# TwittMap Report

Team members: Xiaolong Jiang (N18261025), Guoqing Xiong (N17801169)

## 1. Twitter API

In the Twitter Streaming API part, we applied a account to configure Key and Secret. We used twitter4j package and call TwitterStreamFactory class to get instance of twitterStream.

Then we call the listener in twitterStream. We used Filter Query to search word and override onStatus method to get elements such ad id,coordinates in JSON file. Then we upload these JSON file to cloudSearch platform.

```java
public class twitterTest {


    private static String consumerKey = "BzBCWLBCup4VRJ7EDDwPmixRH",        keys and tokens
            consumerSeceret = "lvkwF356cZpRrKqYToPwlisVind5UxXP2g07P1rkjSIjwM5C0U",
            token = "705579364973469696-OVY1yiCFYCK27n5plvlc3kqZVVT5N4L",
            tokenSecret = "gxTZWC1qkkI0iqSmn0z7YtJqGKnTZib6TuLUBFPPD0Uiw";
    public static void main(String[] args) {
        ConfigurationBuilder cb = new ConfigurationBuilder();
        cb.setDebugEnabled(true)
            .setOAuthConsumerKey(consumerKey)
            .setOAuthConsumerSecret(consumerSeceret)
            .setOAuthAccessToken(token)
            .setOAuthAccessTokenSecret(tokenSecret);

        String[] keyWords = {"oscar","award"};//the relationship between the keywords is or
        FilterQuery fq = new FilterQuery();
        fq.track(keyWords);
        TwitterStream twitterStream = new TwitterStreamFactory(cb.build()).getInstance();
                try {

                StatusListener listener = new StatusListener() {
                    @Override
                    public void onStatus(Status status) {
                        try {
                        if (status.getGeoLocation() != null && status.getLang().equalsIgnoreCase("en")) {
Get attributes we         String createAt=status.getCreatedAt().toString();
need, i.e. created at,    long idStr=status.getId();
text, geoLocation.        String text = status.getText();
                            GeoLocation coordinates = status.getGeoLocation();

                            JSONObject json = new JSONObject();

                            json.put("id_str", idStr);
                            json.put("created_at", createAt);
                            json.put("text", text);
                            json.put("coordinates", coordinates.toString());
                            String index = json.toString();
                            System.out.println(index);
                        }
                } catch (Exception e) {
                    e.printStackTrace();
                }
            }
```
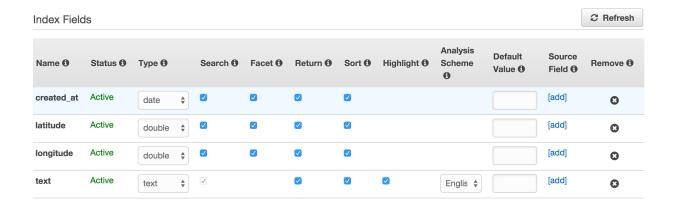
# 2. Deploy AWS CloudSearch

Amazon CloudSearch is a managed service in the AWS Cloud that makes it simple and cost-effective to set up, manage, and scale a search solution for your website or application.

## 2.1 Index Construction

To search files/informations in structured or semi-structured data, we have to build index of files first. In each tweet, the indices we need for now are created_at, latitude, longitude, and text. Besides, tweed ID is needed to uniquely identify each tweet file in CloudSearch.

Figure below is the index fields of CloudSearch.

Index Fields    ⟳ Refresh

| Name ❶ | Status ❶ | Type ❶ | Search ❶ | Facet ❶ | Return ❶ | Sort ❶ | Highlight ❶ | Analysis Scheme ❶ | Default Value ❶ | Source Field ❶ | Remove ❶ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| created_at | Active | date ⏶⏷ | ☑ | ☑ | ☑ | ☑ | | | | [add] | ✖ |
| latitude | Active | double ⏶⏷ | ☑ | ☑ | ☑ | ☑ | | | | [add] | ✖ |
| longitude | Active | double ⏶⏷ | ☑ | ☑ | ☑ | ☑ | | | | [add] | ✖ |
| text | Active | text ⏶⏷ | ☑ | | ☑ | ☑ | ☑ | Englis ⏶⏷ | | [add] | ✖ |

## 2.2 Upload File

Because the file standard formats of Twitter API and AWS CloudSearch don't match. We need a Tweet Format Transfer program, which will be mention at Chapter 3.

## 2.3 Search

After upload file onto AWS CloudSearch, we can search through AWS CloudSearch Console or HTTP Request.

### 2.3.1 Search through AWS CloudSearch Console

Click on 'Run a Test Search' in the navigation bar, and type in search query, then click 'GO' button. The relevant file will be returned by CloudSearch.

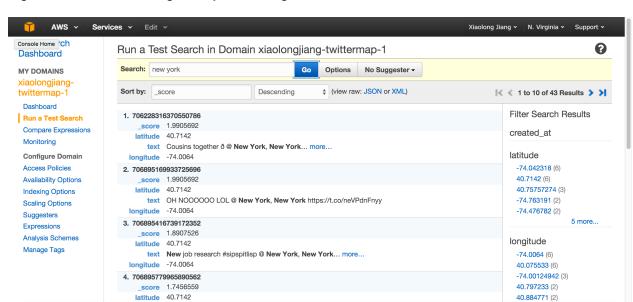Figure below is searching 'new york' through CS console.



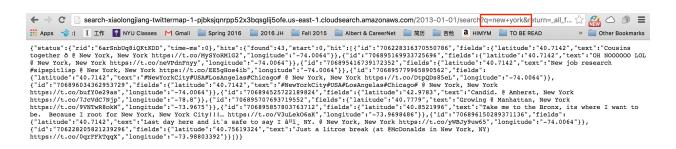### 2.3.2 Search through HTTP Request

Using <"http://" + Search Endpoint + "search?q=" + query> to get response from CloudSearch. It supports multiple words query with "+" between words instead of whitespace.

Figure below is searching "new york" using http request.



# 3. Data Processor from Twitter API JSON File to AWS CloudSearch JSON File

Since file format returned from Twitter API is different from upload format in AWS CloudSearch. We need to write a program to transfer Tweets format. The program reads original tweet file from disk, parses original tweet file in each index and put indices into standard format of CloudSearch, and write the result to disk.

Moreover, S3 can be used in future to process data on the cloud, EC2 can be used in future to run parse program, and the data output from EC2 could upload to CloudSearch directly.

Figure below is main program of Tweet Format Transfer Professor.

```java
public class TweetFormatTransferProcessor {
    public static void main(String[] args) {
        String fileSourceLoc = "/Users/xiaolongjiang/Desktop/cloud_computing/TwitterMap/";
        String fileSourceName = "twitterAPIResult_Mar7_query_Oscar.json";

        String fileDestLoc = "/Users/xiaolongjiang/Desktop/cloud_computing/TwitterMap/";
        String fileDestName = "twitterAPIResult_Mar7_query_Oscar_for_CloudSearch";

        TweetReader tweetreader = new TweetReader(fileSourceLoc, fileSourceName);
        tweetreader.readFromDisk();
        String[] tweets = tweetreader.getTweets();

        TweetsFormatProcessor tweetsProcessor = new TweetsFormatProcessor();
        String formattedTweets = tweetsProcessor.processTweets(tweets);

        TwitterWriter tweetWriter = new TwitterWriter();
        tweetWriter.writeToDisk(fileDestLoc, fileDestName, formattedTweets);
    }
}
```
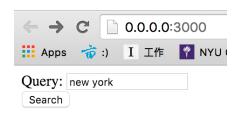
Figure below is Tweets returned from Twitter API



Figure below is Tweet file after ProcessTwitterFile Program.

# 4. Web Interface Using Ruby on Rails

## 4.1 User Interface

We are using Ruby on Rails to build the UI of twitter map. The UI is simple and clear— typing query into text filed of 'query', and then click the 'search' button. Then a web page of relevant tweets, and a map of tweets' text content and location are.

Figure below is searching User Interface. User could type query into textfield and then click search.

## 4.2 Send Request to CloudSearch and Get Result

At the UI, we get User's query input. First, we store the query into database. Then concatenate query string with CloudSearch Endpoint. After that, we can send a HTTP Request to CloudSearch and get response.

Figure below is ruby code of Request Controller.

```ruby
class AsksController < ApplicationController
  def new
    @ask = Ask.new
  end

  def create
    require "net/http"
    require "uri"
    require "json"

    @ask = Ask.new(ask_params)
    link_to_cloudsearch = "http://search-xiaolongjiang-twittermap-1
    uri = URI.parse(link_to_cloudsearch)

    # Get Response of HTTP Request
    @response_cloud_search = Net::HTTP.get_response(uri)

    # Parse Request into JSON file.
    @parsed = JSON.parse @response_cloud_search.body

    if @ask.save
      render 'cloud_search/request_to_cloudsearch'
      # redirect_to :back
    end
  end

  private

  def ask_params
    params.require(:ask).permit(:query)
  end
end
```

# 5. Google Map API and Rails

We are using Javascript to access Google Map API. The UI get query from user, and then sends query http request to CloudSearch, the returned JSON file is parsed to a Hash data structure by ruby. Then the tweet text and geological information is passed to Google Map API. We initialized Google map object and configured its view center, zoom degree and mapType. Then we inserted tweet makers into map by its coordinates and inserted infoWindow into map by its text. We do it to all tweet data set in a for loop. At last, we call map's event to addDomListener to load all these data.

Since Rails is embedded with HTML and Javascript, we can put HTML/Ruby/Javascript code together.

# 6. Deploy Rails Application on Elastic Beanstalk

To deploy Rails web application, we need to:

(1)Setup git repository
$ git init
$ git add -A
$ git commit -m "default rails project"

(2) Configure the EB CLI
$ eb init
$ git commit -am "updated .gitignore"

(3) add a gem 'puma' into Gemfile
$ git commit -am "Add Puma to Gemfile"

(4) Deploy the project
$ eb create rails-beanstalk-env

Figure below is Rails application successfully deployed on Elastic Beanstalk
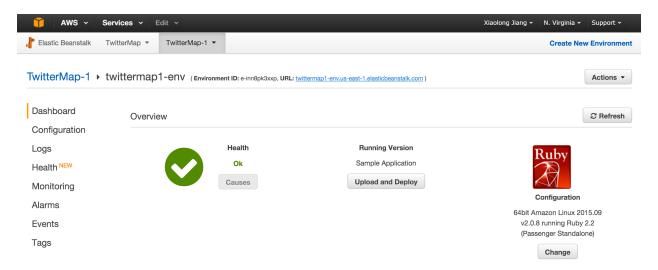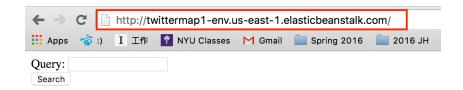


Figure below is accessing Rails application through EB

## 7. Refinement

In this assignment, we could implement S3 and EC2 into Twitter crawling process and uploading files onto CloudSearch process. If we can do crawling and update in a time period, say 1 hour per update, then real-time twitter information crawl/search platform can be established purely on AWS cloud.

Moreover, we're crawling twitter in a single thread program. If we can build a multi-threads crawling program, we can do much faster than single thread since the processor wouldn't be idle during the response waiting time.

## 8. Conclusion

AWS is very useful and powerful platform. We can build projects based on the services AWS provided. In this assignment, we used AWS CloudSearch and AWS Beanstalk in this project. AWS CloudSearch is more like a powerful search engine for structured and semi-structured file, and Beanstalk can deploy Web application in a second. These services could save a lot time and money for startup even mid-sized companies.