. Collecting roughly 1 GB of data takes ~3 hours on an AMD Ryzen 5 2600 CPU when using a single thread. Adding the url_title field to that JSON data takes another ~3 hours, and uses up to 16 processes in a producer-consumer concurrency model.

The producer is the intermediate data file called tweets.json while the consumer processes the urls found in each JSON of the tweets.json file, visits that page, and parses the 'title' tag of the html received from the get request. It then appends the 'url_title' field with that parsed title to the json data, which it then writes to the file 'tweets_url.data'

**Duplicate Tweets:**

To account for duplicates, we have a check.py script which adds each Tweet from the 'tweets_url.data' file into a set to filter out any duplicates. The script then writes every element of the set back into the file 'tweets_url_no_dup.data'

# Data Structures:

The streaming process is as follows:
- Stream Tweet JSON Data -> Write JSON data to tweets.data
- For each JSON in tweets.data, construct a new JSON object with the 'url_title' field
- Write that new JSON object to tweets_url.data

Our project uses files to store each intermediate step. This is so no data needs to be re-streamed in case of an error when the 'url_title' field is added. For simplicity, we convert back and forth between JSON strings and Python dictionaries to make the implementation easier when adding the 'url_title' field.

One issue we came across is a *MemoryError* when distributing the jobs to the 16 processes to add the 'url_title' field. To counteract this, we split the total JSON into chunks, and passed in a chunk of jobs at a time instead of all at once.

Our producer-consumer concurrency model implicitly uses a queue to administer jobs. To check for duplicate Tweets in our check.py script, we use a set that stores each of the JSON data.

## Limitations:

**Not Automated:** The main drawback to this project is the lack of automation. The application is contingent on the user entering input to determine when Tweet Streaming should end. Another drawback is the lack of automated duplication checking since duplication checking was more of an afterthought than something planned during architecture.

**Lack of Testing:** While the application works in the basic scenario, events that are not accounted for include missing out on stream data during disconnects and more robust disconnection handling.

## Deployment:

**Requirements:** Python 3.6, Pip installed on Python 3.6

To deploy, run the following steps:
- Make sure you are on the Tweets directory
- Run ./install.sh - This will install the Python virtual environment with all the required packages
- Alternatively, you can manually create the virtual environment through the virtualenv package and use the command: pip install -r requirements.txt
- Run ./crawler.sh <CONSUMER_KEY> <CONSUMER_SECRET> <ACCESS_TOKEN> <ACCESS_SECRET>

CONSUMER_KEY = Twitter API Consumer Key
CONSUMER_SECRET = Twitter API Consumer Secret
ACCESS_TOKEN = Twitter API Access Token
ACCESS_SECRET = Twitter API Access Token Secret

You can retrieve these keys by first making a Twitter development account and creating an app.

The process is documented here:
https://developer.twitter.com/en/docs/basics/developer-portal/overview

# Result Screenshots:

# Part 2: Lucene Index

Assigned: Siddharth Menon

## Architecture:

We implemented the indexing of the input twitter data using Java Lucene Core. This section of the project consists of:

The Indexer: this module takes care of taking in the input and building an inverted index, which will make searching based on aggregated points/ranking easier.

The Searcher: this module uses the index built by the indexer to search based on a query and return the ranked results ready to be viewed by the user.

## Index Structures:

The input twitter data is of the json format. The indexer first parses this json data using the Java's org.json library. The parsed data gives us access to different fields in the input data, such as "Time created", "Username", "Tweet text", "Location Coordinates", and "URL".

Our Indexer takes each of these fields and builds a unique Lucene document object. This document object is then added into the associated lucene index.

Once all the data is parsed and added into the index, it is ready to be used by the Searcher to return results.

## Search Algorithm:

The search algorithm uses the previously created inverted lucene index. Using the parsed twitter text data as a way to traverse the index. It is able to take in a user's input and then search the index for relevant data.

The relevant data is returned in terms of the lucene documents (as stored in the indexer structure). The relevant information can be taken out of the documents to display to the user.

## Limitations:

Currently the searcher is only able to return relevant results based on the inherent Lucene ranking. However, it should be possible to gather more information during the indexing process to do some more custom ranking and give more sophisticated relevancy information.

## Deployment:

To deploy this section of the project, it is required to install Java Lucene Core and download the org.json library as well. To test it in your own environment, the paths to essential Lucene .jar files and the org.json .jar file needs to be included in the environment variables. This also requires access to a data file that holds the input data.

# Part 3: Geolocations

Assigned: Edward Zabrensky

## Description:

For the front end portion of this project we used Flask to run our web server. The web server will wait for a user to put in a query. Once a query is requested our Flask web server will call our Lucene Indexer search function that was made in java. This call will search the index that was created and return the top 20 results. These results are then put into a list and sent back to the web server. Once the web server receives the results it will center the map around the user's location based on the browser. Then it will create google maps markers from the locations received from the indexer. The user can click on these markers and it will show the username, tweet, and the title of the link if there is a link.

We also included functions to calculate the tweet center and to calculate the distance of a tweet from the user's location. The calculated the tweet center was created because tweets are given in a bounding box, so to estimate the center of the tweet we created this function. The calculation of the distance from a tweet to the user's location was created so that we can limit the amount of tweets shown. However, since the number of tweets is only 1GB we decided to show all the tweets anywhere in the world.
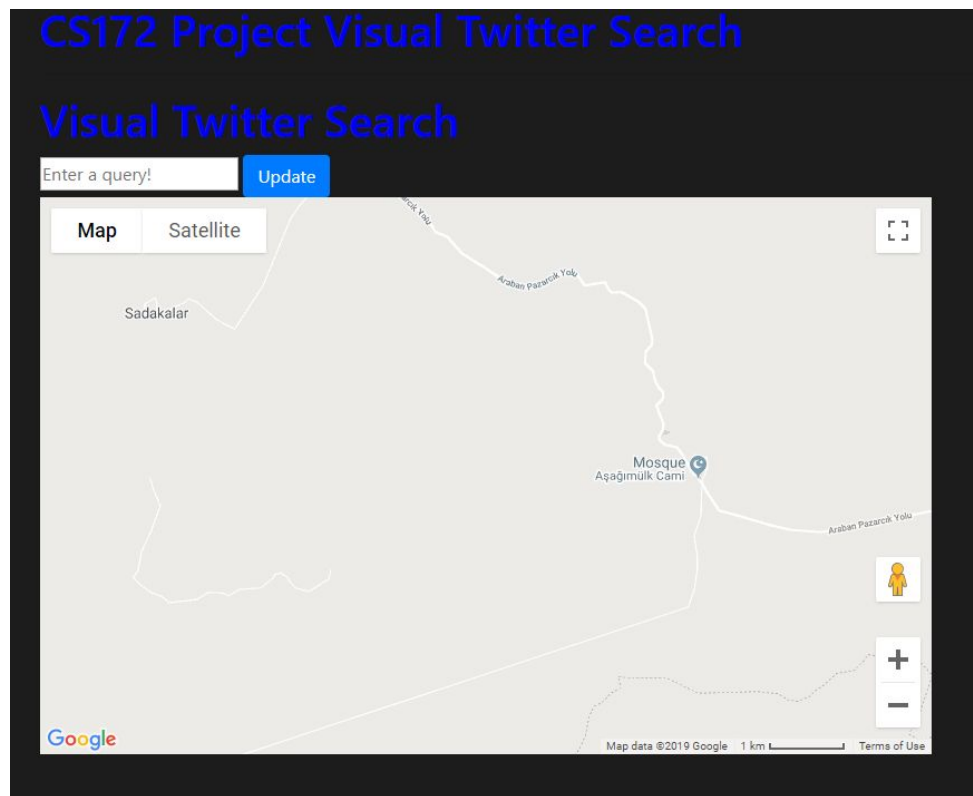
## Limitations:

One limitation was that our indexer was created in java while our web server was created in python. To resolve this we used jython which allowed our python web server to call methods in our Java. Also to make this work we needed to create methods in the Lucene Indexer file that returned python objects instead of java ones so that we could call these methods from our python script. We also had to update the jython file itself to include the classpaths of our java dependencies because when our web server would call our Lucene Indexer more than once it would only include the classpaths for the first call to Lucene Indexer.

Another limitation was the emojis. When we tried sending the text of the tweet as a python string to the server it would always crash. We figured out that we needed to encode the text and url title of tweets as python unicode so that it would know how to encode these pieces of data.

## Deployment:

To deploy this front end we need google maps api key, flask, jython 2.7.1 and java's org.python.core dependency. More about deployment is on the frontend section of the readme.
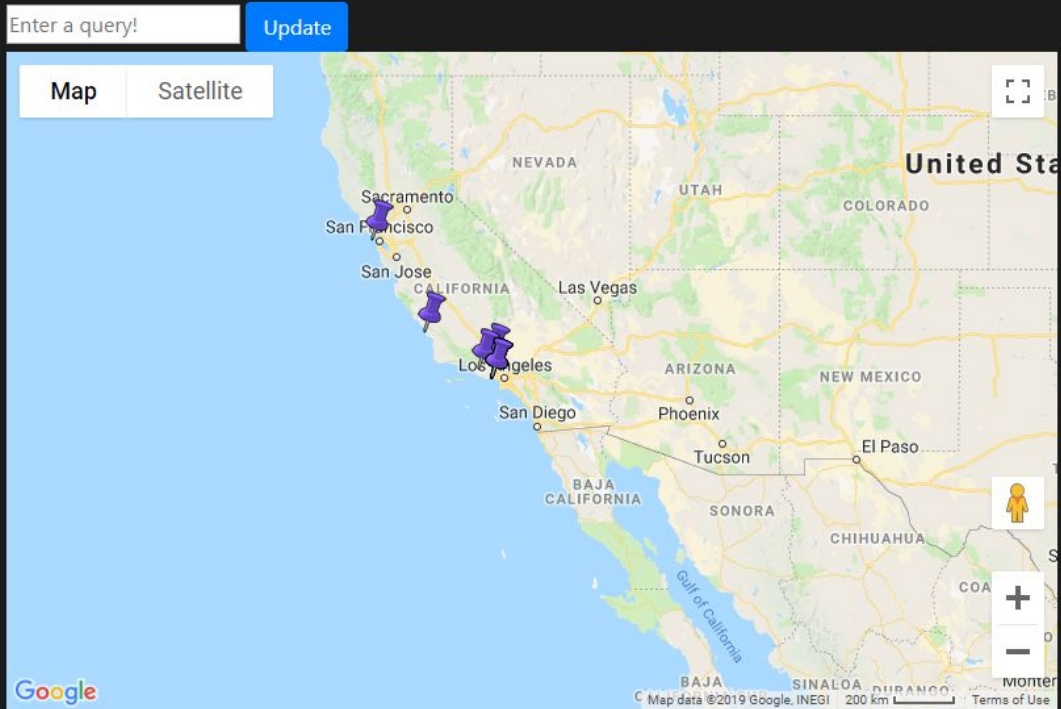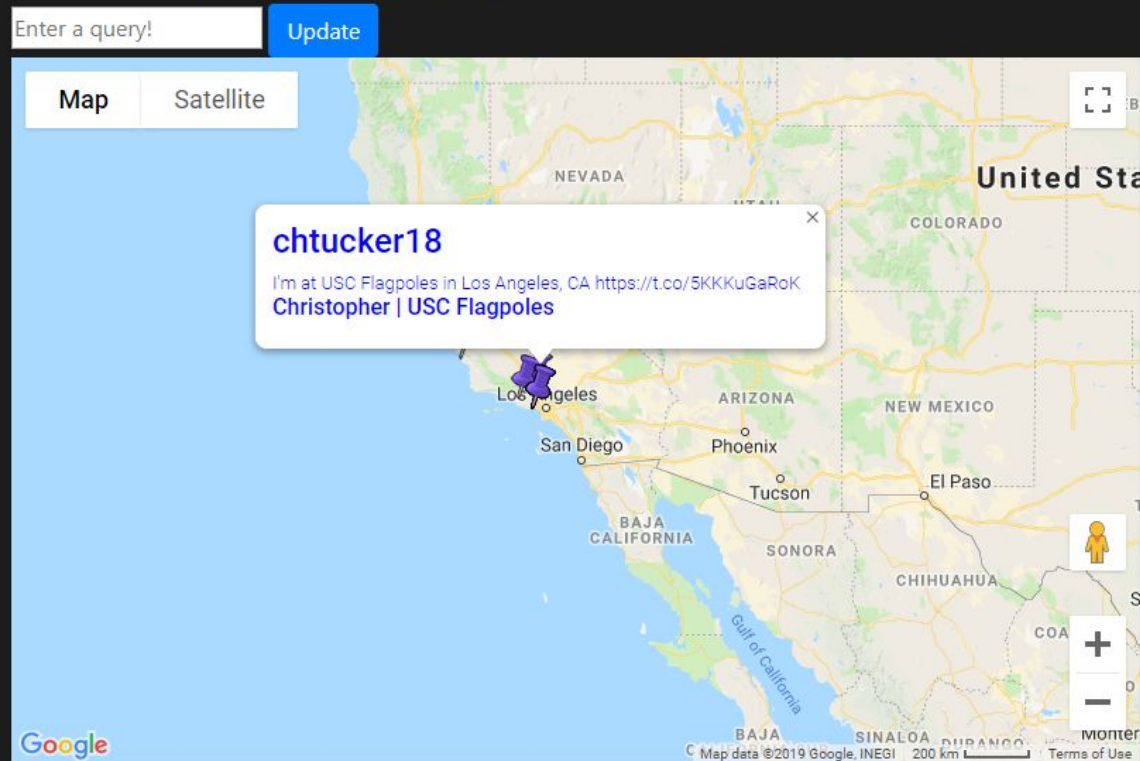
# Result Screenshots:



Search for 'Angeles'

After clicking on a pin

Shows the username, the tweet, and the title of the link if a link is available.