



KIV/IR Semester project

# Search Engine

Jakub Šilhavý  
A21N0072P  
silhavyj@students.zcu.cz

19. 03. 2022

# Contents

|          |                                 |           |
|----------|---------------------------------|-----------|
| <b>1</b> | <b>Introduction</b>             | <b>1</b>  |
| 1.1      | User manual . . . . .           | 1         |
| <b>2</b> | <b>Implementation</b>           | <b>2</b>  |
| 2.1      | Technology . . . . .            | 2         |
| 2.1.1    | Programming language . . . . .  | 2         |
| 2.1.2    | Build system . . . . .          | 2         |
| 2.1.3    | Dependencies . . . . .          | 2         |
| 2.2      | Architecture . . . . .          | 3         |
| 2.3      | Input data . . . . .            | 3         |
| 2.3.1    | BBC News . . . . .              | 4         |
| 2.3.2    | TREC . . . . .                  | 4         |
| 2.4      | Preprocessing . . . . .         | 4         |
| 2.4.1    | Language detection . . . . .    | 4         |
| 2.4.2    | Stemming . . . . .              | 5         |
| 2.5      | Inverted Index . . . . .        | 5         |
| 2.6      | Search query . . . . .          | 6         |
| 2.6.1    | Query evaluation . . . . .      | 6         |
| 2.7      | Ranking document . . . . .      | 8         |
| 2.7.1    | Bag of words . . . . .          | 8         |
| 2.7.2    | Cosine similarity . . . . .     | 8         |
| 2.7.3    | TF-IDF . . . . .                | 8         |
| <b>3</b> | <b>Testing</b>                  | <b>9</b>  |
| 3.1      | Operation . . . . .             | 9         |
| 3.2      | Query . . . . .                 | 9         |
| 3.3      | Search . . . . .                | 9         |
| <b>4</b> | <b>Conclusion</b>               | <b>10</b> |
| 4.0.1    | Implemented features . . . . .  | 10        |
| 4.0.2    | Achieved performance . . . . .  | 10        |
| 4.0.3    | Possible improvements . . . . . | 10        |

# 1 Introduction

This application represents a simple search engine that was developed as a semester project within the KIV/IR module at the University of West Bohemia. Using a boolean expression, it allows the user to search for specific documents within a given set of documents that was inputted into the application. The results of a search can be ranked either by cosine similarity or tf-idf. It also supports automatic language detection. However, for the sake of simplicity, the user can only search for documents in Czech or English.

## 1.1 User manual

The main purpose of this document is to explain the details of how this assignment was tackled. If you are looking for a user manual on how to use this application, please visit

<https://github.com/silhavyj/KIV-IR-Search-Engine>.

## 2 Implementation

### 2.1 Technology

#### 2.1.1 Programming language

The entire application is written in **Java**. This decision was made based on previous experience with this programming language as well as its wide range of library support. The project was developed on Linux using `openjdk 11.0.14` 2022-01-18.

#### 2.1.2 Build system

As for the build system, the project uses **Maven**, primarily because of its convenience when it comes to importing third-party libraries. The version of **Maven** used within this project is `3.6.3`

#### 2.1.3 Dependencies

The entire list of dependencies along with their versions could be found in the `pom.xml` file, which is located in the root directory of the project structure.

- `org.openjfx` (JavaFX GUI library)
- `junit` (unit testing)
- `org.apache.opennlp` (stemmer - preprocessing)
- `com.github.pemistahl` (language detection)
- `org.json` (working with the JSON format)
- `org.jsoup` (fetching down articles from the internet)

## 2.2 Architecture

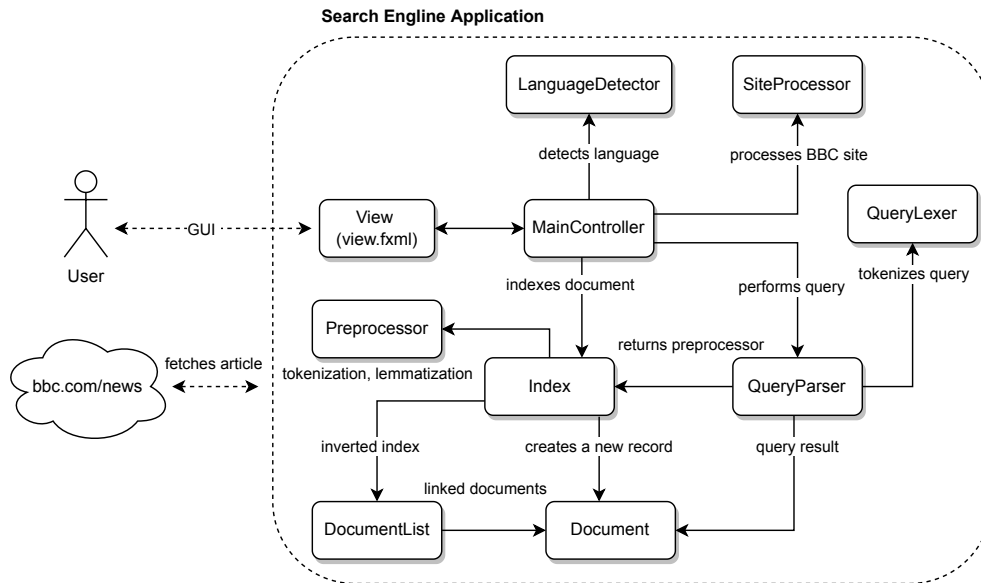


Figure 2.1: Graph of dependencies within different modules

As you can see in the figure above, the user interacts with the system via a GUI interface. The `MainController` class handles all events triggered by the user, e.g. inserting a document, executing a query, etc. With each event, it uses the backend logic, which is independent of the GUI, to perform the desired task. For more information about each class, please see the [Javadoc](#) documentation.

## 2.3 Input data

The application expects the document to be in the JSON format of the following structure.

```
{ "datetime": "",
  "author": "",
  "title": "",
  "article": "",
  "url": ""
}
```

It is worth mentioning that the only obligatory fields are the title and article. The rest is just metadata that will be displayed along with each result of a search. The application distinguishes the documents by their location on the disk, not by their contents. If the user attempts to insert multiple instances of the same file (document), the first document will be indexed and the rest will be discarded.

All documents that come with the application can be found hosted over at <https://cutt.ly/5SzcWSG>.

### **2.3.1 BBC News**

These documents represent articles that were fetched down using a technique called crawling. In total, there are 8,486 articles. All the articles were downloaded from <https://www.bbc.com/news>. Each document contains a URL that leads to its the original location on the website.

### **2.3.2 TREC**

The documents stored in the trec folder were provided with the assignment as a serialized binary file. In order to keep things consistent, they were converted into individual JSON files which have the same structure as the documents downloaded from the BBC News. The overall number of the documents is 81,735.

## **2.4 Preprocessing**

Before a document is indexed, it needs to be preprocessed first. This process involves language detection, tokenization, accent removal, removing stop-words, and stemming.

### **2.4.1 Language detection**

In order to support multiple language, I needed to take advantage of the following library <https://github.com/pemistahl/lingua>. This library supports many different language. However, because each language needs to have its own stemmer and tokenizer, the number of languages supported in this application was reduces to two - Czech and English.

During debugging, it was discovered that the Czech language gets sometimes confused with the Slovak language. Therefore, the Slovak language and the Czech language are treated as the same.

## 2.4.2 Stemming

The process of stemming is done separately for both languages. For each language, there is a file containing stopwords of that language. All stop words are excluded from a document before it is indexed. The files can be found in the root folder of the project structure - `stopwords-cs.txt` and `stopwords-en.txt`.

### English

As far as English is concerned, I used the following library to do the job <https://github.com/apache/opennlp>. In particular, I used the `PorterStemmer` class.

### Czech

As for the Czech language, I used the same class that was used in the second practical of the KIV/IR module. This class can be found over at <https://github.com/dundalek/czech-stemmer>.

## 2.5 Inverted Index

There is a separate index for each language detected in the documents. Which index will be used for searching is determined by the language detected in the query. Inverted index is represented as a `HashMap` of `LinkedList`s. The key of each chain of documents is a term. The value represents a list of documents where the term occurs. The `Index` class also holds a map of locations of individual files on the disk.

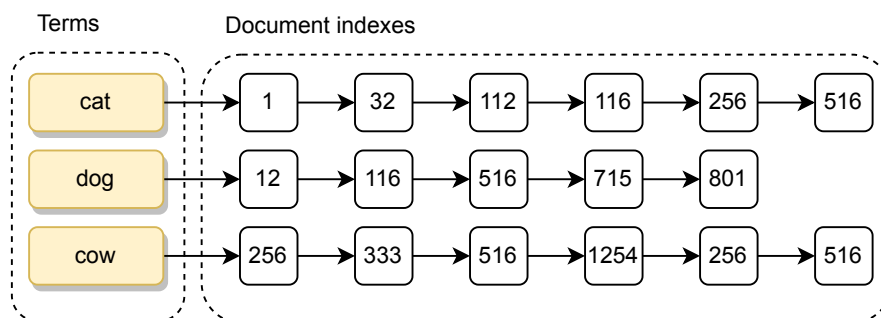


Figure 2.2: Structure of an inverted index

## 2.6 Search query

### 2.6.1 Query evaluation

In order to search for documents, I needed to evaluate a boolean expression the user inputs into the application. This expression, as required, supports the AND, OR, and NOT operators. Operands represents individual key words the user wants to find in documents. If two words are not explicitly separated by AND or OR, they are automatically evaluated as `term1 AND term2`.<sup>1</sup>

In terms of implementation, I decided to represent the operators as '&', '|', and '!' respectively. The user can also further specify the priority using a pair of parentheses. Examples of the query syntax can be seen in the highlighted box below.

```
> price & (gas | petrol) & (increase | high) & Northern Ireland  
> Biden & !Trump  
> Christmas & !(omicron | coronavirus)
```

As for query evaluation, I first converted the infix notation into a postfix notation, which is easier to parse. An explanation of this algorithm can be seen, for example, at <https://www.geeksforgeeks.org/stack-set-2-infix-to-postfix/>. An expression written in a postfix notation can be easily evaluated using a stack. An algorithm for this can be also found over at <https://www.geeksforgeeks.org/stack-set-4-evaluation-postfix-expression/>. Essentially, it all comes down to being able to evaluate a simple expression of one the following forms: `term1 & term2`, `term1 | term2`, or `!term1`.

### AND

Calculating AND for two given sets of documents means finding the documents that appear in both sets. Given the fact that two sets are represented as sorted linked lists, using two pointers, we can calculate the intersection in linear time.

---

<sup>1</sup>Each word found in a query is also preprocessed using the same technique that is used for indexing a document.



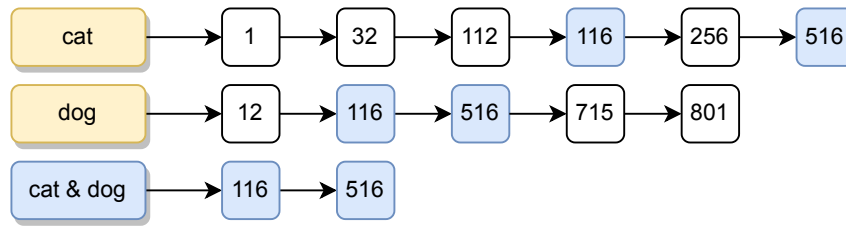


Figure 2.3: Example of an AND operation

The final list of intersecting documents is also sorted, so it can be reused later on.

## OR

The OR operation works as a union of two given sets of documents. This operation needs to make sure that there are no duplicities in the final result. With the assumption that the two sets are sorted linked lists, the OR operation can be calculated, also, in linear time.

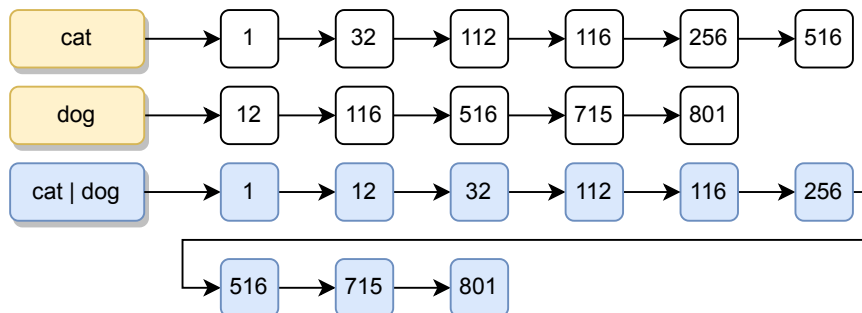


Figure 2.4: Example of an OR operation

## NOT

The NOT operation returns a complement of a given set of documents. It uses the Index class to retrieve all documents that have been indexed. In other words, it returns all documents except for those which are passed in as a parameter.

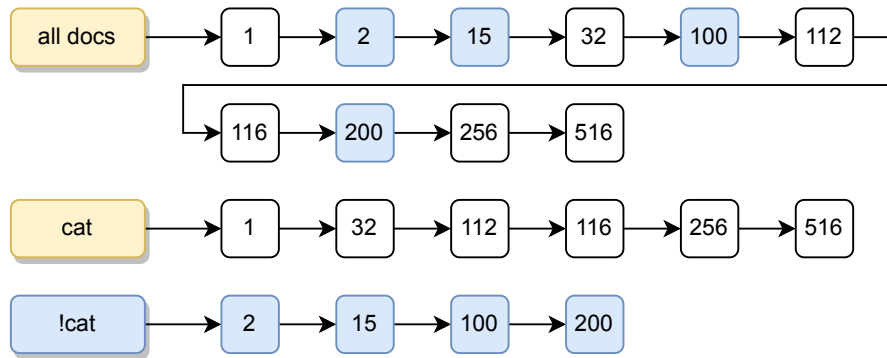


Figure 2.5: Example of an NOT operation

## 2.7 Ranking document

### 2.7.1 Bag of words

There is also an in-memory **bag-of-words** representation of each document held in the `Index` class. Essentially, there were two options. The first option was to create a **bag-of-words** representation for specific documents based on a given query. The advantage of this is that it does not require as much space for the representation to be created. It is simply created upon evaluation of a query, and destroyed by the garbage collector afterwards. The trade off is that ranking might be a bit slower.

### 2.7.2 Cosine similarity

### 2.7.3 TF-IDF

## **3    Testing**

### **3.1    Operation**

### **3.2    Query**

### **3.3    Search**

## 4 Conclusion

4.0.1 Implemented features

4.0.2 Achieved performance

4.0.3 Possible improvements