# Project Report

Topic: Tweet Clustering System
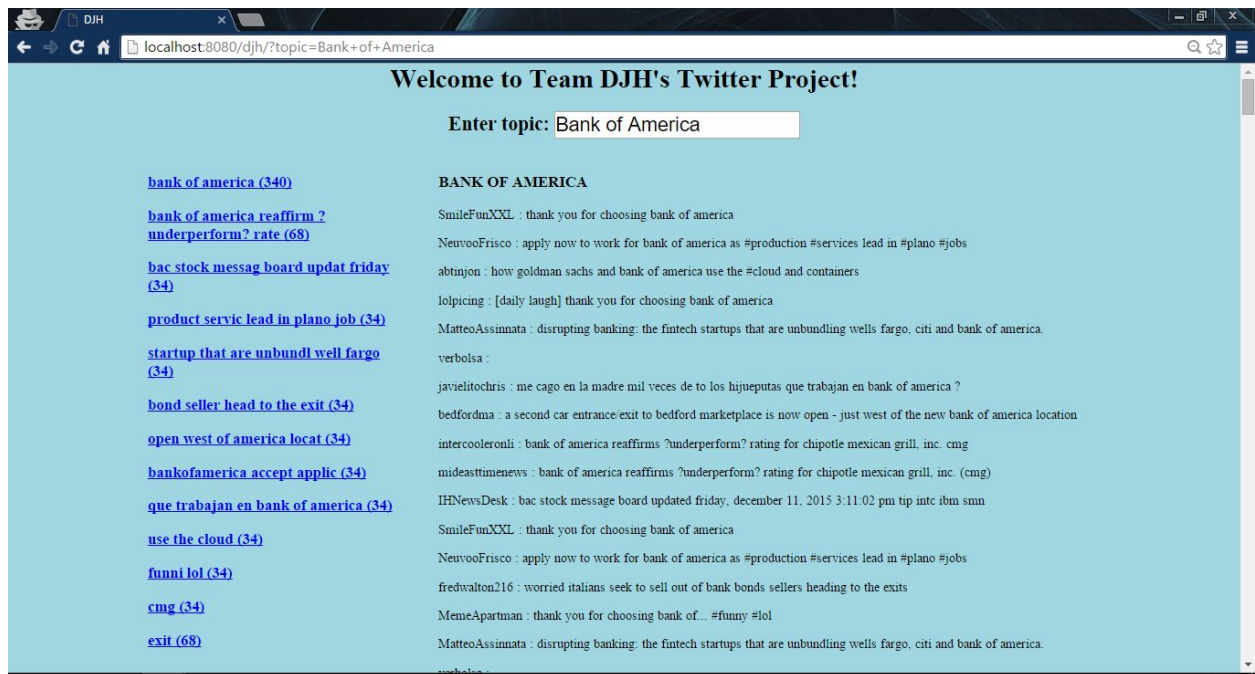
Link: https://github.com/TeamDJH/DJH

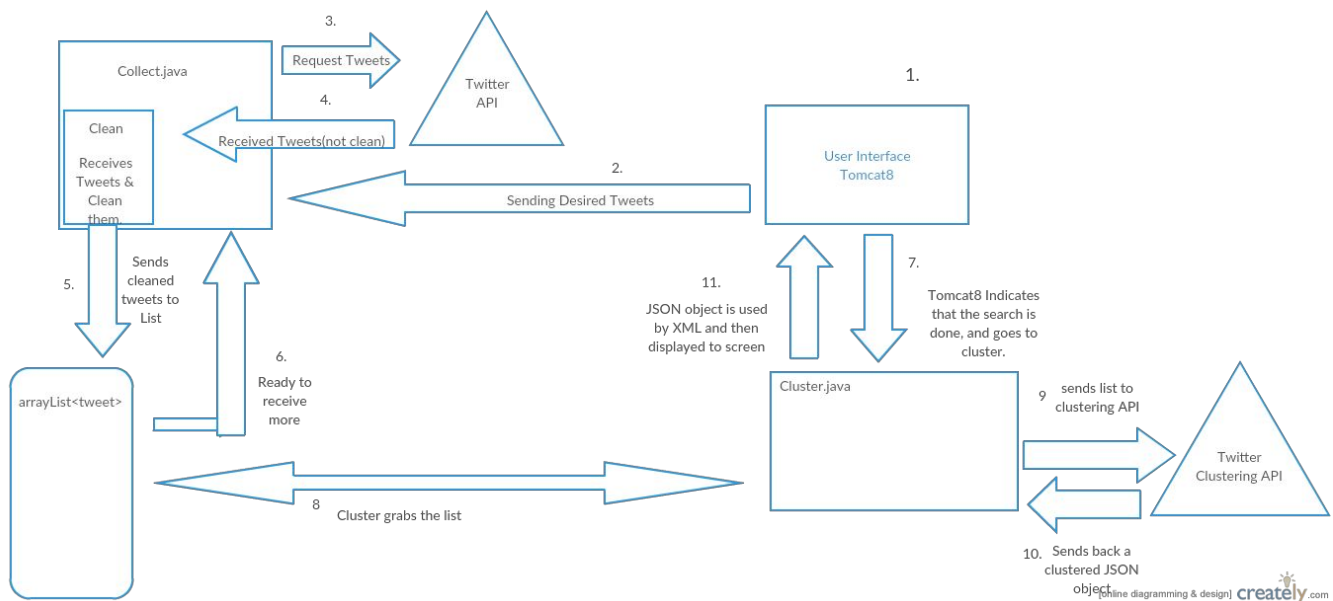Team DJH Members: David Gwin, Jacob Reid, Helam Franco



**Overview**:

This program allows you to search for tweets relating to a specific topic. It accesses the Twitter database and returns the results clustered into relevant topics. In the example below we searched for Bank of America.

**Program Architecture:**

We split our project into three main classes: Collect, Cluster and Tweet. The Tweet class we wrote to be able to use tweet objects. In the Collect class we connected to the Twitter API using the Twitter4j library, and collected and cleaned the tweets. The Cluster class takes the cleaned tweets and clustered them using the RxNLP's Sentence Clustering API. We then built a servlet to give our program a user interface using Apache TomCat.



**Modules used:**

Apache TomCat: used to make servlet to be able to make a user interface with a front-end and back-end.

GitHub: used to manage code changes and keep everyone's code up to date.

Slack: used to manage communication between team members.

RxNLP Sentence Clustering API: used to cluster tweet text into topics.

Twitter Streaming API: used to collect tweets from Twitter data base.

Twitter4j: used to make using Twitter Streaming API easier to use.

Mashape: used to get a key to access the RxNLP Sentence Clustering API.

XML: used to display clustered tweets as per instructors request.

JSON: used to pass data to the RxNLP Sentence Clustering and also was the format that the

data was in upon return from the RxNLP Sentence Clustering.

HTML: used to format and display the clustered tweets.

CSS: used to style the display of the clustered tweets.

JSP: used to call and run the program and to get the results and pass it to the JavaScript and

HTML.

JAVA: used to write the core of the program; Collect, Cluster and Tweet classes.

JavaScript: used to show selected topic in the user interface.

**Classes:**

Collect class accesses the Twitter streaming API using Twitter4j. We registered our application with Twitter and received the proper tokens and keys. The collect class takes a string parameter of topic. it passes that topic to the Twitter streaming API through Twitter4j. It then loops through requesting tweets till it reaches 500. Each tweet that is read is stored in a an array list of tweets. The list is then passed to the clean function to remove duplicates and clean up the text and beep out profanity.

Tweet is a class we write to make a tweet object. Our tweet object has two String fields: 'username' and 'text'.

Cluster is passed the ArrayList of tweets from the Collect class. It takes objects from the arraylist and formats them into a JSONObject and sends it to the clustering API. The clustering API returns a JSONObject of all the cluster topics and clusters. The cluster class then maps the clustered tweets back to the original tweets in the arraylist and adds the username to the text. We found a function online that then formats the JSONObject into XML to be displayed in the browser.(http://stackoverflow.com/questions/25864316/pretty-print-xml-in-java-8/33541820#33541820)

**Extra Credit:**

We were able to implement a profanity filter that blocks a lot of those swear words. We also used TomCat make a servlet to be able to make a user interface with a front-end and back-end. We used JSP, HTML, CSS and JavaScript to display the user interface.

**Division of Labor:**

Jacob: Made Cluster and Tweet classes and setup RxNLP's Sentence Clustering. Setup servlet to display frontend using TomCat, JSP, JavaScript, HTML and CSS. Created mashape account to get API key. Also found a XML formatter function to turn the JSONObject into XML. (http://stackoverflow.com/questions/25864316/pretty-print-XML-in-java-8/33541820#33541820) Used CSS and javascript to display different tweets based on topics along with showing the XML upon click. Set up Git repository on GitHub for code collaboration. Figured out how to import all the external libraries via adding the jar files to build path. Wrote documentation and setup instructions using GitHub Wiki. Set up Slack communication tool to improve communication. Kept everyone informed and on the same page by communicating with team members often.

Helam: Implemented the cleaning part of the Twitter program. This included removing duplicate tweets, lowercase text, remove unwanted text (such as @,!,hyperlinks, etc). This was achieved by using regular expressions. Regex was a challenge to implement particularly the URL Regex and figuring out how to remove them from any string. The initial approach to solve the Regex url problem was making a huge 'if, else' statement that would convert the string to a char array, and search for the http, but after some help from other friends, it was finalized that this regex worked on a URL :

"(https?|http?|ftp|file)://[a-zA-Z0-9+&@#/%?=~_!:,.;]*[-a-zA-Z0-9+&@#/%=~_|]"

Allowed us to remove the substring of the url and everything after it until a whitespace.  Also implemented a profanity filter.

David: Wrote getProperties, getAccess and writeProperties classes in the collect class. Created Twitter account, registered the application with Twitter. Read through the Twitter API. Helped implement the Twitter4j. Maintained communication though slack communication tool.

**Issues:**

We struggled adding the jar files for libraries we needed to use into our project. We also struggled working on the same project together, we fixed that problem by using GitHub to manage our code changes. We struggled with communication but we fixed that by setting up Slack communication tool. We also struggled because a member of our group (Jorge) dropped the class and did not tell us. He had been assigned to work on the clustering class and we thought he was working on it but we found out a week before the presentation that he had dropped the class and had not done any work on the project. This sent us into panic mode as we tried to catch up and get the project together.