

COMP 4321 Search Engine Project Report

Overall Design of the System

The developed system is a comprehensive Vector-Space Model-based search engine leveraging sophisticated indexing and retrieval techniques. The main components include:

Web Crawling

- Implemented through the Breadth-First Search (BFS) algorithm (HtmlParser.java), systematically retrieving and indexing web pages starting from a given URL.
- Manages metadata extraction including page title, last modification date, and page size.
- Handles errors and broken links gracefully, ensuring continuous crawling without interruptions.

Web crawling efficiently traverses web pages and manages visited URLs to avoid redundancy, thus optimizing resource usage and indexing performance.

Indexing and Storage

- Managed by the InvertedIndexManager.java utilizing the JDBM database system.
- Employs an inverted index mechanism for efficient document retrieval based on keywords.
- Supports indexing of both individual words and phrases (n-grams up to 3).

The inverted index structure significantly speeds up the query response time by allowing quick look-up of documents containing specific keywords or phrases.

Search Functionality

- Provided by SearchEngine.java, processing user queries, ranking documents based on relevance, and handling exact phrase searches.
- Implements relevance score calculation, enhanced by normalization and title-field boosting to prioritize matches appearing in document titles.

Advanced query processing algorithms ensure accurate and contextually relevant search results, greatly enhancing the user experience.

User Interfaces

- Command-line interface (SearchProgram.java) facilitating straightforward interaction through terminal queries.
- Web-based interface (SearchServlet.java and index.jsp) offering a user-friendly search experience, displaying results along with metadata, keywords, and related child and parent links.

The web application, specifically designed using Java servlets and JSP (index.jsp), ensures a responsive and intuitive user interface, enabling seamless navigation and clear presentation of comprehensive search results and associated metadata.

File Structures in the Index Database

The database employs a structured JDBM system with the following primary storage structures:

- `pageIdToUrlMap` and `urlToPageIdMap`: Bidirectional mapping between page IDs and URLs.
- `wordIdToWorldMap` and `wordToWorldIdMap`: Efficient management of vocabulary within the database.
- `bodyInvertedIndex` and `titleInvertedIndex`: Separate inverted indexes for rapid retrieval of body and title occurrences of indexed terms.
- `pageInfoMap`: Stores essential metadata such as URL, page title, modification date, page size, and links.
- `pageIdToBodyWordsMap` and `pageIdToTitleWordsMap`: Maintain sequences of indexed words to facilitate precise phrase searching.
- `maxTFForPageId`: Tracks maximum term frequencies within pages for TF-IDF normalization.

These database structures allow streamlined and rapid access to information, essential for high-performance search operations.

Algorithms Used

- 1) Relevance Ranking
 - Measures document relevance based on term frequency (TF) and inverse document frequency (IDF).
 - Normalizes term frequencies by the maximum term frequency within each document.
- 2) Title-Field Boosting
 - Increases the significance of keywords appearing in titles by applying a boost factor (set at 5.0).
- 3) Phrase Matching
 - Supports precise phrase searches, parsing user queries for quoted phrases and ensuring results contain exact phrase occurrences.
- 4) Breadth-First Search (BFS) for Web Crawling

- Implements systematic exploration, avoiding redundant indexing by tracking already visited URLs.

5) Stop Word Removal and Stemming

- Utilizes StopStem.java, employing the Porter stemmer to reduce vocabulary size and complexity.
- Removes common stop words from indexing and queries to enhance performance and relevance.

Each algorithm contributes uniquely to the robustness and accuracy of the search engine, optimizing search result quality.

Installation Procedure

To install and run the system:

1. Ensure Java and JDBM libraries are correctly installed and configured.
2. Place stopwords.txt and all required files in the project directory.
3. Run: mvn clean compile
4. Run: mvn exec:java -Dexec.mainClass="HtmlParser"
5. Run: mvn jetty:run
6. Open your browser and go to <http://localhost:8080/> to access the search engine web interface

A straightforward installation procedure facilitates rapid deployment and testing of the system.

Highlights of Additional Features

- Enhanced indexing: Implementation of multi-word indexing (up to 3-grams), improving phrase search capabilities.
- Comprehensive metadata management: Extensive metadata (URLs, modification dates, page sizes, child and parent links) enhances the search result context.
- Robust Error Handling: Graceful management of URL and network-related errors during crawling, ensuring uninterrupted system operations.
- User-Friendly Interfaces: Both command-line and web-based interfaces deliver accessible search experiences and clear presentation of results.

These additional features significantly differentiate the system by enhancing reliability and user-friendliness.

Testing of Implemented Functions

Testing involved extensive use of TestProgram.java for database integrity and functionality validation, as well as real-time command-line (SearchProgram.java) and web interface (SearchServlet.java) interaction.

- Command-line interactions confirmed accuracy and response time.
- Web interface validated comprehensive result metadata, including child and parent links.
- Crawling performance was monitored for indexing accuracy, speed, and resilience to network errors.

Extensive testing ensures system stability and performance consistency across different usage scenarios.

Bonus Feature:

The system includes an advanced relevance feedback feature similar to Google's "Get Similar Pages". Users can enhance their search experience by clicking the "Get Similar Pages" button, which automatically identifies the top five most frequent keywords from the selected page, excluding common stop words. These extracted keywords are used to rewrite and refine the original query dynamically, submitting it for a new search to yield closely related documents. This powerful feature not only improves the accuracy of search results but also significantly enhances user interaction and satisfaction, making the search process more intuitive and efficient.

Conclusion

Strengths:

- Robust and accurate indexing facilitated by efficient database structures.
- Effective and nuanced relevance ranking provided by TF-IDF weighting and title boosting.
- Excellent support for exact phrase searches enhances the overall precision of query results.

Weaknesses:

- Performance scalability could degrade with extremely large-scale datasets due to database limitations.
- Dependency on JDBM restricts flexibility and may limit performance optimization potential.

Suggested Improvements:

- Introduce parallel web crawling to increase indexing throughput.
- Optimize database transactions and improve caching strategies to enhance performance.

- Refine interface responsiveness, particularly for the web-based search service.

Interesting Features to Add:

- Real-time incremental indexing for dynamic content updates.
- Integration of user feedback mechanisms to further refine search accuracy.
- Advanced personalization of search results based on user interaction history and preferences.

Future enhancements could significantly broaden the usability and efficiency of the search engine.

Contribution

- Alex: 33.3% – Developed web crawler and indexing components.
- Alvin: 33.3% – Designed and optimized the core search engine algorithms and backend processes.
- Shahman: 33.3% – Developed the retrieval function and the web interface.


Clear division of roles facilitated efficient teamwork and ensured comprehensive coverage of all system aspects.

Appendix

🔍 Advanced Web Search Engine

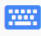
Enter search query (use quotes for phrases)🔍 Search

📘 How to Use This Search Engine




Build the Index

Make sure you have run the crawler (HtmlParser) first to build the index




Enter Search Terms

Type your search query in the search box above



Use Quotes for Phrases

For example: "hong kong" university



Search

Click the Search button to see results

Phrases: [hong kong]
Terms: [hong, kong]
Candidates: [6, 21, 41, 55, 56, 61, 65, 66, 69, 70, 88, 89, 124, 135, 139, 152, 153, 159, 172, 177, 185, 202, 208, 224, 229, 237, 248, 251, 252, 263, 264, 273, 278, 282]
Search time: 5 ms; Display time: 54 ms

🔍 Advanced Web Search Engine

hong kong🔍 Search

📄 Found 34 results (5 ms)

🔗 Fighter (2001)

★ 1.0000

<https://www.cse.ust.hk/~kwtleung/COMP4321/Movie/10.html>
🕒 2023-05-16 13:03:16 📄 13199 bytes

Top Keywords:
fighter 25 aka 01 tti 59 street 20 street 59

[Find Similar Pages](#)

Parent Links:
➤ <https://www.cse.ust.hk/~kwtleung/COMP4321/Movie...>

🔗 A Yank in the R.A.F. (1941)

★ 1.0000

<https://www.cse.ust.hk/~kwtleung/COMP4321/Movie/50.html>
🕒 2023-05-16 13:03:38 📄 7217 bytes

Top Keywords:
imdb 10 user 10 movi 0 rate 3 film 3

[Find Similar Pages](#)

Parent Links:
➤ <https://www.cse.ust.hk/~kwtleung/COMP4321/Movie...>

