

# Web Search Engines HW2

You

October 22, 2024

## Abstract

This report describes a search system that creates an inverted index from large document datasets, and returns ranked search results through user queries.

## 1 Introduction

This search system is designed in C/C++ to process queries efficiently through creating and querying an inverted index. The entire program consists of multiple sub components which parse the MS MARCO data set, build a compressed inverted index, and processes user queries with the BM25 ranking function. This program aims to balance memory and time efficiency by using multithreading, blocked varbyte compression techniques, and document-at-a-time (DAAT) query processing. This system is designed to handle several gigabytes of documents, and provides conjunctive and disjunctive search queries.

## 2 Program Overview

The search system is a multi-component architecture designed to efficiently index a collection of documents and provide fast retrieval capabilities based on user queries. The system is composed of three main components: the **Parser and Indexer**, the **Merger**, and the **Query Processor**. Together, these components build up the whole search system.

### 2.1 Component Descriptions

#### 2.1.1 Parser and Indexer

The Parser and Indexer is responsible for reading the input document collection, extracting relevant terms, and generating an inverted index. It leverages multi-threading to enhance performance, allowing simultaneous processing of multiple passages. Key functionalities include:

- **Multi-Threaded Processing:** Utilizing a thread pool to parse documents concurrently, significantly reducing indexing time.
- **Term-Document Pair Generation:** For each processed passage, it generates term-document pairs that are stored in temporary files for later merging.

#### 2.1.2 Merger

The Merger component consolidates the temporary files generated by the Parser and Indexer into a final inverted index and lexicon. This component ensures that the data is structured for efficient retrieval. Key features include:

- **File Merging:** Reads multiple temporary files and consolidates their content into a single index file using a priority queue for efficient merging.
- **Lexicon Construction:** Maintains a lexicon that maps terms to their metadata, including document frequency and offset information in the index.
- **Compression:** Implements variable-byte encoding to compress document IDs, optimizing storage space in the index file.

### 2.1.3 Query Processor

The Query Processor is the interface through which users interact with the search system. It processes user queries and retrieves relevant documents based on their input. Key functionalities include:

- **BM25 Ranking:** Implements the BM25 algorithm to score documents based on their relevance to the query terms, considering term frequency and document length.
- **Flexible Query Processing:** Supports both conjunctive (AND) and disjunctive (OR) queries, allowing users to specify their search criteria.
- **Document Retrieval:** Uses a page table to map document IDs to document names, enabling easy identification of retrieved documents.
- **User Interaction:** Provides a user-friendly interface for entering queries and receiving results, with a loop for continuous querying until the user opts to exit.

## 2.2 System Workflow

1. **Document Ingestion:** The system begins with the Parser and Indexer, which reads a collection of documents, extracting terms and generating term-document pairs. This process is optimized through multi-threading.
2. **Temporary File Generation:** As documents are processed, term-document pairs are written to temporary files for efficient storage.
3. **Index Merging:** The Merger component consolidates these temporary files into a final inverted index and lexicon, employing efficient merging techniques and data compression to optimize performance.
4. **Query Execution:** Users can input queries into the Query Processor, which parses the query, retrieves relevant documents from the inverted index, and ranks them using the BM25 algorithm.
5. **Results Presentation:** The top results are displayed to the user, providing document IDs, names, and relevance scores, facilitating quick access to information.

This search system integrates advanced indexing and retrieval techniques, providing a robust solution for handling large document collections. Its modular design allows for efficient processing, merging, and querying, making it suitable for various applications in information retrieval and search engine development.

## 3 Running the Program

To run the program, follow the steps outlined below. The program is implemented in C++ and is built using a Makefile.

### 3.1 Toolkit Preparation

To build and run the program, ensure that you have the following installed:

- **C++ Compiler:** This project uses g++ version **14.2.0**, which is part of the MinGW-W64 distribution. You can check your version by running:

---

```
1 g++ --version
```

---

- **Make Utility:** The project requires GNU Make version **3.81** or higher. You can verify your version with:

---

```
1 make --version
```

---

## 3.2 Dataset preparation

- Download the MS MARCO dataset from <https://msmarco.z22.web.core.windows.net/msmarcoranking/collection.tar.gz>
- Create a folder called *data* under the root folder and extract *collection.tsv* into it.

## 3.3 Building the Program

1. **Navigate to the Source Directory:** Open a terminal and change to the `src` folder where the source code is located:

---

```
1 cd path/to/your/src
```

---

2. **Run the Make Command:** Execute the following command to build the program:

---

```
1 make build_and_run_all
```

---

This command will compile the necessary components and generate the executables in the `build` directory.

The Makefile provided includes the following targets:

- `all`: Compiles the main components: `parser_and_indexer_mt`, `merger`, and `query_processor`.
- `run_parser_and_indexer_mt`: Executes the multi-threaded parser and indexer.
- `run_merger`: Executes the merger component.
- `run_query_processor`: Executes the query processor.

## 3.4 Running the Components

Once the build process is complete, you can run each component as follows:

- **Run the Multi-Threaded Parser and Indexer:**

---

```
1 ../build/parser_and_indexer_mt
```

---

- **Run the Merger:**

---

```
1 ../build/merger
```

---

- **Run the Query Processor:**

---

```
1 ../build/query_processor
```

---

# 4 Parser and Indexer

## 4.1 High-Level Design

The Parser and Indexer has several different sub components:

- **Thread Pool:** Manages concurrent execution of tasks.
- **Parser and Indexer:** Processes input files, generates term-document pairs, and computes term frequencies.

- **Utilities:** Provides helper functions for tokenization, parsing, logging, and directory management.
- **Compression Module:** Compresses data using **varbyte encoding** to store integers compactly.

The parser and indexer uses a producer-consumer architecture, where the parser enqueues tasks into the thread pool for parallel processing. Synchronization between threads is achieved using `mutex` locks and `condition_variables` to ensure thread-safe operations.

## 4.2 Key Data Structures

### 4.2.1 TermDocPair

The `TermDocPair` structure stores essential information about terms and their associated documents to be used by the Merger program.

```

1      struct TermDocPair {
2          std::string term;
3          int docID;
4          float termFScore;
5
6          TermDocPair(const std::string& t, int d, float termFScore)
7              : term(t), docID(d), termFScore(termFScore) {}
8      };

```

The **Term frequency score** `termFScore`, is a precomputed part of the BM25 function, displayed in the Query Processor section.

### 4.2.2 ThreadPool

The `ThreadPool` class manages a pool of threads, supporting task enqueueing and uses mutexes and condition variables to manage thread-safe access to shared resources.

```

1      class ThreadPool {
2      public:
3          ThreadPool(size_t threads = 8, size_t maxThreadsInQueue = 32);
4          ~ThreadPool();
5
6          void enqueue(std::function<void()> f);
7          void waitAll();
8
9      private:
10         void worker();
11         std::vector<std::thread> workers;
12         std::queue<std::function<void()>> tasks;
13         std::mutex queueMutex;
14         std::condition_variable taskAvailable;
15         bool stop;
16         std::atomic<size_t> tasksRemaining;
17     };

```

### 4.2.3 Page Table and Document Lengths

The `pageTable` stores mappings of document IDs to document names, while the `docLengths` stores the lengths of each document. These data structures are synchronized using `mutex` locks to prevent race conditions.

## 4.3 Encoding Techniques

### 4.3.1 Varbyte Encoding

Varbyte encoding is used to compress integers into bytes with optimal length. The most significant bit (MSB) acts as a continuation bit, indicating whether more bytes follow.

```
1      std::vector<unsigned char> varbyteEncode(int number) {
2          std::vector<unsigned char> bytes;
3          while (number >= 128) {
4              bytes.push_back((number % 128) | 128); // Set MSB
5              number /= 128;
6          }
7          bytes.push_back(number); // Final byte with MSB unset
8          return bytes;
9      }
```

## 4.4 Synchronization Mechanisms

The program allows for thread safe access in critical regions.

- **Mutex:** Protects critical sections in the code.
- **Condition Variables:** Used to notify threads when tasks are available.
- **Atomic Variables:** Ensures atomic operations on counters (e.g., file counter).

## 4.5 Critical Functions

The `processPassageMT` function processes individual document passages, calculates term frequencies, and generates term-document pairs. The `TermDocPairs` are stored in a blocked, compressed, binary format. Each temporary file is stored according to the size of `MAXRECORDS`.

```
1      void processPassageMT(int docID, const std::string &passage,
2                          std::vector<TermDocPair> &termDocPairs,
3                          std::mutex &termDocPairsMutex,
4                          std::atomic<int> &fileCounter,
5                          std::unordered_map<int, int> &docLengths,
6                          std::mutex &docLengthsMutex) {
7          auto terms = tokenize(passage);
8          {
9              std::lock_guard<std::mutex> docLengthsLock(docLengthsMutex);
10             docLengths[docID] = terms.size();
11         }
12
13         std::map<std::string, int> frequencyMap;
14         for (const auto &term : terms) {
15             frequencyMap[term]++;
16         }
17
18         std::unique_lock<std::mutex> termDocPairsLock(termDocPairsMutex);
19         for (const auto &termFreqPair : frequencyMap) {
20             termDocPairs.emplace_back(termFreqPair.first, docID,
21                                     calculateTermFreqScore(termFreqPair.second, k1, b, terms.size()));
22         }
23
24         if (termDocPairs.size() >= MAX_RECORDS) {
25             std::vector<TermDocPair> termDocPairsCopy(std::move(termDocPairs));
26             termDocPairs.clear();
27         }
```

```

27         int curFileCounter = fileCounter.fetch_add(1);
28         saveTermDocPairsToFile(termDocPairsCopy, curFileCounter);
29     }
30 }

```

## 5 Merger

### 5.1 High-Level Design

The merger program contains several key components to build an optimized inverted index:

- **Priority Queue for Merging:** Merges temporary files outputted from parser and indexer in sorted manner
- **Index Construction:** Stores compressed postings lists in binary files.
- **Lexicon Management:** Maintains metadata for terms and allows for fast retrieval.
- **Compression Module:** Uses same varbyte encoding/decoding to compress integers

### 5.2 Key Data Structures

#### 5.2.1 LexiconEntry

The LexiconEntry structure stores metadata about terms in the index:

- **offset:** Byte offset where the postings list for the term begins in the index file.
- **length:** byte length of posting list.
- **docFrequency:** Number of documents containing the term.
- **blockCount:** Number of blocks used to store the term's postings.

```

1     struct LexiconEntry {
2         int64_t offset;
3         int32_t length;
4         int32_t docFrequency;
5         int32_t blockCount;
6         std::vector<int32_t> blockMaxDocIDs;
7         std::vector<int64_t> blockOffsets;
8     };

```

#### 5.2.2 TuplePQ

The TuplePQ type is a priority queue that merges entries from temporary files. The queue orders entries based on terms and document IDs, and does incremental merging with the merging sub indexes algorithm.

```

1     typedef std::priority_queue<
2         std::tuple<std::string, int, int, float>,
3         std::vector<std::tuple<std::string, int, int, float>>,
4         std::function<bool(const std::tuple<std::string, int, int, float>&,
5                             const std::tuple<std::string, int, int, float>&)>
6     > TuplePQ;

```

### 5.3 Encoding Techniques

#### 5.3.1 Varbyte Encoding

Varbyte encoding and decoding is used in the merging process as well.

## 5.4 Critical Functions

### 5.4.1 mergeTempFiles

The function merges temporary files by inserting them into priority queues using the merging subindexes algorithm, outputting compressed postings to the index file.

```
1      void mergeTempFiles(int numFiles, std::unordered_map<std::string, LexiconEntry> &lexicon) {
2          std::vector<std::ifstream> tempFiles(numFiles);
3          for (int i = 0; i < numFiles; ++i) {
4              tempFiles[i].open("../data/intermediate/temp" + std::to_string(i) + ".bin", std::ios::
5          }
6
7          std::ofstream indexFile("../data/index.bin", std::ios::binary);
8          TuplePQ pq([](const auto &a, const auto &b) {
9              return std::get<0>(a) > std::get<0>(b) ||
10                 (std::get<0>(a) == std::get<0>(b) && std::get<1>(a) > std::get<1>(b));
11          });
12
13          for (int i = 0; i < numFiles; ++i) {
14              if (tempFiles[i].peek() != EOF) {
15                  readPairToPQ(tempFiles[i], i, pq);
16              }
17          }
18
19          std::string currentTerm;
20          std::vector<std::pair<int, float>> postingsList;
21          int64_t offset = 0;
22
23          while (!pq.empty()) {
24              auto [term, docID, fileIndex, termFreqScore] = pq.top();
25              pq.pop();
26
27              if (term != currentTerm && !currentTerm.empty()) {
28                  saveAndClearCurPostingsList(postingsList, offset, indexFile, lexicon, currentTerm
29              }
30              currentTerm = term;
31              postingsList.emplace_back(docID, termFreqScore);
32
33              if (tempFiles[fileIndex].peek() != EOF) {
34                  readPairToPQ(tempFiles[fileIndex], fileIndex, pq);
35              }
36          }
37
38          if (!postingsList.empty()) {
39              saveAndClearCurPostingsList(postingsList, offset, indexFile, lexicon, currentTerm, cu
40          }
41      }
```

## 6 Query Processor

### 6.1 High-Level Design

The query processor contains the following components:

- **Inverted Index:** This structure efficiently maps terms to their corresponding document identifiers and term frequencies score. It serves as the primary data repository for lookup during query processing.

- **Inverted Index Pointer:** This structure is an iterator for traversing the inverted index. It provides functionality to navigate through the list of document identifiers associated with a particular term, enabling efficient access to the posting lists.
- **Query Parser:** Responsible for processing the user's query and holding necessary data in memory (e.g. page table).

## 6.2 Key Data Structures

The following data structures are integral to the Query Processor's functionality:

### 6.2.1 InvertedListPointer

The `InvertedListPointer` handles the traversal of document IDs and retrieves term frequency scores. Here are some key attribute that `InvertedListPointer` has:

- `indexFile`: A pointer to the index file from which document IDs and term frequencies are read.
- `lexEntry`: An entry from the lexicon that provides metadata about the term being processed.
- `compressedData`: A buffer for storing compressed document ID associated with the term.
- `termFreqScore`: A vector that holds the term frequency scores for each document with that term.

```

1  class InvertedListPointer {
2  public:
3      InvertedListPointer(std::ifstream *indexFile, const LexiconEntry &lexEntry);
4      bool next();
5      bool nextGEQ(int docID);
6      int getDocID() const;
7      float getTFS() const;
8      int getTF() const;
9      bool isValid() const;
10     void close();
11     float getIDF() const;
12 private:
13     std::ifstream *indexFile;
14     LexiconEntry lexEntry;
15     int currentDocID;
16     bool valid;
17     int lastDocID;
18     size_t bufferPos;
19     size_t termFreqScoreIndex;
20     std::vector<unsigned char> compressedData;
21     std::vector<float> termFreqScore;
22 };

```

### 6.2.2 InvertedIndex

The `InvertedIndex` class is responsible for retrieving an inverted index given an index file, lexicon information, and a term. Here are some key attributed:

- `indexFile`: An input file stream for reading the inverted index data.
- `lexicon`: A hash map that stores the lexicon entries, mapping each term to its corresponding metadata, including offsets and other information.



```

1     class InvertedIndex {
2     public:
3         InvertedIndex(const std::string &indexFilename, const std::string &lexiconFilename);
4         bool openList(const std::string &term);
5         InvertedListPointer getListPointer(const std::string &term);
6         void closeList(const std::string &term);
7         int getDocFrequency(const std::string &term); // New method
8     private:
9         std::ifstream indexFile;
10        std::unordered_map<std::string, LexiconEntry> lexicon;
11
12        void loadLexicon(const std::string &lexiconFilename);
13    };

```

### 6.2.3 QueryProcessor

The QueryProcessor class is designed to handle the processing of user queries.

- **invertedIndex:** An instance of the InvertedIndex class that provides access to the document lists and term frequencies needed for query processing.
- **pageTable:** A hash map that maps document IDs to their corresponding document names.

```

1     class QueryProcessor {
2     public:
3         QueryProcessor(const std::string &indexFilename, const std::string &lexiconFilename, const std::string &pageTableFilename, const std::string &docLengthsFilename);
4         void processQuery(const std::string &query, bool conjunctive);
5     private:
6         InvertedIndex invertedIndex;
7         std::unordered_map<int, std::string> pageTable; // docID -> docName
8         std::unordered_map<int, int> docLengths; // docID -> docLength
9         int totalDocs;
10        double avgDocLength;
11
12        std::vector<std::string> parseQuery(const std::string &query);
13        void loadPageTable(const std::string &pageTableFilename);
14        void loadDocumentLengths(const std::string &docLengthsFilename);
15    };

```

## 6.3 Critical Functions

### 6.3.1 QueryProcessor::processQuery

The QueryProcessor::processQuery function is designed to manage search queries against an inverted index. It parses the query into terms, retrieves and validates the corresponding inverted lists, scores the documents based on their relevance using the BM25 algorithm, and ultimately ranks and displays the top results based on the aggregated scores.

```

1     void QueryProcessor::processQuery(const std::string &query, bool conjunctive) {
2         auto terms = parseQuery(query);
3
4         if (terms.empty()) {
5             std::cout << "No terms found in query." << std::endl;
6             return;
7         }
8
9         // Open inverted lists for each term and store term-pointer pairs
10        std::vector<std::pair<std::string, InvertedListPointer>> termPointers;
11        for (const auto &term : terms) {

```

```

12         if (!invertedIndex.openList(term)) {
13             std::cout << "Term not found: " << term << std::endl;
14             continue;
15         }
16         termPointers.emplace_back(term, invertedIndex.getListPointer(term));
17     }
18
19     if (termPointers.empty()) {
20         std::cout << "No valid terms found in query." << std::endl;
21         return;
22     }
23
24     // DAAT Processing
25     std::unordered_map<int, double> docScores; // docID -> aggregated score
26
27     if (conjunctive) {
28         // Initialize docIDs for each list
29         std::vector<int> docIDs;
30         for (auto &tp : termPointers) {
31             auto &ptr = tp.second;
32             if (!ptr.isValid()) { // Check if list is empty
33                 ptr.close();
34                 return; // One of the lists is empty, no results
35             }
36             if (!ptr.next()) { // Move to first element
37                 ptr.close();
38                 return; // Empty list
39             }
40             docIDs.push_back(ptr.getDocID());
41         }
42
43         while (true) {
44             int maxDocID = *std::max_element(docIDs.begin(), docIDs.end());
45             bool allMatch = true;
46
47             // Advance pointers where docID < maxDocID
48             for (size_t i = 0; i < termPointers.size(); ++i) {
49                 auto &ptr = termPointers[i].second;
50                 while (docIDs[i] < maxDocID) {
51                     if (!ptr.nextGEQ(maxDocID)) {
52                         allMatch = false;
53                         break; // Reached end of list
54                     }
55                     docIDs[i] = ptr.getDocID();
56                 }
57                 if (docIDs[i] != maxDocID) {
58                     allMatch = false;
59                 }
60             }
61
62             if (!allMatch) {
63                 // Check if any list has reached the end
64                 bool anyEnd = false;
65                 for (auto &tp : termPointers) {
66                     if (!tp.second.isValid()) {
67                         anyEnd = true;

```

```

68         break;
69     }
70 }
71 if (anyEnd) break;
72 continue;
73 }
74 // allMatch equal to true
75 // All pointers are at the same docID
76 int docID = maxDocID;
77 double totalScore = 0.0;
78 for (auto &tp : termPointers) {
79     auto &ptr = tp.second;
80     const std::string &term = tp.first;
81
82     // Compute BM25 score
83     float bm25Score = ptr.getIDF() * ptr.getTFS();
84     totalScore += bm25Score;
85
86     ptr.next(); // Advance pointer for next iteration
87 }
88 docScores[docID] = totalScore;
89
90 // Update docIDs
91 bool validPointers = true;
92 for (size_t i = 0; i < termPointers.size(); ++i) {
93     auto &ptr = termPointers[i].second;
94     if (ptr.isValid()) {
95         docIDs[i] = ptr.getDocID();
96     } else {
97         // One of the lists has reached the end
98         validPointers = false;
99         break;
100     }
101 }
102 if (!validPointers) break;
103 }
104 } else {
105     // Disjunctive query processing using a min-heap
106     auto cmp = [](std::pair<InvertedListPointer*, std::string> a, std::pair<InvertedListPointer*, std::string> b) {
107         return a.first->getDocID() > b.first->getDocID();
108     };
109     std::priority_queue<
110         std::pair<InvertedListPointer*, std::string>,
111         std::vector<std::pair<InvertedListPointer*, std::string>>,
112         decltype(cmp)> pq(cmp);
113
114     // Initialize heap with the first posting from each term
115     for (auto &tp : termPointers) {
116         auto &ptr = tp.second;
117         const std::string &term = tp.first;
118         if (ptr.isValid() && ptr.next()) {
119             pq.push({&ptr, term});
120         }
121     }
122
123     while (!pq.empty()) {

```

```

124         auto [ptr, term] = pq.top();
125         pq.pop();
126
127         int docID = ptr->getDocID();
128
129         // Compute BM25 score
130         float bm25Score = ptr->getIDF() * ptr->getTFS();
131
132         docScores[docID] += bm25Score;
133
134         // Advance the pointer and re-add to the heap if valid
135         if (ptr->next()) {
136             pq.push({ptr, term});
137         }
138     }
139 }
140
141 // Close all inverted lists
142 for (auto &tp : termPointers) {
143     tp.second.close();
144 }
145
146 if (docScores.empty()) {
147     std::cout << "No documents matched the query." << std::endl;
148     return;
149 }
150
151 // Rank documents by aggregated score
152 std::vector<std::pair<int, double>> rankedDocs(docScores.begin(), docScores.end());
153
154 // Sort by score in descending order
155 std::sort(rankedDocs.begin(), rankedDocs.end(), [](const auto &a, const auto &b) {
156     return a.second > b.second;
157 });
158
159 // Display top 10 results
160 int resultsCount = std::min(10, static_cast<int>(rankedDocs.size()));
161 for (int i = 0; i < resultsCount; ++i) {
162     int docID = rankedDocs[i].first;
163     double score = rankedDocs[i].second;
164     std::string docName = pageTable[docID];
165     std::cout << i + 1 << ". DocID: " << docID << ", DocName: " << docName << ", Score: " <<
166     }
167 }

```

## 7 Web Interface

## 8 Performance Analysis

Metric	Details
collection.tsv	2.85 GB
Intermediate Binary Temporary Files	348 files, each 14.82 MB
Parsing + Indexing Time	128 seconds
Merging Time	1387 seconds

## 8.1 Query Processing

- Query 1: ‘‘hello world + AND’’ *Response: 5 ms*
- Query 2: ‘‘hello world + OR’’ *Response: 222 ms*
- Query 3: ‘‘how to cook good food + AND’’ *Response: 166 ms*
- Query 4: ‘‘how to cook good food + OR’’ *Response: 8156 ms*

## 9 Limitations

Currently our merger is the biggest bottleneck. This could probably be sped up by adding multi-threading to our process.

A web interface is currently in the works as well.

Created two test files **test\_bin\_reader.cpp** and **test\_merger.cpp** to test temp file and binary file format