# Part B: Recipe Search Engine Interface

## Repository

https://github.com/rmace001/RecipeSearchEngine

## Team Members

Rogelio Macedo - rmace001@ucr.edu
Jerry Tan - jtan027@ucr.edu
Erin Wong DeForest - ewong012@ucr.edu
Shiyao Feng - sfeng023@ucr.edu

# Collaboration Details

## Map Reduce/ Data Scripts

Jerry Tan
Entirely responsible for developing multiple scripts to change the formatting of overall recipe data on elastic search index, as well as expressing the hadoop index into key value pairs in a csv file. Furthermore, the scripts handled all data-cleaning and preprocessing, as well as schema conversions in effort to allow for sorting of recipe results by active time and total times.

## Web Interface Backend/Frontend

Rogelio Macedo & Shiyao Feng

### Elastic Enterprise App Search API

Set up a deployment on the Elastic Cloud Service, along with an Enterprise Search API engine

### ReactJS Web Application

Provides an interface for querying the Elastic App Search API.

### Hadoop MapReduce Index Insertion in Cassandra

Instantiation and table creation with Apache Cassandra on an Ubuntu 18.04.5 virtual machine.

### Express, NodeJS Web Application

Provides an interface for querying the Apache Cassandra, along with ranking and displaying of documents.

Hadoop MapReduce Index Generation

Erin Wong DeForest
Used Hadoop on a single cluster and wrote MapReduce program to create the inverted index of recipe terms from the json data. Reformatted the data for input into Cassandra database for use in querying.

# Overview of System

## Architecture

We use Hadoop and MapReduce to create an inverted index of recipe terms where the document IDs are recipe links. We built our search engine frontend using Elastic Enterprise App Search with React as an interface. In the backend, we use Cassandra DB to store our inverted index and recipe data to return and display results upon querying. Since we are creating a search engine for culinary purposes, we have chosen to return queries based on the intersection of search terms. We made this choice because people searching for recipes often do so based on specific ingredients they have available and wish to cook with, for example "rice coconut milk".

## Hadoop MapReduce Indexing

We ran Hadoop on a single cluster to process our dataset of recipes. The program begins by editing the input file(s) to be one line per recipe and removing active and total time and JSON fields and punctuation as well as inserting tabs as delimiters between each field. Afterwards, MapReduce is used to create an inverted index where recipe links are used as document IDs because they are unique. All terms from text fields are used as keys; again, this only excludes active and total time.

While we removed punctuation, we chose not to perform stemming because many recipes contained foreign words which may be unintentionally altered by a stemmer. We handled foreign characters by using UTF-8.

The Mapper receives one recipe at a time, separating the link from the recipe and tokenizing all other text. The tokens are then emitted along with the recipe link. The Reducer compiles the inverted index for each term by appending recipe links.

The runtime of the Hadoop index construction was consistently about 25 seconds. There is no comparison of runtime to the Lucene Index because our ElasticSearch index was created automatically by uploading to the Elastic Enterprise App Search Server.

Figure 0: Partial printout of inverted index.

# Elastic Enterprise App Search with ReactJS

This involved uploading our preprocessed JSON recipe data to the cloud, which allowed for automatic indexing of documents. The cloud service provided an easy web interface for viewing individual documents, as well as other features like changing the schema for the entire dataset with just a few clicks. In addition, it enabled fast and secure queries over the web, and provided an intuitive npm client for issuing queries and neatly-displaying results with ReactJS.
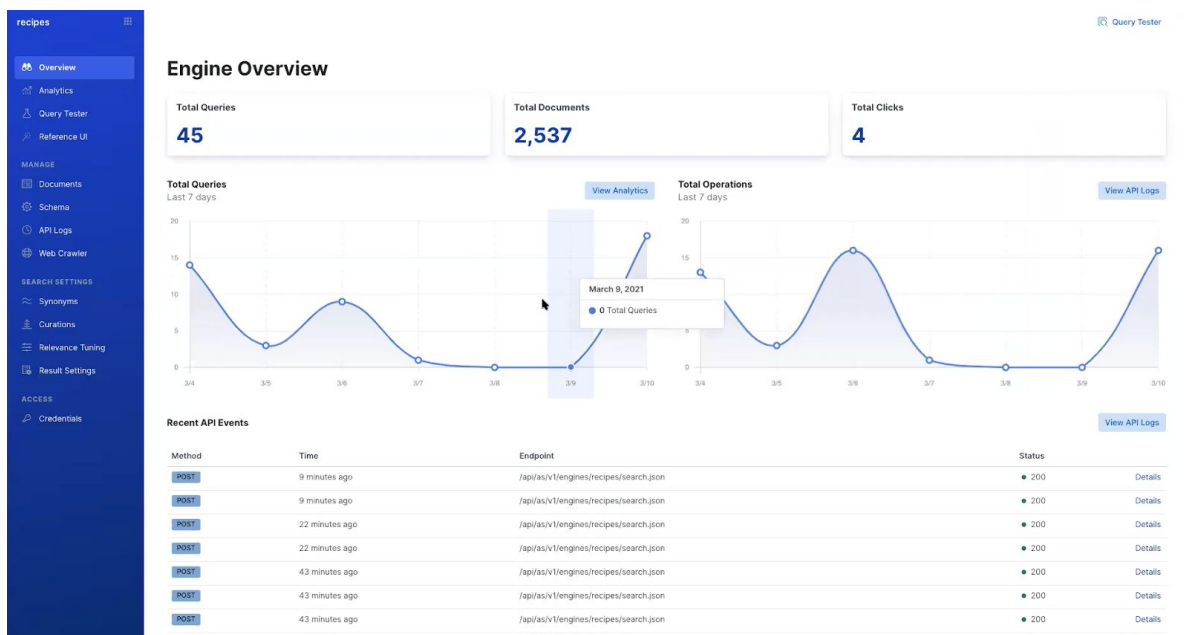

Figure 1: Search Engine UI provided from Elastic Cloud.

The React application was bootstrapped with the `create-react-app` command line interface. Furthermore we used the Elasticsearch npm client to connect to the remote endpoint, allow for filtering of recipe results by active times and total times, and sort the results by relevance scores, or by the active and total times. Elasticsearch, by default, uses BM-25 for ranking.



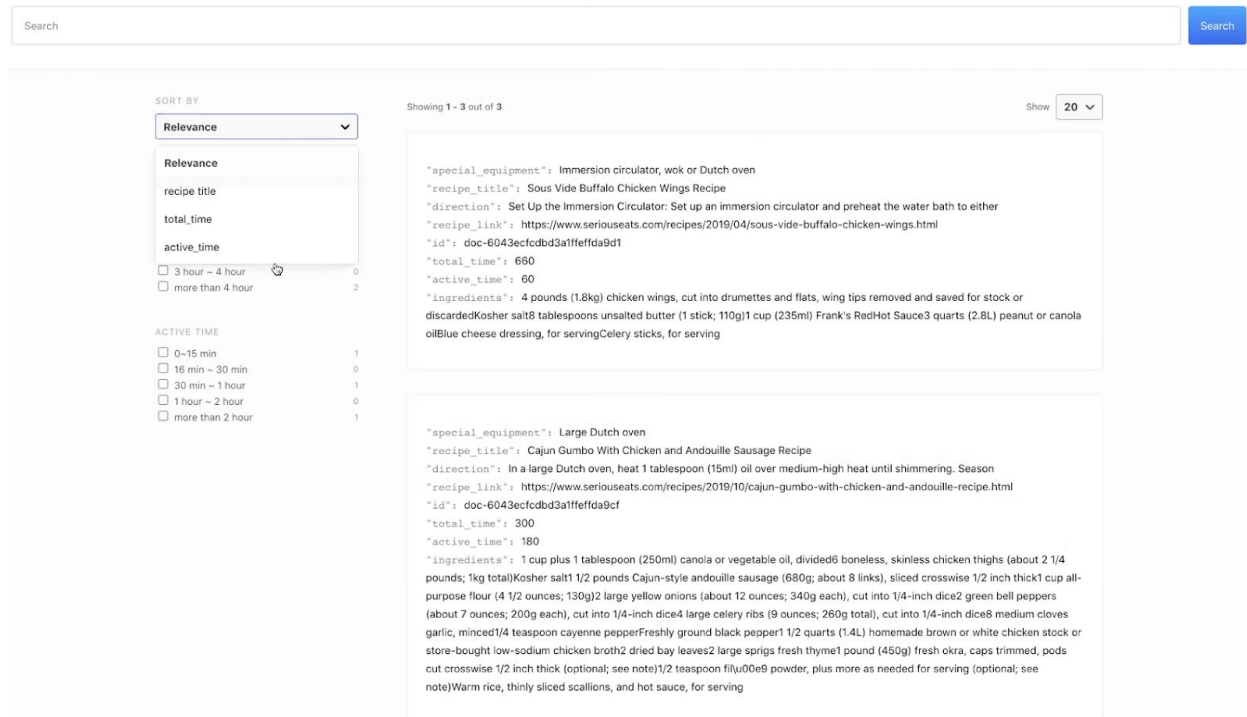Figure 2: ReactJS Search Engine UI showing results, and ability to filter/sort results.

## Apache Cassandra with Express/NodeJS

We setup Apache Cassandra on an Ubuntu 18.04.5 virtual machine, loaded the index into a table, and finally, loaded the entire dataset of documents into another table. We were able to use the Cassandra command line interface to load data into a table from a .csv file. Our index.csv consists of two columns, (1) term, (2) list of links where term is found; the term is the primary key of this table. The document dataset.csv consists of all of the fields we chose to store per recipe document, and the primary key of this table is the recipe link.

After instantiating the database, it was time to connect to it using the Express/NodeJS application. Luckily, the npm library provided a Cassandra driver, where we learned how to connect to the database, and perform queries. The connection to the database was the easy part, as we were able to console.log hardcoded queries from the node.js application. The difficult part was to design a web UI displaying a search box and submit button and serve it using Express and NodeJS. Furthermore, we spend a lot of time figuring out how to take query results and display them in a presentable fashion on the front end.

Figure 4: Screenshot of Express/NodeJS interface after searching for "Pumpkin Pecans", demo of screenshot here:

https://drive.google.com/file/d/1RdDPGTuXpGI7y9YxKl_cBGtXo9Rq68Za/view?usp=sharing

# System Limitations

Limitations regarding the Hadoop Search Engine consisted of choice of ranking due to much more time invested in incorporating MapReduce-generated index into our Cassandra Database and connecting the DB and backend server into our frontend. So we were not able to fully implement an ideal ranking system. We struggled so much with uploading our full dataset into Cassandra, since we had to spend a lot of overhead figuring out how to change ~2500 documents from JSON data to csv. Due to issues such as delimiting data by commas, or strings being excessively long, some of the documents were not added to Cassandra, as Cassandra would fail on these instances.

The bulk of the time was spent on learning and implementing a full-stack web development. Once all of the pieces were connected together (i.e. back end, middle end, front end), we finally tackled results displaying, as well as results ranking. Every keyword in our query needed to be queried, and we did not have time to implement a more effective method of querying such as batch queries. For every user query, it required multiple queries, one per query keyword. Although the method is easily explainable, it is unfortunate that it is also highly inefficient. Furthermore, the ranking is very basic, as we simply just perform the intersection of documents as a method of ranking between search results.

One other important limitation is how we obtain the set of intersected documents. Since we did not have time to figure out how to query efficiently, we simply decided to query Cassandra for every document and the intersected list. This turned out to be highly inefficient, though simple in design and implementation.

# Obstacles and Solutions

We struggled to choose a ranking algorithm, especially leaning towards BM-25 due to its consideration of document and query term weights. Ultimately, we chose to use a non-stemmed intersection-based ranking system because we agreed that the most common use case is searches by ingredients, many of which are foreign words. For example, we expect poor query results if "ricer" were to be stemmed to "rice" (the former referring to a kitchen tool, the latter referring to the popular grain). Furthermore, we concluded that the common case of the singular and plural was a nonissue since a recipe with an ingredient in the title would often refer to it in the opposite form. For example, "Tomato Soup" would surely include "tomatoes" in the ingredients and directions while the same would be true of "Refried Beans" and "bean", negating the benefits of stemming and avoiding the complications of stemming foreign words.

# Deployment Instructions

Please refer to the following text documents located inside the root directory of the repository linked at the top of this document.
- Installing-Cassandra-on-Ubuntu-18.04.5.txt
- ejs-mrsearch/instructions.txt

For Hadoop and MapReduce, please use Maven to run the files included in the Hadoop_MapReduce directory. We have been running it using the following command lines:

```
rm output; mvn package && hadoop jar
target/MapReduce-1.0-SNAPSHOT.jar
edu.ucr.cs.cs242.RecipeSearch.MapReduce
```

# References

https://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html

https://timepasstechies.com/map-reduce-inverted-index-sample/

https://codeburst.io/how-to-build-great-react-search-experiences-quickly-8c69081f328d

https://cassandra.apache.org/doc/latest/getting_started/installing.html

https://www.npmjs.com/package/cassandra-driver

https://github.com/DataStax-Examples/quickstart-nodejs/blob/master/quickstart.js

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