

Rogelio Macedo

Instructor: Mariam Salloum

Alexander Yee

Team: LightSABR

Shiyao Feng

Brett McCausland

18-March-2019

Search Engine

In completing this project, I consulted:

- Mariam Salloum during lectures/office hours
- Class channel on Campuswire

All the important code is original. Unimportant subroutines that are not completely original are sourced:

- <https://sundog-education.com/elasticsearch/>
- <https://www.udemy.com/elasticsearch-6-and-elastic-stack-in-depth-and-hands-on/>
- <https://www.elastic.co/blog/how-to-find-and-remove-duplicate-documents-in-elasticsearch>
- <https://www.elastic.co/guide/en/elasticsearch/client/javascript-api/current/api-reference.html>
- <https://www.elastic.co/guide/en/elasticsearch/reference/current/query-dsl-match-all-query.html>
- http://exploringelasticsearch.com/searching_data.html#ch-searching-data
- <http://okfnlabs.org/blog/2013/07/01/elasticsearch-query-tutorial.html#query-dsl-overview>
- <https://www.youtube.com/watch?v=52G5ZzE0XpY#t=1471>
- <https://sundog-education.com/elasticsearch/>
- <https://towardsdatascience.com/getting-started-with-elasticsearch-in-python-c3598e718380>
- https://www.elastic.co/guide/en/elasticsearch/guide/current/_talking_to_elasticsearch.html
- <https://www.elastic.co/guide/en/elasticsearch/client/index.html>
- <https://stackoverflow.com/questions/24153996/is-there-a-limit-on-the-size-of-a-string-in-json-with-node-js>
- <https://www.elastic.co/guide/en/elasticsearch/reference/current/docs-bulk.html>
- <https://www.elastic.co/guide/en/elasticsearch/reference/current/indices-delete-index.html>
- https://www.elastic.co/guide/en/elasticsearch/reference/6.1/_list_all_indices.html
- [Tutorial on Webscriping and parsing with BeautifulSoup](#)
- <https://stackoverflow.com/questions/1936466/beautifulsoup-grab-visible-webpage-text>
- <https://chartio.com/resources/tutorials/how-to-install-elasticsearch-on-mac-os-x/>

- <https://stackoverflow.com/questions/12451997/beautifulsoup-gettext-from-between-p-not-picking-up-subsequent-paragraphs>
- <https://stackoverflow.com/questions/50045253/if-statement-in-ejs>
- <https://www.elastic.co/guide/en/elasticsearch/reference/current/query-dsl-match-query.html>
- <https://stackoverflow.com/questions/10326950/render-a-variable-as-html-in-ejs>
- <https://stackoverflow.com/questions/50103466/ejs-cannot-read-property-of-undefined-how-can-i-fix-this>
- <https://code.tutsplus.com/articles/introduction-to-parallel-and-concurrent-programming-in-python--cms-28612>
- <https://docs.python.org/2/library/threading.html>
- <https://www.oreilly.com/learning/python-cookbook-concurrency>
- <https://youtu.be/gnsO8-xJ8rs>
- <https://getbootstrap.com>

Introduction

The search engine is made up of three main components: the web crawler, the indexer and the user interface. There is also an intermediate phase in which we create json objects from each of html files crawled in order to load the index. This external phase is crucial for combining the search engine's components.

The web crawler crawls a set of seed .edu urls and extracts any link with the .edu top level domain. We then iterate through all the crawled html files, parsing each one for certain fields and create its instance as a json object. This intermediate phase results in a json file containing all of the web-document objects. This file is then bulk loaded onto the Elasticsearch cluster through an HTTP protocol PUT/POST request.

The interface offers a search bar and a submit button that returns query results based on the title and paragraph fields of the html document. Our json objects have collected paragraph text that is all compiled into one "body" field. Elasticsearch returns the results based on multi-field matching, displaying the matched web-documents and italicized highlights from "title" and "body" matches.

Web Crawler

(a) Architecture

- Multi-threaded design
- Duplicate detection

(b) The Crawling or Data Collection Strategy

- Using multiple threading to get the URL from seed web, and extract links from it to other docs (URLs) without duplicates by concurrent.futures and regular expression.
- Using multiple threading to download each HTML file to a local folder by 16 multi-task.
- In the data folder, the application creates a different level folder to keep the files.
- The application allows user input the number of pages to crawl and number of levels to limit application.

Final Project Report

CS 172: Information Retrieval

- The default setting of application: if the application gets the file over 1GB in some arbitrary level K, the application will stop when it finished web crawling in level K.
- If the user wants to input the limit by page or level, the memory wouldn't be limit by 1 GB
- When web crawling reaches the max level, the application will stop until the job is at max level.
- When web crawling reaches the max pages, the application will stop at $\text{Max} \sim \text{Max} + 15$.

```
fengdeMacBook-Pro:finalproject lightsabr fengs$ python multi-threading.py

-----
Welcome to LightSABR's web crawling
Default: : if application get the file over 1GB in level K, the application will stop when it finished web crawling in level K
Do you want to limit the web crawl by page or level ?
1 For yes OR 0 For No : 0
total memory now is: 2M
Now level is: 2
total memory now is: 4M
Now level is: 3
total memory now is: 10M
Now level is: 4
total memory now is: 17M
Now level is: 5
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 37M
Now level is: 6
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 75M
Now level is: 7
total memory now is: 106M
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 132M
Now level is: 8
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 158M
Now level is: 9
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 192M
Now level is: 10
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 220M
Now level is: 11
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 248M
Now level is: 12
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 277M
Now level is: 13
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 306M
Now level is: 14
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 369M
Now level is: 15
total memory now is: 438M
Now level is: 16
total memory now is: 529M
Now level is: 19
total memory now is: 576M
Now level is: 20
total memory now is: 606M
Now level is: 21
total memory now is: 636M
Now level is: 22
total memory now is: 666M
Now level is: 23
total memory now is: 692M
Now level is: 24
total memory now is: 724M
Now level is: 25
total memory now is: 760M
Now level is: 26
total memory now is: 792M
Now level is: 27
Some characters could not be decoded, and were replaced with REPLACEMENT CHARACTER.
total memory now is: 826M
Now level is: 28
total memory now is: 857M
Now level is: 29
total memory now is: 885M
Now level is: 30
total memory now is: 908M
Now level is: 31
total memory now is: 933M
Now level is: 32
total memory now is: 956M
Now level is: 33
total memory now is: 983M
Now level is: 34
LightSABR web crawling reach memory: 1031376458
LightSABR web crawling reach page: 17628
LightSABR web crawling reach level: 34
LightSABR web crawling spend : 1502.6790297031403 sec
fengdeMacBook-Pro:finalproject-lightsabr fengs$
```

Figure 1 & 2: show the result of the web crawling application default setting.

(c) Data Structures/libraries employed

- Python lists, dictionaries, queues, time, hurry filesize(file size), re(regular expression)
- The BeautifulSoup library
- The [urllib.request](#) module

(d) Limitations

- the multiprocessing of python is hard to control the shared variable - page_number, so the application will stop at $\text{Max} \sim \text{Max} + 15$. The reason why the application will stop at $\text{Max} + 15$ is that I set 16 tasks to run the process.

(e) Structure HTML Data

- Parse .html and convert to .json bulk-load format

```
{
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "39"
  }
}, {
  "id": "39",
  "title": "Highlander Early Start Academy",
  "url": "http://earlystart.ucr.edu/",
  "level": 3,
  "filename": "data/level3/web_39",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "40"
  }
}, {
  "id": "40",
  "title": "R'Side Calendar",
  "url": "http://rside.ucr.edu/",
  "level": 3,
  "filename": "data/level3/web_40",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "41"
  }
}, {
  "id": "41",
  "title": "Citrus Collection Gifts: Home",
  "url": "http://citrusgifts.ucr.edu",
  "level": 3,
  "filename": "data/level3/web_41",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "42"
  }
}, {
  "id": "42",
  "title": "Home | Women's Resource Center",
  "url": "https://wrc.ucr.edu/",
  "level": 3,
  "filename": "data/level3/web_42",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "43"
  }
}, {
  "id": "43",
  "title": "Home | Chicano Student Programs",
  "url": "https://csp.ucr.edu/",
  "level": 3,
  "filename": "data/level3/web_43",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "44"
  }
}, {
  "id": "44",
  "title": "UCR Newsroom: Home",
  "url": "http://newsroom.ucr.edu",
  "level": 3,
  "filename": "data/level3/web_44",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "45"
  }
}, {
  "id": "45",
  "title": "Home | African Student Programs",
  "url": "http://asp.ucr.edu/",
  "level": 3,
  "filename": "data/level3/web_45",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "46"
  }
}, {
  "id": "46",
  "title": "Fleet Services: Home",
  "url": "http://fleet.ucr.edu",
  "level": 3,
  "filename": "data/level3/web_46",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "47"
  }
}, {
  "id": "47",
  "title": "Jobs - University of California",
  "url": "http://jobs.universityofcalifornia.edu/",
  "level": 3,
  "filename": "data/level3/web_47",
  "body": ""
}, {
  "create": {
    "_index": "webdocs",
    "_type": "webdoc",
    "_id": "48"
  }
}, {
  "id": "48",
  "title": "UCR ARTS",
  "url": "http://artsblock.ucr.edu",
  "level": 3,
  "filename": "data/level3/web_48",
  "body": ""
}
```

Figure 1: html web-documents readily available to load into an Elasticsearch Lucene Index.

Indexer

(a) Architecture: Ubuntu Virtual Machine on Virtualbox

```
ubuelasticsearch [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
* Management: https://landscape.canonical.com
* Support: https://ubuntu.com/advantage

System information disabled due to load higher than 1.0

* Ubuntu's Kubernetes 1.14 distributions can bypass Docker and use containerd
  directly, see https://bit.ly/ubuntu-containerd or try it now with

  snap install microk8s --channel=1.14/beta --classic

* Canonical Livepatch is available for installation.
  - Reduce system reboots and improve kernel security. Activate at:
    https://Ubuntu.com/livepatch

29 packages can be updated.
11 updates are security updates.

rogith@vubuntu:~$ curl 127.0.0.1:9200
{
  "name": "32rftA",
  "cluster_name": "elasticsearch",
  "cluster_uuid": "0zCHegQ0TMiqDVfy0Zf3Bw",
  "version": {
    "number": "6.6.1",
    "build_flavor": "default",
    "build_type": "deb",
    "build_hash": "1fd8f69",
    "build_date": "2019-02-13T17:10:04.160291Z",
    "build_snapshot": false,
    "lucene_version": "7.6.0",
    "minimum_wire_compatibility_version": "5.6.0",
    "minimum_index_compatibility_version": "5.0.0"
  },
  "tagline": "You Know, for Search"
}
```

Figure 1: Elasticsearch instance on the virtual Ubuntu server/cluster.

(b) Index Structure: JSON Mapping Structure of a Web-document

```
{  
  
    id: number,  
  
    title: string,  
  
    url: string,  
  
    level: string,  
  
    filename: string,  
  
    body: string  
  
}
```

(c) Search Algorithm

- Pose queries through RESTful API over HTTP requests (HEAD, GET, POST, etc.)
- Retrieve results based on multi-match of web-document fields: title, body

(d) Limitations

- Due to virtual cluster, only the host machine of the cluster may successfully make HTTP requests and thus run the application

Extension : Web Interface Front-End

(a) Frameworks used: Express and Node.js

(b) Ping the cluster when executing the application

(c) Allow user queries through a search bar and submit button

(d) Display ranked list of .edu web page documents

(e) Limitations

- We wish to modify the Embedded JavaScript file (.ejs) to search for the matched words in a query and only highlight those-
- The only way to access the virtual cluster is directly through the host running the virtual machine

```
rogith@DESKTOP-1Q8NDL3:~/finalproject-lightsabr/FrontEnd-master$ node app.js
body-parser deprecated undefined extended: provide extended option app.js:14:2
body-parser deprecated undefined extended: provide extended option app.js:32:2
server is on port 3000
all is well
```

Figure 1: Ping the Elasticsearch cluster when executing application (HEAD request).

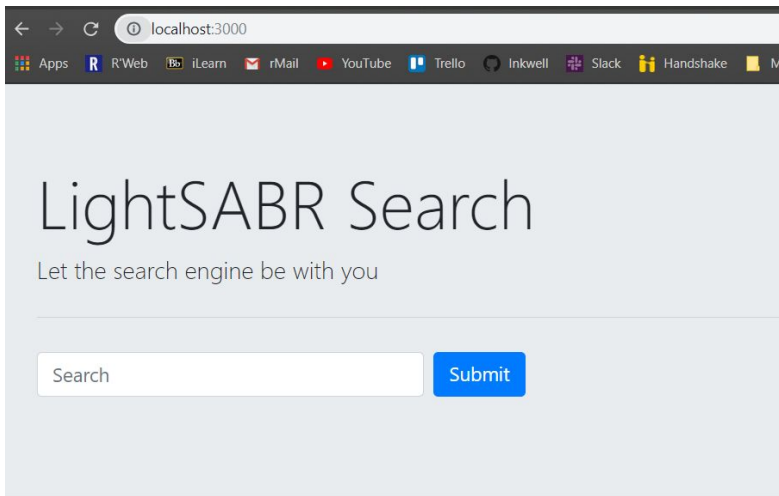


Figure 2: User interface.

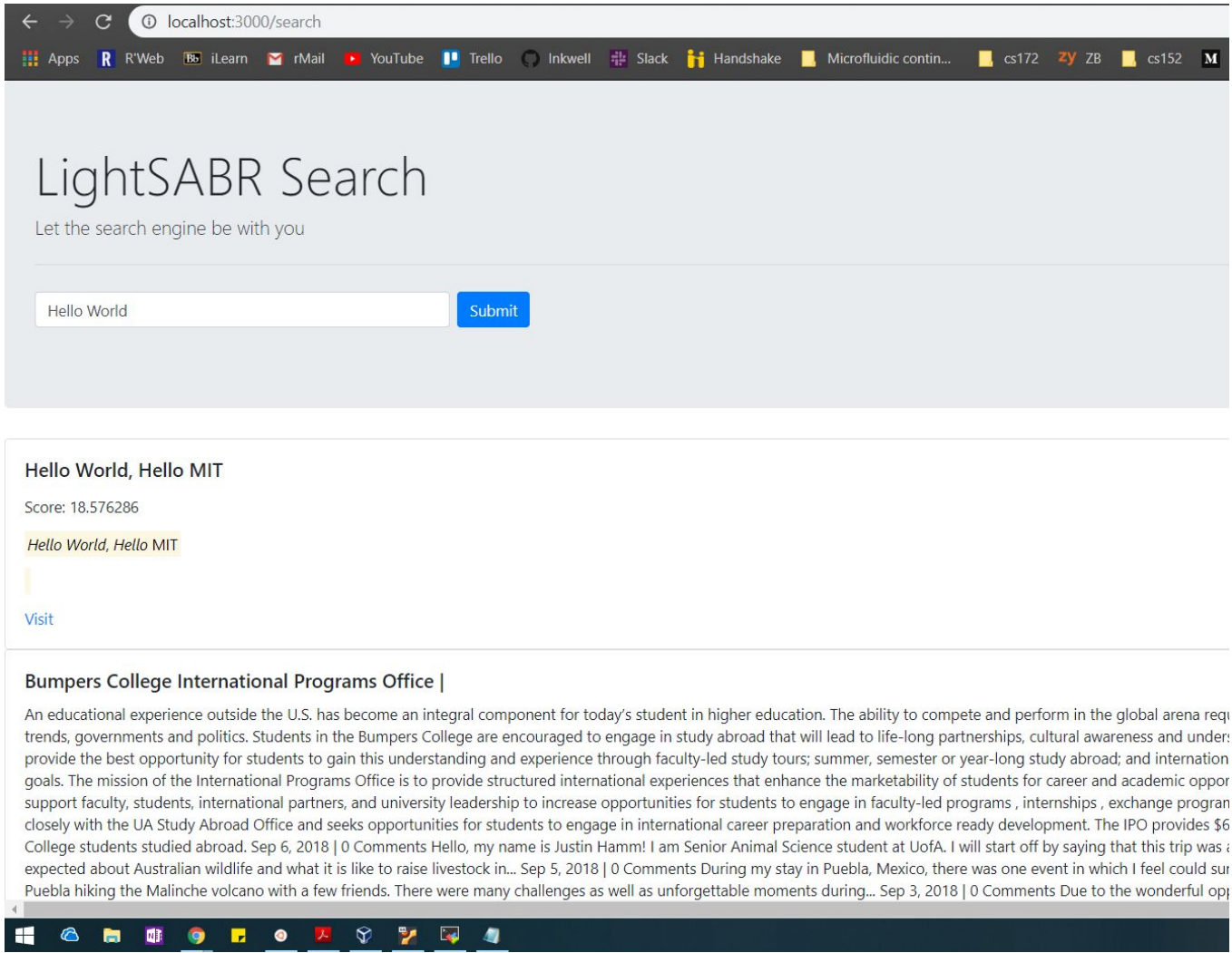


Figure 3: Search results in action.

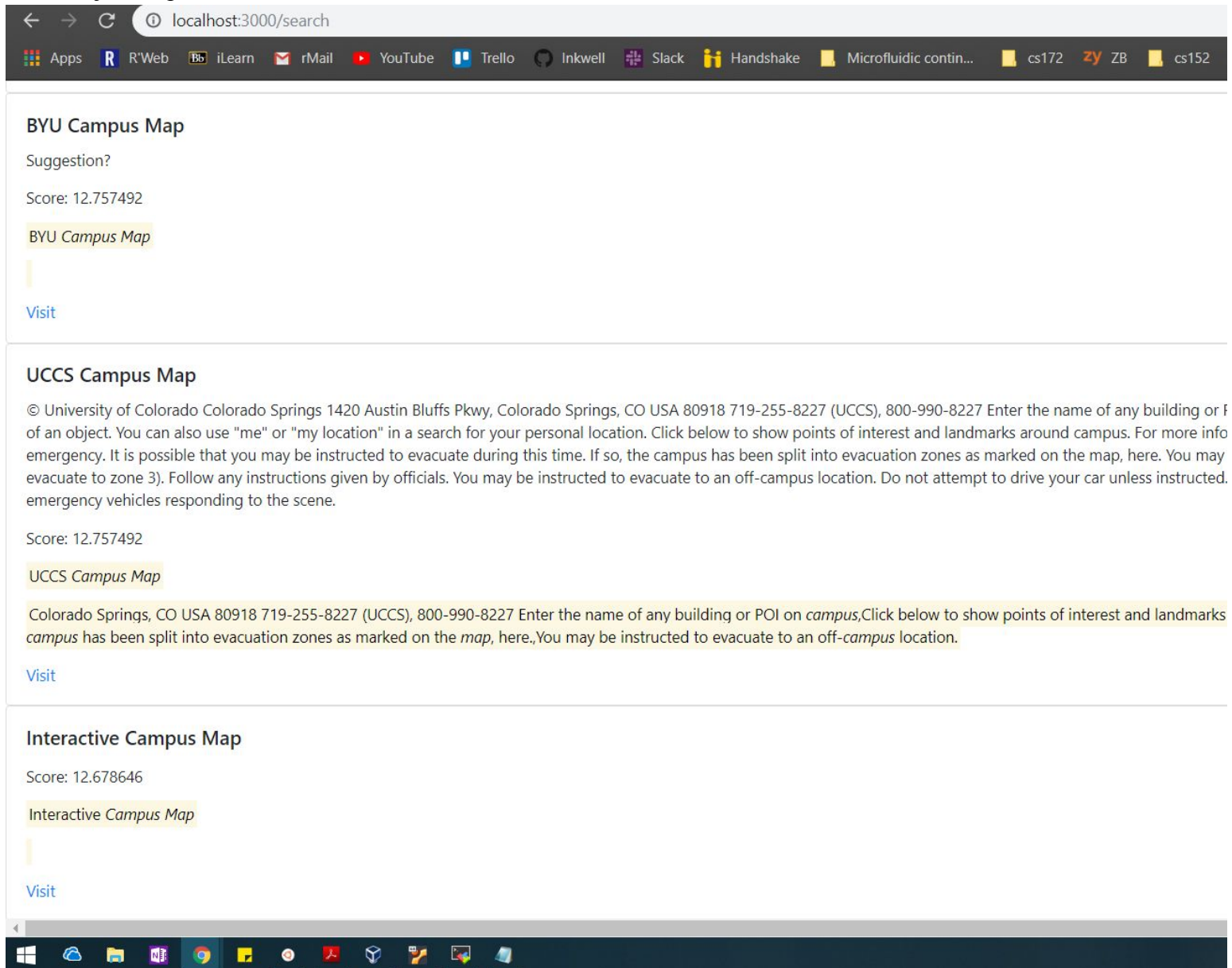


Figure 4: Highlights displayed from title and body matches (query shown: “campus map”).

Sample Test Cases and Data

```
curl -XPUT 127.0.0.1:9200/_bulk -d '

{ "create": { "_index": "movies", "_type": "movie", "_id": "135569" } }
{ "id": "135569", "title": "Star Trek Beyond", "year": 2016, "genre": ["Action", "Adventure", "Sci-Fi"] }
{ "create": { "_index": "movies", "_type": "movie", "_id": "122886" } }
{ "id": "122886", "title": "Star Wars: Episode VII - The Force Awakens", "year": 2015, "genre": ["Action", "Adventure", "Fantasy", "Sci-Fi", "IMAX"] }
{ "create": { "_index": "movies", "_type": "movie", "_id": "109487" } }
{ "id": "109487", "title": "Interstellar", "year": 2014, "genre": ["Sci-Fi", "IMAX"] }
{ "create": { "_index": "movies", "_type": "movie", "_id": "58559" } }
{ "id": "58559", "title": "Dark Knight, The", "year": 2008, "genre": ["Action", "Crime", "Drama", "IMAX"] }
{ "create": { "_index": "movies", "_type": "movie", "_id": "1924" } }
{ "id": "1924", "title": "Plan 9 from Outer Space", "year": 1959, "genre": ["Horror", "Sci-Fi"] } }
```

Figure 1: Test data that was first loaded to the Elasticsearch cluster to render results onto the interface. (Can be found from the Udemy sources above).

Trello Board Project Organizer

For a comprehensive look at the history of our development, please take a look at our public team

Trello board: <https://trello.com/b/lz8MPUvc/cs172-final>

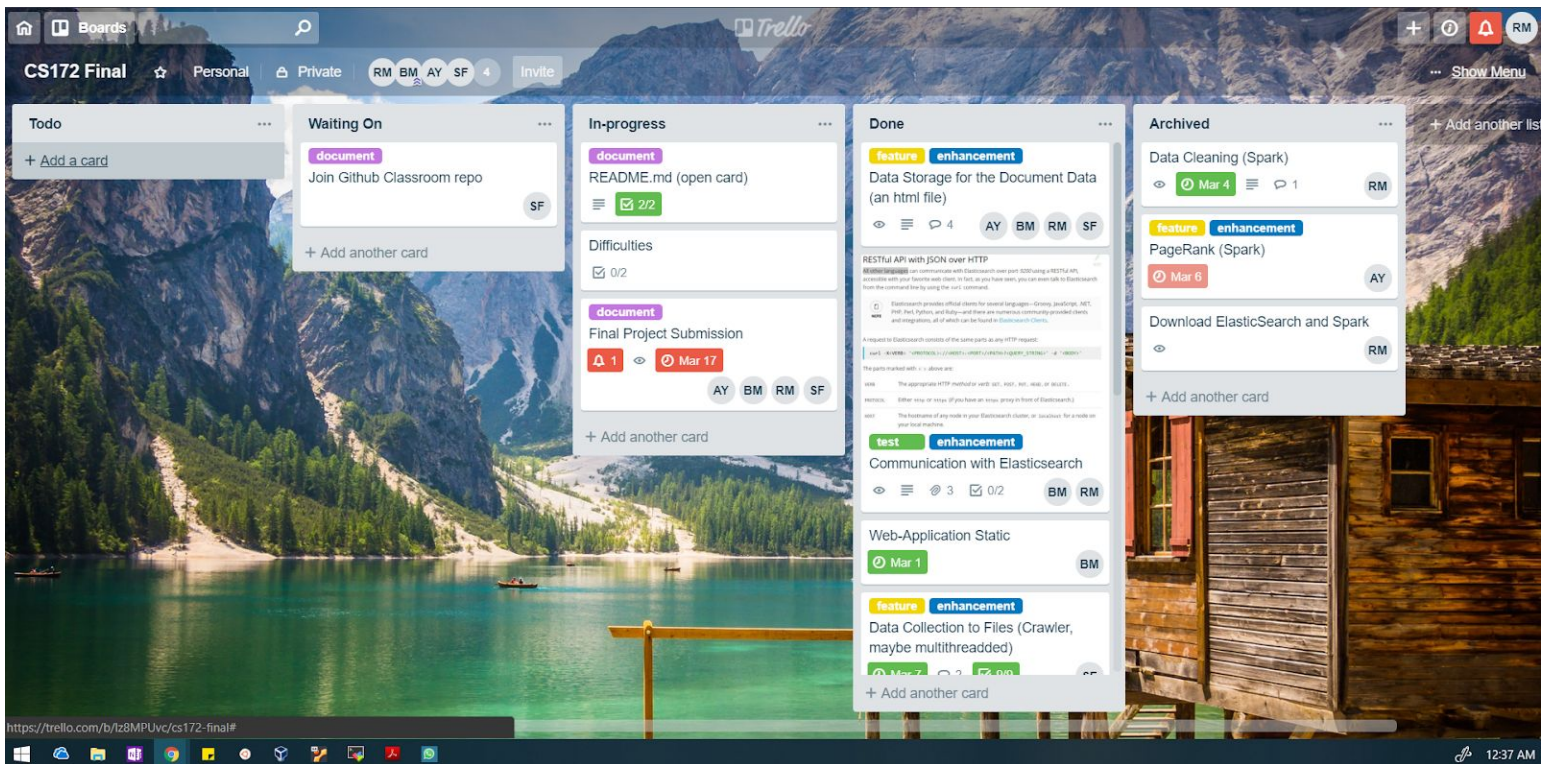


Figure 1: Team LightSABR Trello Board.

Data Storage for the Document Data (an html file)
in list [Done](#)

MEMBERS
AY BM RM SF +

LABELS
feature enhancement +

Description [Edit](#)

How to store the html data and its cleaned version

How should we store the collected data? If we are thinking of working with json files, i think we should also figuring out how to convert data into `.csv` and `.json`. Maybe this is where we can use Spark.

Store the Following Fields (either `.csv` or `.json`)

- Title
- Body
- Creation Date (if available)

```
import csv

with open('mycsv.csv', 'w', newline='') as f:
    fieldnames = ['column1', 'column2', 'column3']
    lilwriter = csv.DictWriter(f, fieldnames=fieldnames)

    lilwriter.writeheader()
    lilwriter.writerow({'column1': '1', 'column2': '2', 'co
```

code to generate csv file in python

Generates `mycsv.csv`

ADD TO CARD

- Members
- Labels
- Checklist
- Due Date
- Attachment

POWER-UPS

[Get Power-Ups](#)

ACTIONS

- Move
- Copy
- Watch ☒
- Archive
- Share

Figure 2: A completed Trello card that led us towards converting html data into json format: Data Storage for the Document Data (an html file).

Conclusion

At first glance, although the final project seemed very intimidating because of how many critical components there are, our group believes that producing a working search engine provides us with valuable knowledge that helps us when working in the industry. This knowledge includes how to build a web crawler, experience using an index API (Elasticsearch and Lucene), experience with working with html and json data, implementing a front-end user interface and overall how to efficiently work with big data.

A difficult challenge our team faced while working on the final project together was integrating all of our work into a finished final product. The fear of integration influenced us to use a scrum-like system to prioritize crucial components for functionality to optimize work efficiency. We implemented a Trello board to easily partition and prioritize each individual task and also use the board to archive a list of tasks that are complete.

Some of the further build-upons are mentioned in each [limitations](#) section above.