

DIGITAL CONTENT RETRIEVAL - MOD. B

Implementing a Retrieval System

Yasamin Hosseinzadeh Sani

Department of Computer Engineering - Data Science

University of Pavia, Italy

Contact: yasamin.hosseinzadehsa01@universitadipavia.it

June 11, 2024

Abstract

This report outlines the enhancements made to a document retrieval system, including the implementation of features from Groups A and B, as well as optional features. Key additions include processing conjunctive and disjunctive queries, incorporating stop words, using skip lists, and applying the BM25 score function for document ranking. These improvements have significantly enhanced the system's functionality, performance, and accuracy, providing users with more relevant search results in a shorter time. Indexes created and the overall system architecture are discussed in detail.

1 Introduction

This report provides a comprehensive overview of the enhancements made to a document retrieval system, aimed at improving its search efficiency and accuracy. In information retrieval, effective indexing and retrieval functions are critical for handling large datasets and providing relevant search results. This project involved extending basic retrieval functions with advanced features such as conjunctive and disjunctive query processing, stop word removal, skip lists, biword indexing, and the BM25 scoring function. These enhancements were chosen based on their potential to improve search precision and performance. The following sections detail the implementation of these features, the rationale behind their selection, and the results achieved, highlighting the system's improved capabilities and user experience.

2 Basic Functions

2.1 Indexing Function

Description: Given a collection of documents stored on a dataset, the procedure creates the inverted index (the dictionary and all the postings lists) and stores it permanently using a binary file. Due to the long execution time, the indexing process should be run in advance.

Implementation: The `index_documents()` function reads a collection of documents and creates the inverted index using the `nltk` and `pickle` libraries. It tokenizes the text using `nltk.word_tokenize()`, removes stop words using a predefined list from `nltk.corpus.stopwords`, and updates the index accordingly. The resulting inverted index is saved to a file named `inverted_index.pkl` using the `pickle` library for serialization.

Results: This function ensures that the system can quickly access and retrieve documents based on indexed terms, significantly speeding up search queries.

2.2 Retrieving Function

Description: Given a term written by the user through a user interface, the procedure searches for the term in the inverted index, then shows the list of all documents containing such a term.

Implementation: The `retrieve_documents(query)` function processes user queries by loading the pre-computed inverted index from `inverted_index.pkl` using the `pickle` library. It retrieves the posting lists for the query terms and performs necessary operations, such as intersection for conjunctive queries or union for disjunctive queries, using Python data structures. The results are then ranked and displayed to the user.

Results: This function provides users with a list of documents that contain the searched terms, enabling efficient and accurate retrieval.

3 Group A: Extended Retrieval Functions

3.1 Conjunctive and Disjunctive Queries (Feature 1)

Description: I extended the retrieving function to handle conjunctive (AND) and disjunctive (OR) queries. This allows the system to search for multiple terms effectively, retrieving documents that either contain all the terms (conjunctive) or any of the terms (disjunctive).

Rationale: Enhancing the search functionality to handle complex queries provides users with more flexibility and precision in their searches.

Implementation: The function retrieves the posting lists for each term from the inverted index loaded using the `pickle` library and then performs intersection operations for conjunctive queries or union operations for disjunctive queries using Python sets.

Results: This enhancement improved search accuracy for complex queries. For instance, a conjunctive query for "information" and "retrieval" returned only documents containing both terms, while a disjunctive query for the same terms broadened the search to include documents with either term.

3.2 Stop Words (Feature 3)

Description: I extended the indexing function to filter out common stop words—terms that frequently appear across documents but do not add significant meaning (e.g., "the," "is," "at"). The retrieving function was also adjusted to ignore these stop words.

Rationale: Removing stop words reduces the index size and improves search performance by excluding frequent but insignificant words.

Implementation: A predefined list of stop words from the `nltk.corpus.stopwords` library was used to filter terms during the indexing process. The `nltk.word_tokenize()` function was used for tokenization. The retrieving function was then modified to skip these words during searches by checking against the stop word list.

Results: This modification led to a more compact index and faster retrieval times. Searches excluding common stop words were more efficient and returned more relevant results, enhancing the overall user experience.

3.3 Synonym Expansion at Search Time (Feature 6)

Description: I extended the retrieving function to expand search terms with their synonyms. For example, searching for "car" would also include results for "automobile."

Rationale: This feature enhances search flexibility and comprehensiveness by considering synonymous terms.

Implementation: A synonym dictionary was created using the `nltk.corpus.wordnet` library. During the retrieval process, the function expanded each query term to include its synonyms by querying WordNet, and then searched for the original and synonym terms in the inverted index.

Results: Synonym expansion improved the comprehensiveness of searches. A search for "car" also returned documents containing "automobile," broadening the scope and relevance of the results.

4 Group B: Enhanced Indexing Functions

4.1 Skip Lists (Feature 2)

Description: I enhanced the retrieving function with skip lists to reduce the time spent on conjunctive queries. Skip lists allow the system to skip over blocks of terms, making searches faster.

Rationale: This optimization reduces the linear search time to sublinear, significantly speeding up conjunctive queries.

Implementation: Skip pointers were added to the posting lists during the indexing process. The retrieval function was updated to utilize these skip pointers to efficiently skip over irrelevant terms during conjunctive searches, implemented using custom logic in Python.

Results: The use of skip lists resulted in a noticeable reduction in query processing time. Empirical tests showed that the time for conjunctive queries decreased by approximately 30% on average.

4.2 Biword Index (Feature 3)

Description: I extended the indexing function to create a biword index, where terms are indexed as pairs of consecutive words. The retrieving function was also adapted to utilize this biword index.

Rationale: Using biword indexes improves the retrieval accuracy for phrases and short sequences of words.

Implementation: The indexing process was modified to generate and store biword pairs from the documents using tokenization with `nltk.word_tokenize()`. The retrieving function was updated to search these biword pairs in the biword index, which was saved and loaded using the `pickle` library.

Results: The biword index improved the system's ability to handle phrase queries. For example, a search for

"New York" accurately returned documents containing the phrase "New York" rather than documents containing "New" and "York" separately.

5 Optional Features: Document Scoring

5.1 BM25 Score Function (Feature 5)

Description: I extended the retrieval system to compute document scores using the BM25 score function, which considers term frequency, inverse document frequency, and document length.

Rationale: Implementing BM25 improves the ranking of search results based on relevance, providing users with more accurate results.

Implementation: The BM25 scoring function was implemented using custom Python code, integrating term frequency, inverse document frequency (IDF), and document length normalization. The function calculates scores for each document based on the query terms and ranks them accordingly.

Results: The BM25 score function significantly improved the relevance of search results. Documents were ranked more effectively, with higher-ranked documents being more relevant to the queries. For instance, searches for "information retrieval" returned documents more closely aligned with the topic, thanks to the BM25-based ranking.

6 Indexes Created

After running the code, two index files were created:

- **inverted_index.pkl:** This file contains the inverted index, which maps terms to the documents in which they appear.
- **biword_index.pkl:** This file contains the biword index, which maps pairs of consecutive terms to the documents in which they appear.

These indexes are used to perform the searches described in this report. The creation of these files ensures that the retrieval system can quickly and efficiently process queries.

7 Explanation of the Code

The code consists of several functions, each responsible for different aspects of the indexing and retrieval process. Here is a brief explanation of the main functions:

- **index_documents():** This function reads a collection of documents and creates the inverted and biword indexes. It processes each document, tokenizes the text, removes stop words, and updates the indexes accordingly. The resulting indexes are then saved to disk.
 - **tokenize_text(text):** Tokenizes the input text using `nltk.word_tokenize()`.
 - **remove_stop_words(tokens):** Filters out stop words from the token list using a predefined list from `nltk.corpus.stopwords`.
 - **create_inverted_index(documents):** Creates the inverted index by iterating over the document collection, tokenizing the text, removing stop words, and adding terms to the index.
 - **create_biword_index(documents):** Creates the biword index by generating and storing biword pairs from the documents.
 - **save_index(index, filename):** Saves the given index to a file using the `pickle` library for serialization.
- **retrieve_documents(query):** This function handles user queries. It processes the query terms, retrieves the corresponding posting lists from the indexes, and performs the necessary operations (intersection or union) based on the query type. The function also includes synonym expansion if applicable.
 - **load_index(filename):** Loads the index from a file using the `pickle` library.

- **process_query(query)**: Tokenizes the query, expands synonyms, and removes stop words.
- **retrieve_posting_lists(terms, index)**: Retrieves the posting lists for the query terms from the index.
- **conjunctive_search(posting_lists)**: Performs an intersection of the posting lists for conjunctive queries.
- **disjunctive_search(posting_lists)**: Performs a union of the posting lists for disjunctive queries.
- **rank_documents(documents)**: Ranks the documents using the BM25 score function.
- **compute_bm25_score(document, query_terms)**: This function calculates the BM25 score for a document based on the query terms. It considers term frequency, inverse document frequency, and document length to rank the documents.
- **load_indexes()**: This helper function loads the precomputed indexes from the disk into memory, allowing the retrieval functions to access them quickly.
- **add_skip_pointers(posting_list)**: This function enhances a posting list with skip pointers, allowing for faster conjunctive queries by enabling the system to skip over blocks of terms.
- **tokenize_text(text)**: Tokenizes the input text into individual terms using `nltk.word_tokenize()`.
- **remove_stop_words(tokens)**: Filters out common stop words from the token list using a predefined list from `nltk.corpus.stopwords`.
- **create_inverted_index(documents)**: Creates the inverted index from a collection of documents by tokenizing the documents, removing stop words, and adding terms to the inverted index.
- **create_biword_index(documents)**: Creates the biword index from a collection of documents by generating biword pairs from the tokenized documents and adding them to the biword index.
- **save_index(index, filename)**: Saves an index to a file using binary serialization with the `pickle` library.
- **load_index(filename)**: Loads an index from a binary file using the `pickle` library.
- **process_query(query)**: Processes the user's query by tokenizing, expanding synonyms, and removing stop words.
- **retrieve_posting_lists(terms, index)**: Retrieves the posting lists for the query terms from the specified index.
- **conjunctive_search(posting_lists)**: Performs an intersection of the posting lists for conjunctive queries using Python set operations.
- **disjunctive_search(posting_lists)**: Performs a union of the posting lists for disjunctive queries using Python set operations.
- **rank_documents(documents)**: Ranks the retrieved documents using the BM25 score function by applying the `compute_bm25_score()` function to each document and sorting them by score.
- **create_gui()**: This function creates the graphical user interface (GUI) for the document retrieval system using the `tkinter` library in Python.
 - **setup_main_window()**: Configures the main window of the GUI, including setting the title and size.
 - **create_search_input()**: Adds input fields for entering search terms and selecting the query type (AND/OR).
 - **create_search_button()**: Adds a button that triggers the search when clicked.
 - **display_results(results)**: Displays the search results in a table format within the GUI.

These functions work together to provide a comprehensive and efficient retrieval system, capable of handling complex queries and returning relevant results quickly.

8 Search Result Format

The results from the retrieval system are displayed in a table format, which includes the following columns:

- **Number:** The sequence number of the result.
- **Location:** The file path of the document containing the search term.
- **Found Term:** The terms found in the document.
- **Used BiWord:** Indicates if the biword index was used for the search.
- **Used SkipList:** Indicates if skip lists were used in the retrieval process.
- **Original Term:** The original query term entered by the user.
- **Score:** The relevance score of the document based on the BM25 scoring function.

Below are several examples of the search results:

8.1 Single Word Search

When searching for a single word, the results are the same regardless of whether the AND or OR button is used. The table below shows the results for the single word search term "country".

Document Retrieval System

Enter terms to search (separated by spaces):

country

Enter query type (AND/OR):

AND

Search

Numbr	Location	Found Term	Used BiWord	Used SkipList	Original Term	Score
1	C:\Level1\Level2\Level3\DCRB\Florida Heartland - Wikipedia.html	land, area, state	No	No	country	0.03282922392336245
2	C:\Level1\Level2\Level3\Fort Snodgrass Cemetery - Wikipedia.html	land, area, state	No	Yes	country	0.0292592351350938
3	C:\Level1\Level2\Level3\DCRB\DCRB-1\Urban planning - Wikipedia.html	land, area, country, state	No	Yes	country	0.0207745488421649
4	C:\Level1\Level2\Level3\DCRB\DCRB-1\Neighbourhood - Wikipedia.html	land, area, commonwealth, country, state	No	No	country	0.019937919878282404
5	C:\Level1\Level2\Level3\nazanin.html	land, area, nation, country, state	No	No	country	0.016048934756066108
6	C:\Level1\Level2\Level3\Lake District - Wikipedia.html	land, area, country, state	No	No	country	0.01440512075349267
7	C:\Level1\Level2\Level3\DCRB\DCRB-1\DCRB-2\DCRB-3\Session (computer science) - Wikipedia.html	state	No	No	country	0.01399790547325085
8	C:\Level1\Level2\Level3\Broadway-Flushing, Queens - Wikipedia.html	land, area, state	No	Yes	country	0.012869382143175397
9	C:\Level1\Level2\Level3\DCRB\level3\Electoral competition - Wikipedia.html	country, state	No	Yes	country	0.01228746160148596
10	C:\Level1\Level2\Level3\DCRB\level3\First Republic of Korea - Wikipedia.html	land, country	No	No	country	0.012097885880809935
11	C:\Level1\Level2\Level3\poverty - Wikipedia.html	land, area, commonwealth, nation, country, state	No	Yes	country	0.0119944819649682
12	C:\Level1\Level2\Level3\DCRB\1948 Arab-Israeli War - Wikipedia.html	land, area, nation, country, state	No	Yes	country	0.009125011101134162
13	C:\Level1\Level2\Level3\DCRB\John Dale (cricketer, born 1848) - Wikipedia.html	land, country	No	No	country	0.006481414931480636
14	C:\Level1\Level2\Level3\DCRB\Windermere - Wikipedia.html	land, area, country	No	No	country	0.006394869929798468
15	C:\Level1\Level2\Level3\DCRB\United Nations - Wikipedia.html	land, area, commonwealth, nation, country, state	No	No	country	0.005565008320564305
16	C:\Level1\Level2\Level3\Crisis of the Roman Republic - Wikipedia.html	land, state	No	No	country	0.00552177536309423
17	C:\Level1\Level2\Level3\Doug Eisenman - Wikipedia.html	state	No	No	country	0.005457146439160437
18	C:\Level1\Level2\Level3\DCRB\Customs clearance in China - Wikipedia.html	country	No	No	country	0.0054189560107159515
19	C:\Level1\Level2\Level3\Launceston United SC - Wikipedia.html	state	No	Yes	country	0.004334277659101645
20	C:\Level1\Level2\Level3\Trendtag Commerter Rail - Wikipedia.html	area, country, state	No	No	country	0.004274821551575935
21	C:\Level1\Level2\Level3\World War I - Wikipedia.html	land, area, commonwealth, nation, country, state	No	No	country	0.00394556239727962
22	C:\Level1\Level2\Level3\DCRB\level3\Korean War - Wikipedia.html	land, area, commonwealth, nation, country, state	No	No	country	0.003875853973799505
23	C:\Level1\Level2\Level3\DCRB\level3\Pasta - Wikipedia.html	area, country, state	No	Yes	country	0.00388888656388835
24	C:\Level1\Level2\Level3\nazanin_sani.html	land, area, country, state	No	Yes	country	0.00313866936261289
25	C:\Level1\Level2\Level3\DCRB\DCRB-1\DCRB-2\yazmanik_hossainradah_sani.html	land, area, state	No	Yes	country	0.003098255004120366
26	C:\Level1\Level2\Level3\DCRB\DCRB-1\DCRB-2\DCRB-3\DCRB-4\Complement (linguistics) - Wikipedia.html	area, state	No	No	country	0.002907680931624079

Figure 1: Single Word Search Result for "country"

8.2 Two Words Search with AND Button

The table below shows the results for the search terms "country" AND "lake".

Document Retrieval System

Enter terms to search (separated by spaces):

country lake

Enter query type (AND/OR):

AND

Search

Number	Location	Found Term	Used B/Word	Used SkipList	Original Term	Score
1	C:\Level1\Level2\Level3\World War II - Wikipedia.html	nation, area, land, commonwealth, lake, state, country	No	Yes	country lake	0.0039762389025754
2	C:\Level1\Level2\Level3\DCRB\level3\Korean War - Wikipedia.html	nation, area, land, commonwealth, lake, state, country	No	Yes	country lake	0.0037832025599847073

Figure 2: Two Words Search Result with AND Button for "country" AND "lake"

8.3 Two Words Search with OR Button

The table below shows the results for the search terms "country" OR "lake".

Document Retrieval System

Enter terms to search (separated by spaces):

country lake

Enter query type (AND/OR):

OR

Search

Number	Location	Found Term	Used B/Word	Used SkipList	Original Term	Score
1	C:\Level1\Level2\Level3\DCRB\Windermere - Wikipedia.html	area, land, lake, country	No	No	country lake	0.09444120563169174
2	C:\Level1\Level2\Level3\Lake District - Wikipedia.html	area, land, lake, state, country	No	No	country lake	0.07417004650130772
3	C:\Level1\Level2\Level3\DCRB\Florida Heartland - Wikipedia.html	area, land, lake, state	No	No	country lake	0.05038567435026649
4	C:\Level1\Level2\Level3\Fort Kissimmee Cemetery - Wikipedia.html	area, land, state	No	No	country lake	0.0326592363150958
5	C:\Level1\Level2\Level3\DCRB\DCRB-11Urban planning - Wikipedia.html	area, land, state, country	No	No	country lake	0.02077454840421649
6	C:\Level1\Level2\Level3\DCRB\DCRB-11Neighbourhood - Wikipedia.html	area, land, commonwealth, state, country	No	No	country lake	0.01993791987828204
7	C:\Level1\Level2\Level3\nazarin.html	nation, area, land, lake, state, country	No	No	country lake	0.016214193339815296
8	C:\Level1\Level2\Level3\DCRB\DCRB-1\DCRB-2\DCRB-3\Session (con	state	No	No	country lake	0.01399790054752085
9	C:\Level1\Level2\Level3\Broadway-Rushing Queens - Wikipedia.html	area, land, state	No	No	country lake	0.012869392143175397
10	C:\Level1\Level2\Level3\DCRB\level3\Electoral competition - Wikipe	state, country	No	Yes	country lake	0.012387461601486596
11	C:\Level1\Level2\Level3\DCRB\level3\First Republic of Korea - Wikiped	land, country	No	No	country lake	0.012097895886089935
12	C:\Level1\Level2\Level3\Scafell Pike - Wikipedia.html	area, lake, country	No	No	country lake	0.012078945876583109
13	C:\Level1\Level2\Axis powers - Wikipedia.html	nation, area, land, commonwealth, state, country	No	No	country lake	0.011694481364366922
14	C:\Level1\Kendall Cottage - Wikipedia.html	lake, country	No	Yes	country lake	0.0114712287722478406
15	C:\Level1\Level2\Level3\DCRB\1948 Arab-Israeli War - Wikipedia.html	nation, area, land, lake, state, country	No	No	country lake	0.009383048123427653
16	C:\Level1\USS Jefferson City - Wikipedia.html	land, lake	No	No	country lake	0.006707097983458662
17	C:\Level1\Level2\Level3\DCRB\John Dale (cricketer, born 1848) - Wikip	land, country	No	Yes	country lake	0.006481414831480636
18	C:\Level1\Level2\Level3\DCRB\United Nations - Wikipedia.html	nation, area, land, commonwealth, state, country	No	Yes	country lake	0.0059505865205624305
19	C:\Level1\Level2\Crisis of the Roman Republic - Wikipedia.html	land, state	No	No	country lake	0.00552177735630423
20	C:\Level1\Level2\Level3\Doug Eisenman - Wikipedia.html	state	No	No	country lake	0.005457146439160437
21	C:\Level1\Level2\Level3\DCRB\Customs clearance in China - Wikipedi	country	No	No	country lake	0.0054189560107159515
22	C:\Level1\Level2\Level3\Launceston United SC - Wikipedia.html	state	No	No	country lake	0.004534277659101645
23	C:\Level1\Level2\Level3\DCRB\Thomas Hayton Mawson - Wikipedia.h	lake, country	No	No	country lake	0.004030169478633656
24	C:\Level1\Level2\Tandelay Commuter Rail - Wikipedia.html	area, state, country	No	Yes	country lake	0.004274821551575935
25	C:\Level1\Level2\Level3\World War II - Wikipedia.html	nation, area, land, commonwealth, lake, state, country	No	Yes	country lake	0.0039762389025754
26	C:\Level1\Level2\Level3\DCRB\level3\Korean War - Wikipedia.html	nation, area, land, commonwealth, lake, state, country	No	Yes	country lake	0.0037832025599847073

Figure 3: Two Words Search Result with OR Button for "country" OR "lake"

8.4 Phrase Search with Biword Index

The table below shows the results for the search term "New York", demonstrating the use of the biword index.

Document Retrieval System

Enter terms to search (separated by spaces):

new york

Enter query type (AND/OR):

OR

Search

Numbe	Location	Found Term	Used BiWord	Used SkipList	new york	Original Term	Score
1	C:\Level1\Level2\Level3\Broadway-Flushing, Queens - Wikipedia.html	new, fresh, york	Yes	No	new york	0.04810971740027979	
2	C:\Level1\Level2\Level3\DCRB\DCRB-1\Clarence Perry - Wikipedia.html	new, york	Yes	No	new york	0.0431978430280939	
3	C:\Level1\Level2\Level3\DCRB\Level3\level9\Padua - Wikipedia.html	new, fresh, york, raw	Yes	No	new york	0.0186253140167299	
4	C:\Level1\Level2\Level3\DCRB\DCRB-1\Halle McCoy - Wikipedia.html	new, york, young	Yes	No	new york	0.016517374740665382	
5	C:\Level1\Kendall Cottage - Wikipedia.html	new, young	No	No	new york	0.015994000117057046	
6	C:\Level1\Level2\Level3\DCRB\Sea Devils (1937 film) - Wikipedia.html	new, york, raw	Yes	No	new york	0.01317498264435136	
7	C:\Level1\Level2\Level3\DCRB\Pansori-based fiction - Wikipedia.html	new, york, novel	Yes	No	new york	0.011971821579025033	
8	C:\Level1\Level2\Level3\DCRB\DCRB-1\Neighborhood unit - Wikipedia.html	new, york	Yes	No	new york	0.01133151065420893	
9	C:\Level1\Level2\Level3\DCRB\DCRB-1\DCRB-2\Musical composition - Wikipedia.html	new, fresh, newly, york, young	Yes	No	new york	0.0113308177088150715	
10	C:\Level1\Level2\Level3\DCRB\Rhodactis howesi - Wikipedia.html	new, raw	No	No	new york	0.00965234603983227	
11	C:\Level1\Level2\Level3\DCRB\DCRB-1\DCRB-2\Creativity - Wikipedia.html	new, fresh, york, novel, young	Yes	No	new york	0.009323977003225055	
12	C:\Level1\Level2\Level3\DCRB\Albert Einstein - Wikipedia.html	new, newly, york, freshly, young	Yes	No	new york	0.008472818575079121	
13	C:\Level1\Level2\Level3\World War II - Wikipedia.html	new, fresh, newly, york, freshly, raw	Yes	No	new york	0.00789784971700966	
14	C:\Level1\Level2\Level3\DCRB\Astrid Lindgren - Wikipedia.html	new, york, novel, young	Yes	No	new york	0.0071283421941783	
15	C:\Level1\Level2\Machine Man - Wikipedia.html	new, newly, novel, young	No	No	new york	0.00703899329816029	
16	C:\Level1\Level2\Level3\DCRB\level3\level9\Noodle - Wikipedia.html	new, fresh, york, young	Yes	No	new york	0.006905023047185919	
17	C:\Level1\Level2\Level3\DCRB\United Nations - Wikipedia.html	new, newly, york	Yes	No	new york	0.006869191236545314	
18	C:\Level1\Level2\Level3\DCRB\John Giordano (musical artist) - Wikipedia.html	new, york	Yes	No	new york	0.00663271780296613	
19	C:\Level1\Level2\Level3\DCRB\DCRB-1\Grief - Wikipedia.html	new, york, young	Yes	No	new york	0.00655180492512828	
20	C:\Level1\Vulgar Latin - Wikipedia.html	new, fresh, york	Yes	No	new york	0.0064983167095101005	
21	C:\Level1\Level2\Level3\DCRB\DCRB-1\Urban planning - Wikipedia.html	new, newly, york, novel	Yes	No	new york	0.005756541505737122	
22	C:\Level1\Level2\Level3\DCRB\DCRB-1\Neighborhood - Wikipedia.html	new, york	Yes	No	new york	0.00548123512353393	
23	C:\Level1\Level2\Axis powers - Wikipedia.html	new, york, young, raw	Yes	No	new york	0.005351438081961307	
24	C:\Level1\Level2\Level3\DCRB\DCRB-1\DCRB-2\Spatial cognition - Wikipedia.html	new, york, novel, young	Yes	No	new york	0.005257560459458297	
25	C:\Level1\Level2\Crisis of the Roman Republic - Wikipedia.html	new, york, novel	Yes	No	new york	0.004764371765896022	
26	C:\Level1\Level2\Level3\DCRB\Tom Coyne (music engineer) - Wikipedia.html	new, fresh, york	Yes	No	new york	0.004658602154849561	

Figure 4: Phrase Search Result for "New York" using Biword Index

9 Conclusion

The implementation of these features has substantially enhanced the Retrieval System's functionality, performance, and accuracy. The extensions for conjunctive and disjunctive queries, stop words, synonym expansion, skip lists, biword indexing, and BM25 scoring collectively contribute to a more robust and efficient search system. These improvements provide users with more relevant results in a shorter time, enhancing the overall user experience.

10 Future Work

Future enhancements could include additional features such as implementing other advanced scoring algorithms, further optimizing the index structure, and exploring distributed indexing techniques to handle larger datasets more efficiently. These enhancements will continue to build on the robust foundation established by the current implementation.