# 1 Search Engine

## 1.1 Dataset

The data was acquired from crawling simple.wikipedia.com pages using wikipedia_crawler_2.py script.

The crawling script was run for around 5 hours and collected **89 535 articles** about various topics.

The starting point for crawling was page "World War 2" with max depth set to 4.

Many of the articles are related to history, but due to the depth of the search there are lots of articles which have zero connection to history (e.g. the dataset contains "(You Drive Me) Crazy" article, which is a song by Britney Spears).

## 1.2 Data processing

### 1.2.1 Raw articles processing

After collecting articles, the raw text was preprocessed using ArticleProcessor.py script, which stripped the sections at end of the article (references and other websites section), then it used ntlk library to remove stop words and punctuation and them, used Porter stemming algorithm to stem the words. The processed articles were saved into new folder to process further.

This was done with multi threading to speed the process up.

### 1.2.2 Dictionary creation

The processed articles were used by DictionaryCreator.py script, which collected every word used in articles into a set, it also keep the count of each word occurrence. It also removed every word which had length lower than 3, and words which occurred less than 15 times.

This was done with multi threading to speed the process up.

The final size of the dictionary was **35 030 words**

### 1.2.3 Article lookup map creation

The processed articles were used by ArticleLookupMapCreator.py, which assigned an index to each article. This map will be used later by the engine to map the resulting index to a specific article name. This is done by inserting the article name into an array.

### 1.2.4 Creating term-by-document matrix

The ArticleVectorizer.py script utilizes the processed articles and the dictionary. It creates a MxN matrix, where M is the number of articles and N is the number of words in dictionary. Each row $R_i$ represents how many occurrences of the dictionary words are in a article whose index is i. For example cell $C_{i,j}$ represents how many times word with index j occured in article with index i. The created matrix is sparse, so we handle it using scipy.sparse methods. After creating the matrix, each column is multiplied by inverse document frequency (IDF) value, which is calculated for each word:

$$IDF(w) = log\frac{M}{n_{\mathrm{w}}} \tag{1.1}$$

Where $M$ is the count of articles, and $n_w$ is the count of articles contacting at least one word $w$. We do that to reduce common word significance (word occurring in only one document will have very high IDF value, while word occurring in every document will have IDF equal to 0) After calculating, the matrix is normalized by $L_2$ norm and then saved on the disk to be loaded later.

### 1.2.5 Low rank approximation

The created term-by-document matrix is loaded by ArticleVectorSVD.py script, and then it uses sklearn.decomposition.TruncatedSVD class to calculate matrices for low rank approximation of our original matrix.

My computer could not handle computing the full matrix approximation(U*S*Vh), so the script returns two matrices (U*S) and Vh. I tried some normalization techniques, and i managed to get very good results by using normalized US matrix by rows and normalized Vh matrix by columns. The matrices are then saved into a file to be loaded later.

I calculated approximations of the matrix using values of k: 500,1000,2500,5000

## 1.3 Querying the data

The data can be queried using the class inside search_engine.py file. The SearchEngine loads the computed dictionary, article lookup map and search matrices(without low rank approximation and with low rank approximation). To query the data, search method is called. It requires user string input, which is sanitized the same way as in 1.2.1 section, then the query vector is created. The query vector has size N, and each cell $C_i$ represent how many times word with index i occurred in the query, then the vector is normalized by $L_2$ norm. The search method then performs a search matrix multiplication with created query vector, the resulting vector

has size M, and each cell $C_i$ represent how similar the article with index i, is to the provided query, the similarity is based on cosine similarity.

# 2 Frontend application

The frotend layout is served by Flask server, and the frontend files were created using React.
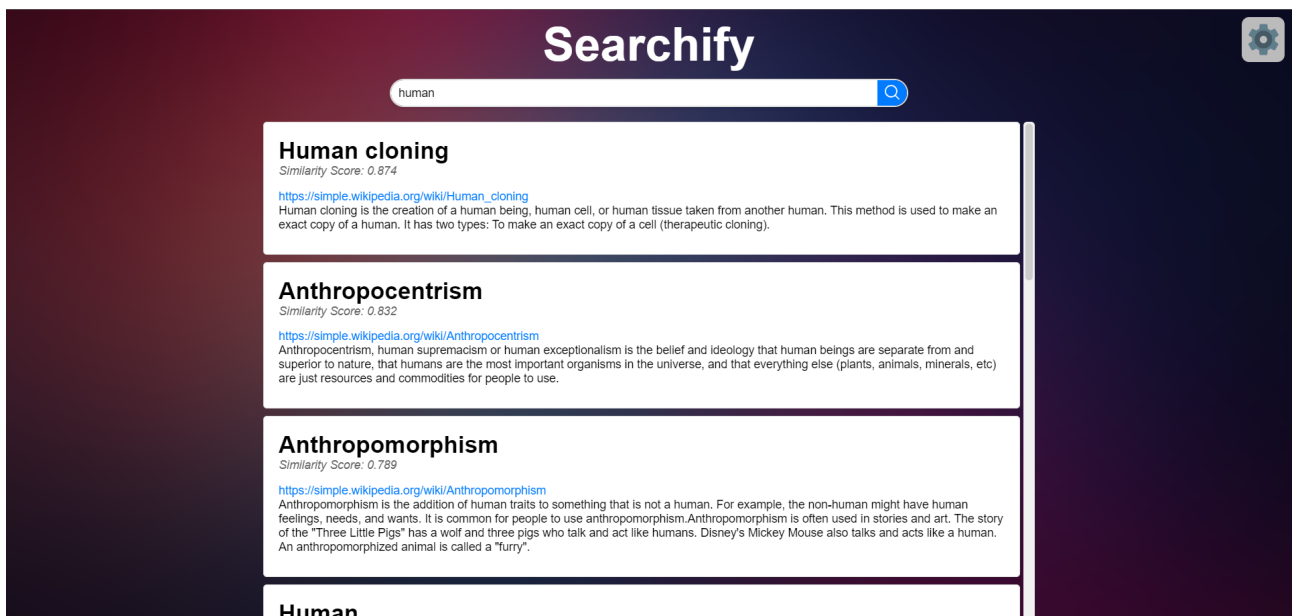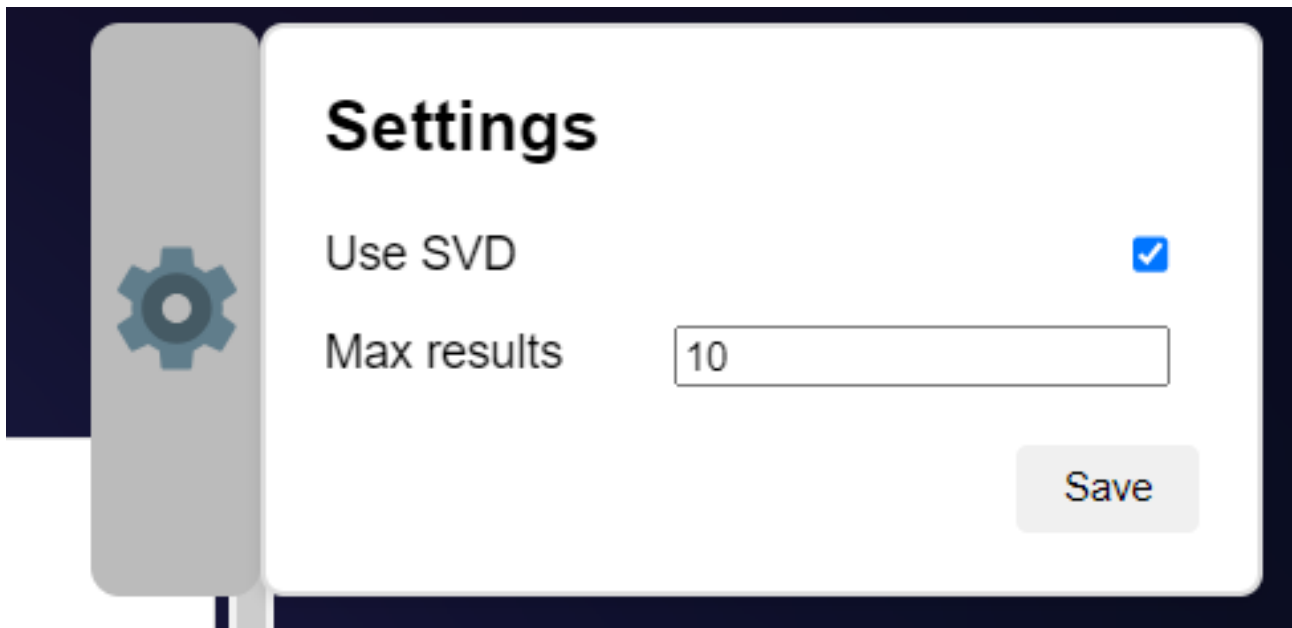


Figure 2.1: Search page



Figure 2.2: Search results

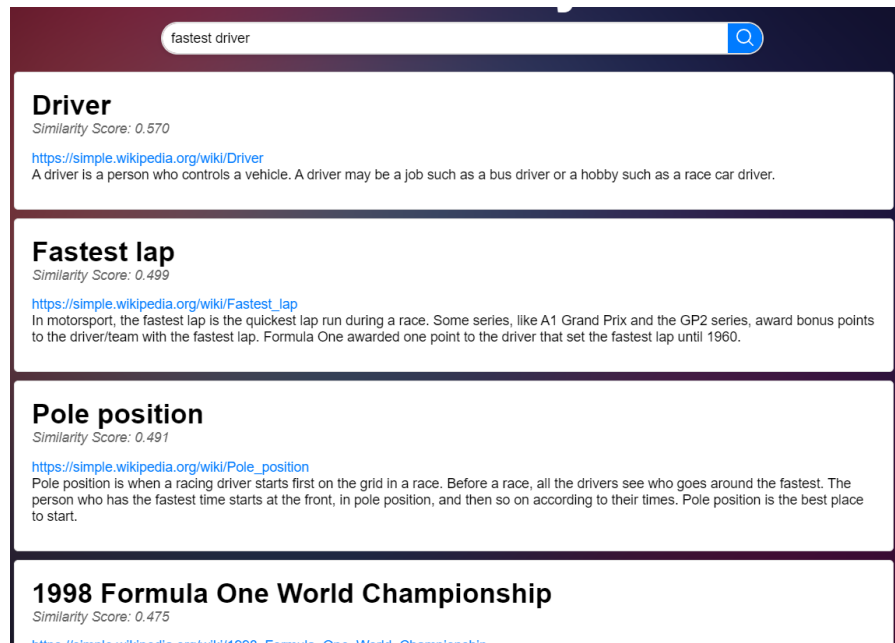Figure 2.3: Settings

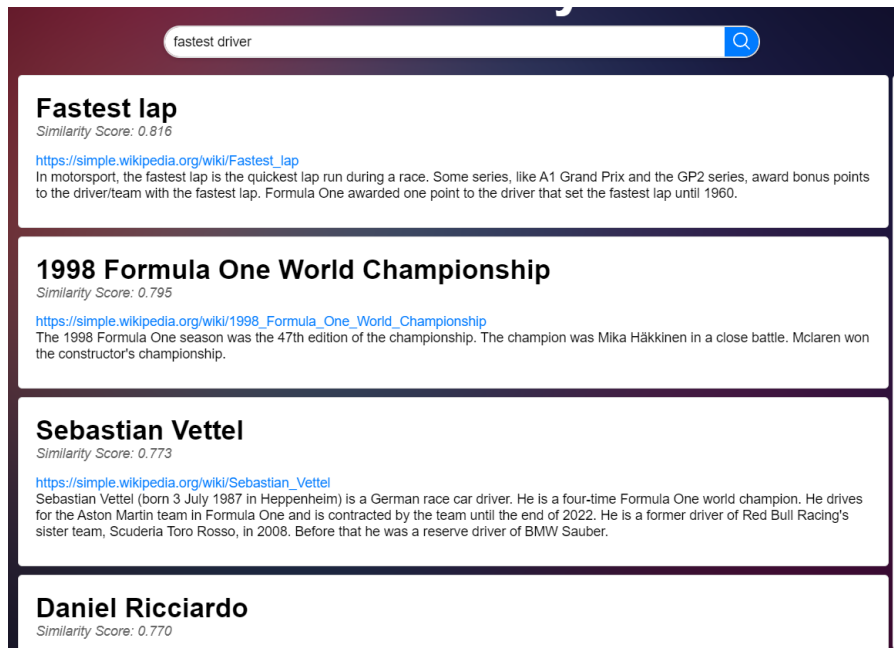# 3 Search results

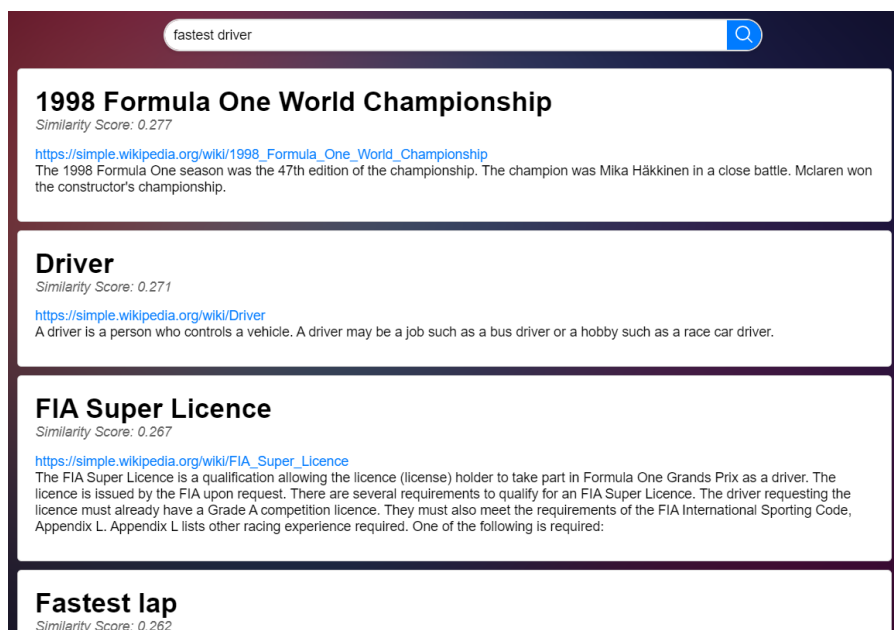## 3.1 fastest driver



Figure 3.1: no SVD



Figure 3.2: SVD, k=5000

Figure 3.3: SVD, k=1000



Figure 3.4: SVD, k=500

## 3.2 bacteria killer



Figure 3.5: no SVD



Figure 3.6: SVD, k=5000

Figure 3.7: SVD, k=1000



Figure 3.8: SVD, k=500

## 3.3 american singer
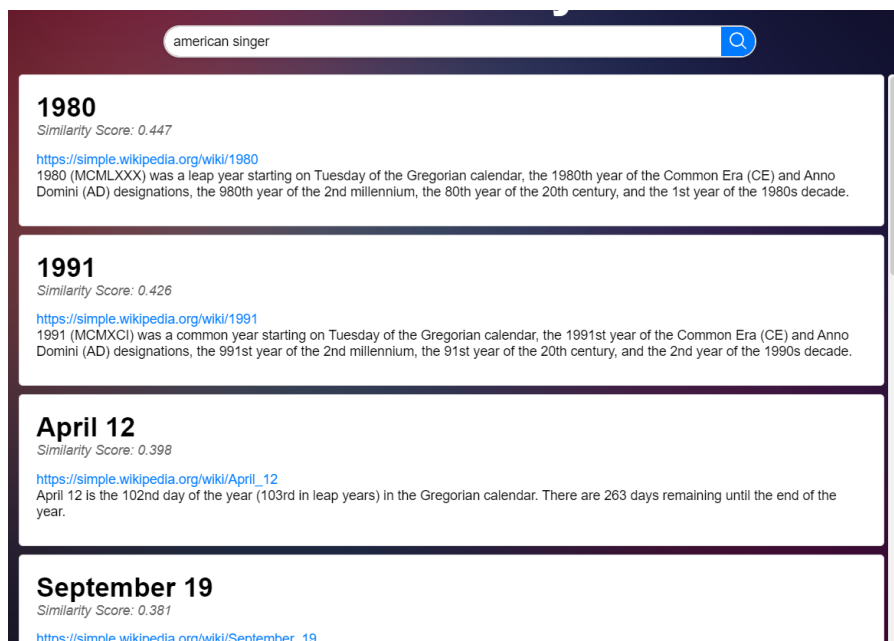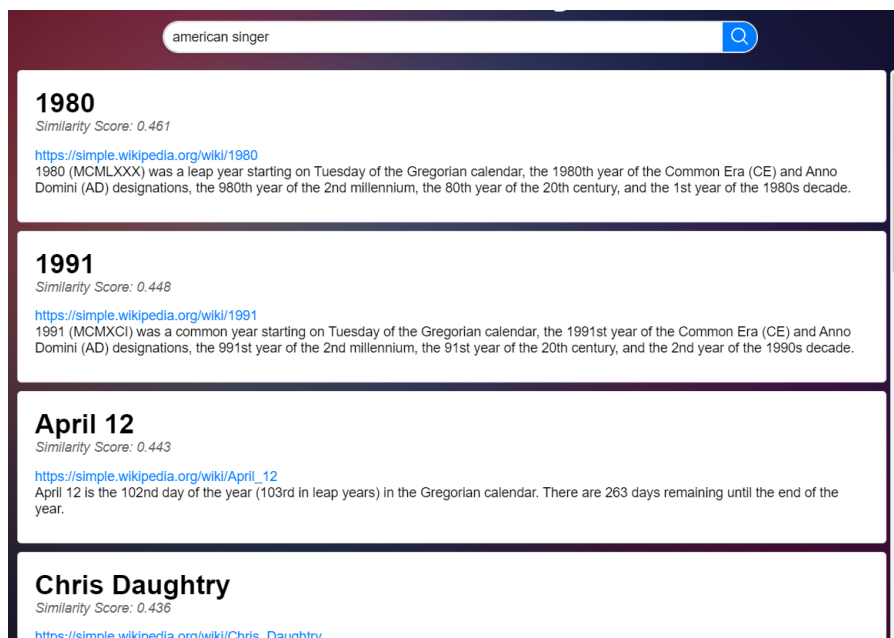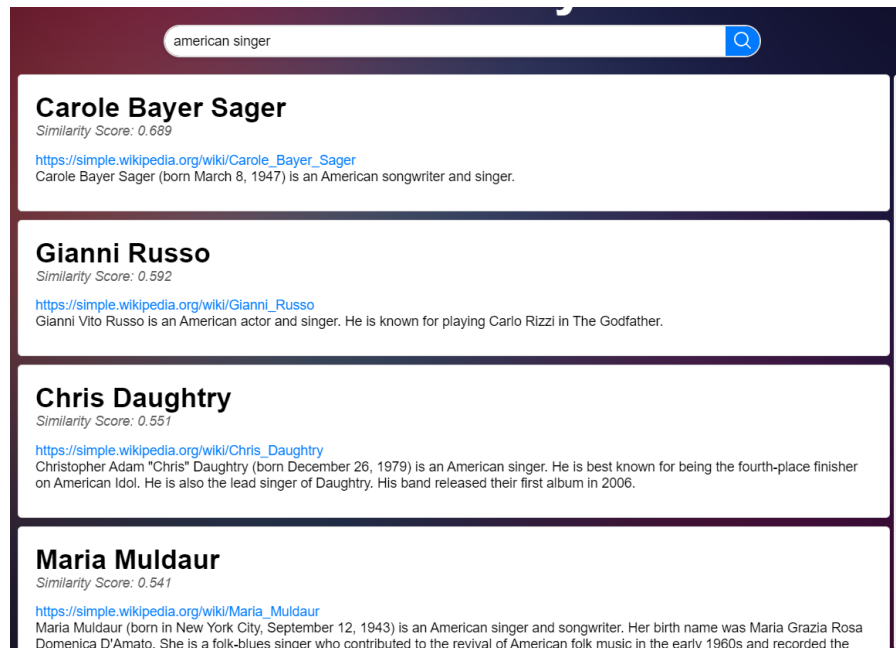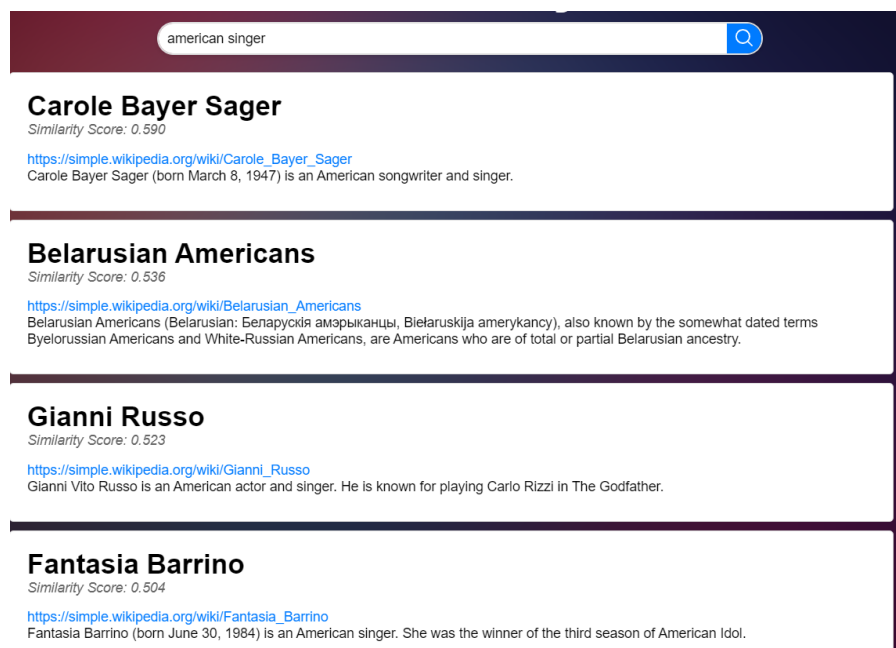


Figure 3.9: no SVD



Figure 3.10: SVD, k=5000

Figure 3.11: SVD, k=1000



Figure 3.12: SVD, k=500

## 3.4 summary

Testing showed that when we use SVD for searching, especially with high values of k, the results are pretty similar to searches without SVD. Sometimes, though, SVD can understand

the search better. But if we use low values of k, the results are often not very related to what we're searching for. In my opinion, the k value of 1000 performed the best.