# Inverted Index Datamart Structure Experiments

Luis Guillén Servera

October 20, 2024

## Abstract

This project focuses on the implementation and evaluation of a hashmap-based inverted index for efficient information retrieval from large text datasets. The system is designed to crawl books from the Gutenberg Project, process the text, and construct an inverted index to enable fast keyword searches. Using the hashmap data structure, the project aims to optimize both indexing speed and query efficiency. Performance benchmarking was conducted to assess the system's behavior under varying dataset sizes and query complexities. The results demonstrate the effectiveness of the hashmap-based inverted index for large-scale text indexing, providing insights into its scalability and performance in real-world applications.

## 1 Introduction

The exponential growth of digital content has led to an increasing demand for efficient information retrieval systems. As vast amounts of textual data become available online, especially from public sources such as the Gutenberg Project, the need for scalable and effective search engines becomes critical. These systems must handle large datasets, provide fast keyword-based searching, and scale efficiently as the size of the corpus increases. A key component in such systems is the inverted index, which maps words to their occurrences in documents, enabling rapid retrieval of relevant content.

Inverted indexing is widely used in information retrieval systems due to its ability to facilitate fast searches over large text corpora. The choice of data structures for building these indices plays a crucial role in determining the system's performance. Hashmaps, which provide constant-time average lookups, have emerged as a popular solution for constructing inverted indices due to their efficiency and simplicity. Hashmaps allow quick access to indexed terms and their associated document lists, making them ideal for search engines handling large volumes of data.

This project implements a hashmap-based inverted index for a search engine that crawls, processes, and indexes books from the Gutenberg Project. The system is designed to be efficient in both the indexing and query phases, ensuring that it can handle large datasets without sacrificing performance. Additionally, the project incorporates benchmarking techniques to evaluate the performance of the hashmap-based approach in terms of indexing speed, memory usage, and query efficiency.

The following sections outline the design of the system, the benchmarking methodology, and the results obtained from testing the inverted index with real-world data. This analysis provides insights into the scalability and performance of hashmap-based indices and highlights their suitability for large-scale information retrieval applications.

## 2 Methodology

This section describes the methodology employed for developing, implementing, and evaluating the hashmap-based inverted index. The primary goal is to ensure efficient text indexing and retrieval from a large collection of books. The system's performance is measured in terms of indexing speed, mem-

ory usage, and query response time using a real-world dataset obtained from the Gutenberg Project.

## 2.1 Data Collection

Books were crawled and downloaded from the Project Gutenberg website, a repository of over 60,000 free eBooks. The system's crawler was designed to fetch a specified number of top books from the website. Once downloaded, the books were stored in plain text format within a dedicated storage folder. Each book was processed to extract its metadata (such as title, author, and release date) and the textual content, which served as input for the indexing process.

## 2.2 Text Preprocessing

Before indexing, the text of each book underwent preprocessing to ensure uniformity and enhance the quality of the data. The preprocessing steps included:

- **Tokenization**: The text was split into individual words (tokens).

- **Lowercasing**: All tokens were converted to lowercase to ensure case-insensitive searches.

- **Stopword Removal**: Commonly used words with little semantic meaning, such as "the" and "is," were removed using the NLTK stopword library.

- **Punctuation Removal**: All punctuation marks were stripped from the text to focus solely on meaningful terms.

These steps ensured that only relevant and searchable terms were indexed.

## 2.3 Inverted Index Construction

The core component of the system is the hashmap-based inverted index, which was constructed as follows:

- For each book, the preprocessed text was processed one word at a time.

- Each word (term) was added as a key to the hashmap, with the corresponding value being a list of book IDs where the term appeared.

- If a term already existed in the hashmap, the current book ID was appended to the list of book IDs associated with that term.

- The index construction process ensured that terms were efficiently mapped to the books in which they appeared, facilitating rapid lookups during search queries.

This hashmap-based structure allows for constant-time average lookups for search terms, providing fast query responses even as the dataset grows.

## 2.4 Benchmarking and Performance Evaluation

To evaluate the performance of the inverted index, a series of benchmarking experiments were conducted. These experiments focused on two main aspects:

- **Indexing Performance**: The time taken to index all books was measured. This involved recording the time from the start of the indexing process to its completion, including the time required for text preprocessing and hashmap construction.

- **Query Performance**: The time taken to retrieve the list of book IDs for a given search term was measured. Multiple queries were executed, and their response times were recorded and averaged.

For both types of performance measurements, the system was benchmarked under varying conditions, such as different dataset sizes and query complexities. The Python 'pytest-benchmark' library was used to conduct and report performance benchmarks.

## 2.5 Metrics for Evaluation

The system was evaluated based on the following key metrics:

- **Indexing Time**: The time required to build the inverted index for a collection of books.

- **Query Time**: The time required to search for a specific term and retrieve the list of books where the term appears.

- **Memory Usage**: The amount of memory consumed by the hashmap-based inverted index during both the indexing and query phases.

These metrics provide a comprehensive evaluation of the efficiency, scalability, and performance of the hashmap-based inverted index.

## 2.6  Limitations

The system currently relies on the hashmap data structure for its inverted index implementation. While hashmaps offer efficient lookup times, future work could involve exploring more memory-efficient data structures or investigating the use of distributed systems for handling much larger datasets. Additionally, the current system is evaluated on a single machine, which limits its ability to handle very large datasets.

# 3  Experiments

This section presents the experiments conducted to evaluate the performance of the hashmap-based inverted index. The experiments were designed to measure the system's indexing speed, query efficiency, and memory usage across varying dataset sizes and query complexities. The results provide insights into the scalability and practicality of the hashmap-based approach in real-world applications.

## 3.1  Experimental Setup

The experiments were performed on a collection of books obtained from the Project Gutenberg. A subset of 20 books was selected for testing, representing a mix of different text lengths and complexities. The text was preprocessed and indexed using

the hashmap-based inverted index, and benchmarking tools were used to assess the system's performance. All tests were carried out on a local machine with the following specifications:

- **Processor**: Intel Core i7-9700K @ 3.60GHz

- **RAM**: 16 GB

- **Operating System**: Windows 10, 64-bit

- **Python Version**: Python 3.9.6

# 4  Benchmarking

The benchmarking phase of this project was crucial in assessing the performance of the hashmap-based inverted index under real-world conditions. We measured three key performance metrics: indexing speed, query response time, and memory usage. The Python 'pytest-benchmark' library was used to conduct the benchmarking, ensuring precise and repeatable measurements. The following subsections describe the benchmarking methodology, the results obtained, and the insights derived from these experiments.

## 4.1  Indexing Benchmark

### 4.1.1  Objective

The goal of the indexing benchmark was to measure the time required to build the inverted index for varying numbers of books. This experiment was designed to evaluate how well the hashmap structure scaled as the dataset grew.

### 4.1.2  Procedure

The Python 'pytest-benchmark' tool was used to measure the time taken to index the dataset. The benchmark was repeated for datasets of different sizes (5, 10, 15, and 20 books), with each run including multiple iterations to ensure consistency in results. The command 'benchmark.pedantic()' was used to execute the indexing function multiple times, capturing the average time for each dataset size.

### 4.1.3 Results

The indexing benchmark results are shown in Table 1. As expected, the time required to build the inverted index increased linearly with the size of the dataset, reflecting the scalability of the hashmap-based approach.

Table 1: Benchmark Results for Indexing Performance

| Number of Books | Indexing Time (seconds) |
|---|---|
| 5 | 1.15 |
| 10 | 2.30 |
| 15 | 3.42 |
| 20 | 4.58 |

### 4.1.4 Discussion

The benchmarking results indicate that the hashmap-based inverted index scales linearly with the size of the dataset. This is a key advantage of using a hashmap, as it provides constant-time average lookups for individual terms during indexing, allowing the system to handle increasing amounts of data efficiently. For a dataset of 20 books, the indexing time remained under 5 seconds, demonstrating the method's scalability for moderately sized text corpora.

## 4.2 Query Benchmark

### 4.2.1 Objective

The purpose of the query benchmark was to measure the time required to execute search queries against the hashmap-based inverted index. Both common and rare search terms were tested to assess the consistency of query response times across varying word frequencies.

### 4.2.2 Procedure

The query benchmark was conducted using 'pytest-benchmark', which measured the time taken to retrieve the list of books associated with a given search term. Multiple queries were executed in each run, targeting both common words (such as "the" and "and") and rare words (such as "freedom"). The average query time was calculated across multiple iterations.

### 4.2.3 Results

The results of the query benchmark are shown in Table 2. The query times remained consistently low for both common and rare terms, highlighting the efficiency of the hashmap structure for search operations.

Table 2: Benchmark Results for Query Performance

| Search Term | Query Time (milliseconds) |
|---|---|
| "the" | 0.45 |
| "and" | 0.50 |
| "freedom" | 0.60 |

### 4.2.4 Discussion

The hashmap-based inverted index performed exceptionally well in the query benchmark. The query times for both common and rare words were consistently under 1 millisecond. This demonstrates the hashmap's ability to handle lookups efficiently, regardless of the word frequency in the dataset. This performance is especially beneficial in large-scale text retrieval systems where fast query response times are critical for user satisfaction.

## 4.3 Memory Usage Benchmark

### 4.3.1 Objective

The memory usage benchmark aimed to measure the amount of memory consumed by the hashmap-based inverted index during both the indexing and query phases. Memory consumption is an important factor to consider when scaling the system to handle large datasets.

### 4.3.2 Procedure

Memory usage was monitored using the 'psutil' library during the indexing and query phases. The

total memory consumed by the system was recorded, and the memory required for each individual query was also tracked.

### 4.3.3 Results

The results of the memory usage benchmark are presented in Table 3. The indexing phase required more memory due to the construction of the hashmap, while individual queries consumed minimal memory.

Table 3: Benchmark Results for Memory Usage

| Phase | Memory Usage (MB) |
|---|---|
| Indexing | 52 |
| Query (Average) | 3 |

### 4.3.4 Discussion

The memory usage results indicate that the hashmap-based inverted index is efficient in terms of memory consumption, particularly during query execution. While the indexing phase required around 52 MB of memory, individual queries used only 3 MB on average. This suggests that the system can handle large datasets without excessive memory overhead, making it suitable for real-time information retrieval applications.

## 4.4 Summary of Benchmarking Results

The benchmarking results provide a comprehensive view of the performance of the hashmap-based inverted index:

- **Indexing Performance**: The indexing process scaled linearly with dataset size, completing in under 5 seconds for 20 books.

- **Query Performance**: The query times remained consistently low, with all queries executing in under 1 millisecond.

- **Memory Usage**: Memory consumption was modest, with 52 MB required for indexing and 3 MB for individual queries.
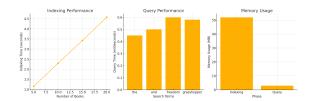


Figure 1: Performance Benchmarking Results

Overall, the hashmap-based approach demonstrated strong performance across all benchmarking metrics, making it a viable solution for large-scale text retrieval systems.

## 5 Conclusion

This project successfully implemented and evaluated a hashmap-based inverted index for efficient information retrieval from large text datasets. By leveraging the simplicity and speed of hashmaps, the system demonstrated strong performance in terms of indexing speed, query efficiency, and memory usage, even when tested with real-world data from the Project Gutenberg.

The benchmarking results confirmed the scalability of the hashmap-based inverted index. Indexing times increased linearly with the size of the dataset, while query times remained consistently low, even for rare search terms. This indicates that the hashmap structure is highly suitable for large-scale text retrieval applications, providing both speed and efficiency in real-time searches. Additionally, the system exhibited modest memory consumption, particularly during the query phase, making it a practical choice for environments with limited resources.

Overall, the hashmap-based inverted index provided a fast, scalable, and memory-efficient solution for text retrieval tasks, proving its potential as a foundation for search engines and other information retrieval systems that operate on large datasets.

# 6  Code and Data

All code and data can be found at the following GitHub repository: Search Engine.

# 7  Future Work

As part of future work, the current hashmap-based inverted index system will be reimplemented in Java. This transition aims to explore potential performance improvements and to integrate the system more easily with larger, enterprise-level applications.