

# CS172 Project: Index and Web Search

## Project Report

Team Members - Jackie Chan (861158919)

Abe Hu (861148832)

Daniel Nguyen (861156246)

### Collaboration Details:

Jackie Chan:

- Implemented indexing and search functions

Abe Hu:

- Created the backend Server endpoints
- Created frontend user interface

Daniel Nguyen:

- Implemented indexing and search functions
- Complete 4-5 page report

### Overall Collaboration Splits:

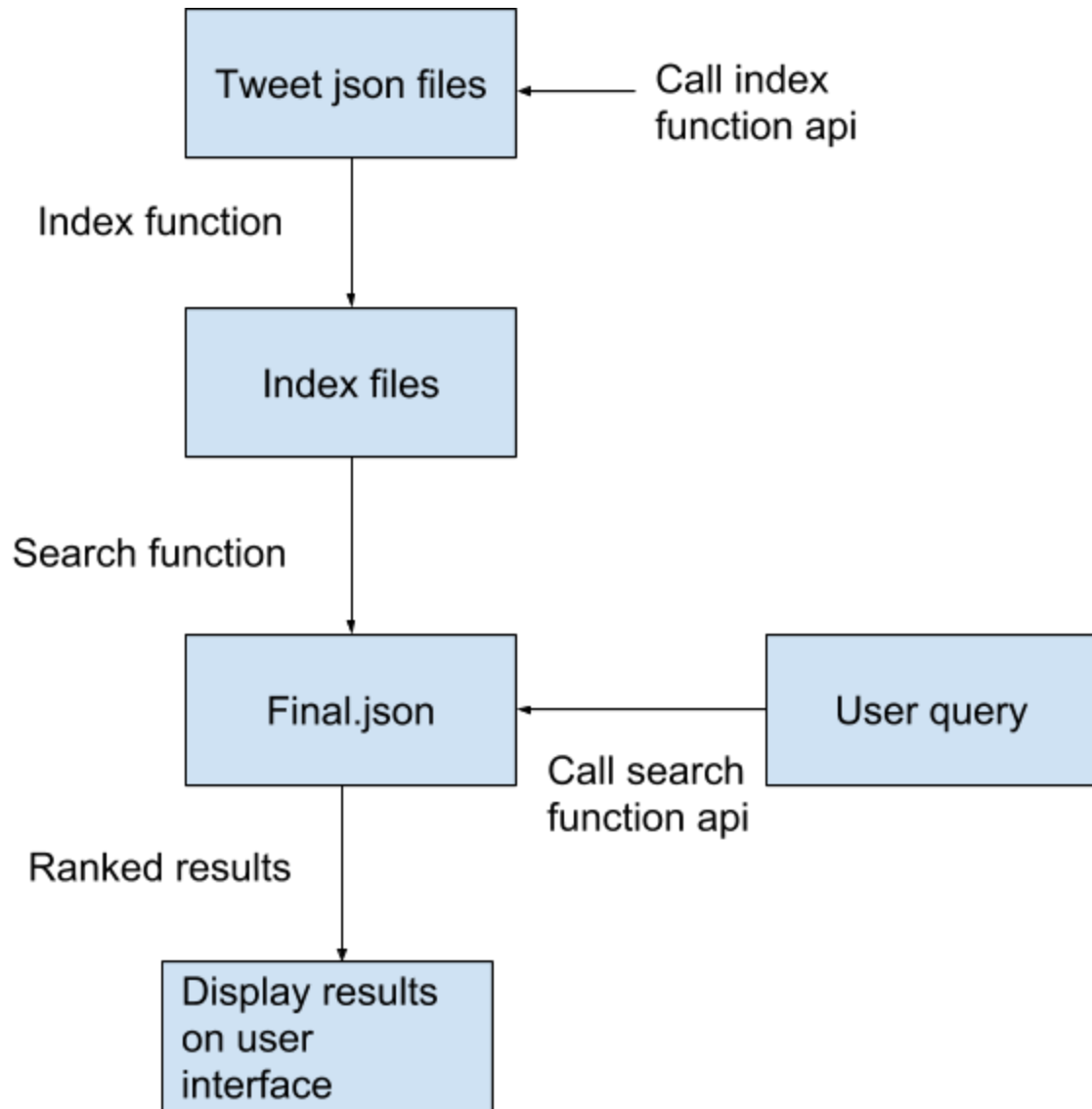
Jackie Chan - 33.33%

Abe Hu - 33.33%

Daniel Nguyen - 33.33%

## Overview of systems:

### a. Architecture:



b. Data indexing and searching:

We used pylucene to iterate through every json file in the data directory and index them. The fields that we included in the index were username, text, time, and location. For the search, we used the standard analyzer, to do a very simple search. After inputting a query, the function will return the top 10 results.

c. Data Structures and Architecture Explained:

For the frontend user interface we chose to use react and redux. Bootstrap was used for styling the DOM elements.

The server is designed to run continuously. We created two endpoints: one for indexing, and one for searching. The index route is called automatically every time the server runs. It will only need to be run once since the index files will be saved to storage after the first call. Entering a query and pressing the submit button will take the search endpoint and run the function.

The index function grabs each line and writes it to a document with specific fields. The new created files will be stored in a directory called `lucene_index`.

The search function scores each element based on the text's relevance to the query. Once it has completed, it will return the top ten results as a json object to the front end. The client side user interface will then display the user, text, time, location, and other information.

### **Limitations of Systems:**

- We did not implement stemming
- Limited in the user relevance of results, due to the fact that we only used the standard analyzer
- Did not implement the extra credit gps coordinate portion

### **Instructions for system deployment:**

1. Install python
2. Install pip (python package manager)
3. Install all python packages
  - a. pylucene
  - b. Flask
  - c. Requests
  - d. React
  - e. Redux
4. To start the backend
  - a. Npm start
  - b. Run React app
  - c. Go to src/public directory
  - d. To run flask input - `FLASK_APP=server.py flask run`
  - e. Input query

## Index function

```
Number of indexed documents: 232592
Finished
```

```

Enter a blank line to quit.
Searching for Information retrieval
10 total matching documents in 0:00:00.002219:

{'username': 'u:sarasalimi', 'text': 'Information not given accurately and completely.', 'score': 11.135669708251953, 'time': 'u'Thu Jun 07 02:53:27 +0000 2018'}

{'username': 'u:disasteraware', 'text': 'u'New Alert: Tsunami Information (Hawaiian Island Off The Hamakua Coast of The Big Island - 4.1, Severity: INFORMATION', 'score': 9.9610932922363, 'time': 'u'Wed Jun 06 20:40:22 +0000 2018'}

{'username': 'u:andrewshaylie', 'text': 'u'If anyone would like to give blood, HMU for the information !\u2764ufe0f\u2764U01f489\u2764ufe0f', 'score': 9.577493667602539, 'time': 'u'hu Jun 07 03:06:08 +0000 2018'}

{'username': 'u:LazyEyyez', 'text': 'u'@Solo_Kalin Aye, u not telln no false information \', 'score': 9.577493667602539, 'time': 'u'Thu Jun 07 03:49:02 +0000 2018'}

{'username': 'u:wirtzbill', 'text': 'u'@icjklr All the information is in the article you dn't read.', 'score': 9.577493667602539, 'time': 'u'Thu Jun 07 03:57:53 +0000 2018'}

{'username': 'u:dipeshkhadye', 'text': 'u'@vtasindia2018 new data classification - Veritas Information Classifier', 'score': 9.577493667602539, 'time': 'u'Thu Jun 07 05:57:11 +0000 18'}

{'username': 'u:JohnDioriod4', 'text': 'u'@renato_mariotti @realDonaldTrump I suppose you ve access to government information', 'score': 8.39596176147461, 'time': 'u'Wed Jun 06 21 9:52 +0000 2018'}

{'username': 'u:richrad', 'text': 'u'there is some information perhaps !\u2019d rather not have https://t.co/vJg1FQVZ', 'score': 8.39596176147461, 'time': 'u'Thu Jun 07 04:06:45 +00 2018'}

{'username': 'u:StormMae', 'text': 'u'@CoachT_SFU For your information, I got salted caral. https://t.co/LWAARFcW9d', 'score': 8.39596176147461, 'time': 'u'Thu Jun 07 05:02:36 +0 2018'}

{'username': 'u:SEASHORE46', 'text': 'u'@markgstrong @realDonaldTrumpFan Where did you get yo information from playboy', 'score': 8.39596176147461, 'time': 'u'Wed Jun 06 22:13:53 +00 2018'}

```