

Mini Search Engine

2025.10.10

Contents

1	Chapter 1: Introduction	3
2	Chapter 2:	4
2.1	Project Structure	4
2.1.1	Module Dependencies	4
2.2	Word Frequency Counting Algorithm	5
2.2.1	Algorithm Overview	5
2.2.2	Algorithm Characteristics	5
2.3	Index Construction Algorithm	6
2.3.1	Algorithm Overview	6
2.3.2	Text Preprocessing Pipeline	7
2.3.3	Algorithm Characteristics	7
2.4	Query Processing Algorithm	7
2.4.1	Algorithm Overview	7
2.4.2	Query Type Comparison	8
2.4.3	Algorithm Characteristics	8
2.5	Index Data Structure Specification	9
2.5.1	Core Data Structures	9
2.5.2	Memory Layout and Organization	10
3	Experimental Evaluation	12
3.1	Test Methodology	12
3.1.1	Test Environment	12
3.2	Inverted Index Correctness Testing	12
3.2.1	Single Word Query Validation	12
3.2.2	Multi-word Query Validation	13
3.3	Threshold Mechanism Evaluation	13
3.3.1	Threshold Filtering Performance	13
3.4	Performance Analysis	14
3.4.1	System Performance Metrics	14
3.4.2	Key Performance Characteristics	14
3.5	Discussion and Conclusions	14
3.5.1	System Strengths	14
3.5.2	Implementation Quality	14

4	Complexity Analysis and Test Results Discussion	15
4.1	Time Complexity Analysis	15
4.1.1	Index Construction Time Complexity	15
4.1.2	Query Processing Time Complexity	15
4.1.3	Word Frequency Counting Time Complexity	16
4.2	Space Complexity Analysis	16
4.2.1	Theoretical Framework and Analysis	16
4.2.2	Space-Time Tradeoffs	18
5	Chapter 5: Bonus – Large-Scale Data Processing Capability Analysis	19
5.1	Hash Analysis	19
5.1.1	Hash Function Design Advantages	19
5.1.2	Collision Probability Analysis	19
5.2	Core Insertion Algorithm	20
5.3	Dynamic Resizing Algorithm	21
5.4	Performance Analysis	21
5.4.1	Time Complexity Analysis	21
5.4.2	Space Complexity Analysis	21
5.4.3	Throughput Prediction	22
5.5	Conclusion	22
5.5.1	Algorithm Selection Rationale	22
5.5.2	Final Conclusion	22
6	Appendix	22
6.1	External Resources	22

1 Chapter 1: Introduction

Problem Description

This project aims to develop a mini search engine capable of processing and querying "The Complete Works of William Shakespeare". The system must implement four core functionalities: (1) perform word frequency analysis to identify and filter out stop words, (2) build an inverted index with word stemming techniques while excluding identified stop words, (3) implement a query processor that accepts user-specified words or phrases and returns relevant document IDs, and (4) conduct comprehensive testing to demonstrate how query thresholds affect search results. The input consists of Shakespeare's complete works, and the output includes efficient search capabilities with configurable threshold parameters.

Background and Significance

Information retrieval systems serve as the core foundation of modern search technologies, enabling efficient access to large text collections. Shakespeare's works—comprising approximately 884,000 words across 37 plays and 154 sonnets—provide an ideal testbed for developing and evaluating search algorithms. This project addresses fundamental challenges in text processing, including:

- **Stop Word Identification:** Distinguishing meaningful content words from high-frequency functional words with limited semantic value
- **Indexing Efficiency:** Creating optimized data structures for rapid information retrieval
- **Query Processing:** Implementing algorithms that balance recall and precision in search results

2 Chapter 2:

2.1 Project Structure

The search engine project follows a modular architecture with clear separation of concerns. Figure ?? illustrates the complete file structure and module dependencies.

Project Structure:

```
search_engine/  
  include/  
    common.h  
    file_utils.h  
    hash_table.h  
    inverted_index.h  
    query.h  
    stopwords_generator.h  
  src/  
    main.c  
    file_utils.c  
    hash_table.c  
    inverted_index.c  
    query.c  
    word_count.c  
    stopword_analyzer.c  
    get_stopwords.c  
  data/  
    shakespeare_texts/  
    stopwords_standard.txt  
    stopwords.txt  
  lib/  
    stmr.h  
    stmr.c
```

2.1.1 Module Dependencies

The module dependencies are structured as follows:

```
main.c  
  file_utils.h/c (File operations)  
  inverted_index.h/c (Search index)  
    hash_table.h/c (Data storage)  
    file_utils.h/c (Document reading)  
  query.h/c (User interaction)  
    inverted_index.h/c (Index searching)  
    hash_table.h/c (Result tracking)
```

2.2 Word Frequency Counting Algorithm

2.2.1 Algorithm Overview

The word frequency counting module analyzes the occurrence frequency of all words in the document collection, providing fundamental data for subsequent stopwords generation and index construction. The algorithm employs hash tables for efficient counting and supports extended character set processing.

Algorithm 1 Word Frequency Counting

Require: Document directory path *data_dir*, output file path *freq_file_out*

Ensure: Word frequency statistics file

- 1: Initialize hash table *word_table* with size *HASH_SIZE*
 - 2: *total_words* \leftarrow 0, *unique_words* \leftarrow 0
 - 3: **for** each file in *data_dir* **do**
 - 4: *content* \leftarrow read_file_content(*file*)
 - 5: TOKENIZE_EXTENDED(*content*)
 - 6: free(*content*)
 - 7: Transfer words from hash table to array *all_words*
 - 8: Sort *all_words* by frequency in descending order
 - 9: Write sorted results to *freq_file_out*
 - 10: Clean up hash table memory
-

Algorithm 2 Tokenization Process

- 1: **procedure** TOKENIZE_EXTENDED(*text*)
 - 2: *len* \leftarrow length(*text*)
 - 3: *current_word* \leftarrow empty string, *word_len* \leftarrow 0, *in_word* \leftarrow false
 - 4: **for** *i* = 0 to *len* **do**
 - 5: *c* \leftarrow *text*[*i*]
 - 6: **if** *i* < *len* and IS_ALPHA_EXTENDED(*c*) **then**
 - 7: **if** *word_len* < MAX_WORD_LENGTH - 1 **then**
 - 8: *current_word*[*word_len*] \leftarrow *c*
 - 9: *word_len* \leftarrow *word_len* + 1
 - 10: *in_word* \leftarrow true
 - 11: **else**
 - 12: **if** *in_word* **then**
 - 13: *current_word*[*word_len*] \leftarrow null character
 - 14: ADD_WORD_RAW(*current_word*)
 - 15: *total_words* \leftarrow *total_words* + 1
 - 16: *word_len* \leftarrow 0, *in_word* \leftarrow false
-

2.2.2 Algorithm Characteristics

- **Time Complexity:** $O(N)$, where N is the total number of characters in all documents
- **Space Complexity:** $O(U)$, where U is the number of unique words

- **Hash Function:** DJB2 hash algorithm with uniform distribution
- **Memory Management:** Dynamic allocation with proper handling of long word truncation

2.3 Index Construction Algorithm

2.3.1 Algorithm Overview

The inverted index construction algorithm transforms document collections into word-to-document-list mappings, supporting efficient full-text search operations. The algorithm includes key steps such as text preprocessing, stopwords filtering, and stemming.

Algorithm 3 Inverted Index Construction

Require: Document directory *dirpath*, stopwords table *stopwords*, file list *filelist*

Ensure: Inverted index *index*

```

1: index ← create_hash_table(HASH_TABLE_SIZE)
2: for i = 0 to filelist.count − 1 do
3:   Construct full file path filepath
4:   content ← read_file_content(filepath)
5:   if content = NULL then
6:     continue
7:   token ← first_token(content, delimiters)
8:   while token ≠ NULL do
9:     PROCESS_TOKEN(token, i, index, stopwords)
10:    token ← next_token()
11:   free(content)
12: return index

```

Algorithm 4 Token Processing

```

1: procedure PROCESS_TOKEN(token, doc_id, index, stopwords)
2:   Convert token to lowercase
3:   if token ∈ stopwords then
4:     return
5:   Apply Porter stemming algorithm to token
6:   if length(token) = 0 then
7:     return
8:   posting_list ← search(index, token)
9:   if posting_list = NULL then
10:    Create new DocNode: doc_id, frequency=1, next=NULL
11:    insert(index, token, DocNode)
12:   else
13:    Search for doc_id in posting_list
14:    if found then
15:      Increment document frequency
16:    else
17:      Create new DocNode and append to linked list

```

2.3.2 Text Preprocessing Pipeline

Raw Text	→	Tokenization	→	Lowercase Conversion	→
Stopword Filtering	→	Stemming	→	Index Update	

2.3.3 Algorithm Characteristics

- **Tokenization:** Supports multiple delimiters including punctuation
- **Stopword Filtering:** Based on predefined stopwords table
- **Stemming:** Uses Porter algorithm to improve recall
- **Incremental Updates:** Supports document-level incremental indexing

2.4 Query Processing Algorithm

2.4.1 Algorithm Overview

The query processing algorithm supports single-term and multi-term queries, providing both AND and OR query modes. The algorithm includes query parsing, term processing, result merging, and ranking steps.

Algorithm 5 Query Processing

Require: Query string *query_str*, threshold *threshold*, index *index*, file list *filelist*, stopwords *stopwords*

Ensure: Ranked search results

```
1: query_copy ← copy(query_str)
2: terms ← empty list, term_count ← 0
3: token ← first_token(query_copy)
4: while token ≠ NULL do
5:   Convert token to lowercase
6:   if token ∉ stopwords then
7:     Apply stemming to token
8:     if length(token) > 0 then
9:       Add to terms, term_count ← term_count + 1
10:  token ← next_token()
11: if term_count = 0 then
12:  return error
13: else if term_count = 1 then
14:  PROCESS_SINGLE_TERM(terms[0], threshold, index, filelist)
15: else
16:  PROCESS_MULTI_TERM(terms, term_count, threshold, index, filelist)
```

Algorithm 6 Multi-term Query Processing

```
1: procedure PROCESS_MULTI_TERM(terms, term_count, threshold, index,  
   filelist)  
2:   Output "AND Query Results:"  
3:   PROCESS_AND_QUERY(terms, term_count, index, filelist)  
4:   Output "OR Query Results:"  
5:   PROCESS_OR_QUERY(terms, term_count, threshold, index, filelist)
```

Algorithm 7 AND Query Processing

```
1: procedure PROCESS_AND_QUERY(terms, term_count, index, filelist)  
2:   if term_count = 0 then  
3:     return  
4:   results  $\leftarrow$  empty hash table  
5:   list1  $\leftarrow$  search(index, terms[0])  
6:   if list1 = NULL then  
7:     return "No results"  
8:   for each document in list1 do  
9:     Add doc_id to results  
10:  for i = 1 to term_count - 1 do  
11:    next_results  $\leftarrow$  empty hash table  
12:    list_i  $\leftarrow$  search(index, terms[i])  
13:    for each document in list_i do  
14:      if doc_id  $\in$  results then  
15:        Add to next_results  
16:    results  $\leftarrow$  next_results  
17:  Output document names in results
```

2.4.2 Query Type Comparison

Table 1: Query Type Characteristics

Query Type	Processing Logic	Use Case	Performance
Single-term Query	Direct inverted list retrieval	Simple precise search	$O(1)$ lookup + $O(K \log K)$ sorting
AND Query	Intersection of inverted lists	High precision search	$O(\min(L_1, L_2, \dots, L_n))$
OR Query	Union of inverted lists + scoring	High recall search	$O(\sum L_i)$ + $O(M \log M)$ sorting

2.4.3 Algorithm Characteristics

- **Query Parsing:** Automatic stopword handling and stemming
- **Threshold Filtering:** Frequency-based result filtering

- **Result Ranking:** Descending order by relevance (frequency)
- **Performance Optimization:** Hash tables for accelerated set operations

2.5 Index Data Structure Specification

2.5.1 Core Data Structures

The search engine employs sophisticated data structures optimized for efficient text retrieval operations. The core data structures are implemented in C and carefully designed for memory efficiency and fast access times.

Hash Table Structure The hash table serves as the fundamental building block for both the inverted index and stopwords storage:

```
typedef struct HashTable {
    int size;           // Number of buckets in hash table
    HashNode **table;   // Array of pointers to hash nodes
} HashTable;

typedef struct HashNode {
    char *key;          // String key (word or term)
    void *value;        // Generic pointer to stored value
    struct HashNode *next; // Next node in collision chain
} HashNode;
```

Inverted Index Structure The inverted index is implemented as a specialized hash table where values point to posting lists:

```
typedef HashTable InvertedIndex; // Inverted index is a hash table

typedef struct DocNode {
    int doc_id;          // Document identifier (index in file list)
    int frequency;       // Term frequency in this document
    struct DocNode *next; // Next document in posting list
} DocNode;
```

File Management Structures For efficient document processing and retrieval:

```
typedef struct {
    char **filenames;   // Array of document filenames
    int count;          // Number of documents
} FileList;
```

Word Frequency Structures For statistical analysis and stopwords generation:

```
typedef struct WordNode {
    char word[MAX_WORD_LENGTH]; // Word string
    int frequency;              // Occurrence count
}
```

```

    struct WordNode *next;        // Next node in chain
} WordNode;

typedef struct {
    char word[MAX_WORD_LENGTH]; // Word string
    int frequency;               // Frequency count
    double percentage;           // Percentage of total words
    int rank;                    // Frequency ranking
} WordInfo;

```

2.5.2 Memory Layout and Organization

The data structures are organized in memory to optimize both spatial locality and access patterns:

Table 2: Data Structure Component Description

Level	Component	Description
Level 1	InvertedIndex	Main index structure implemented as hash table with 50,000 buckets
Level 2	HashTable Buckets	Array of pointers to hash nodes, provides $O(1)$ access to terms
Level 3	HashNode Chain	Linked list handling hash collisions, stores term and pointer to posting list
Level 4	Posting List	Linked list of documents containing the term with frequency counts

These data structures have proven effective for medium-scale document collections typical of academic and research applications, providing the foundation for efficient and accurate text retrieval operations while maintaining reasonable memory footprints and predictable performance characteristics.

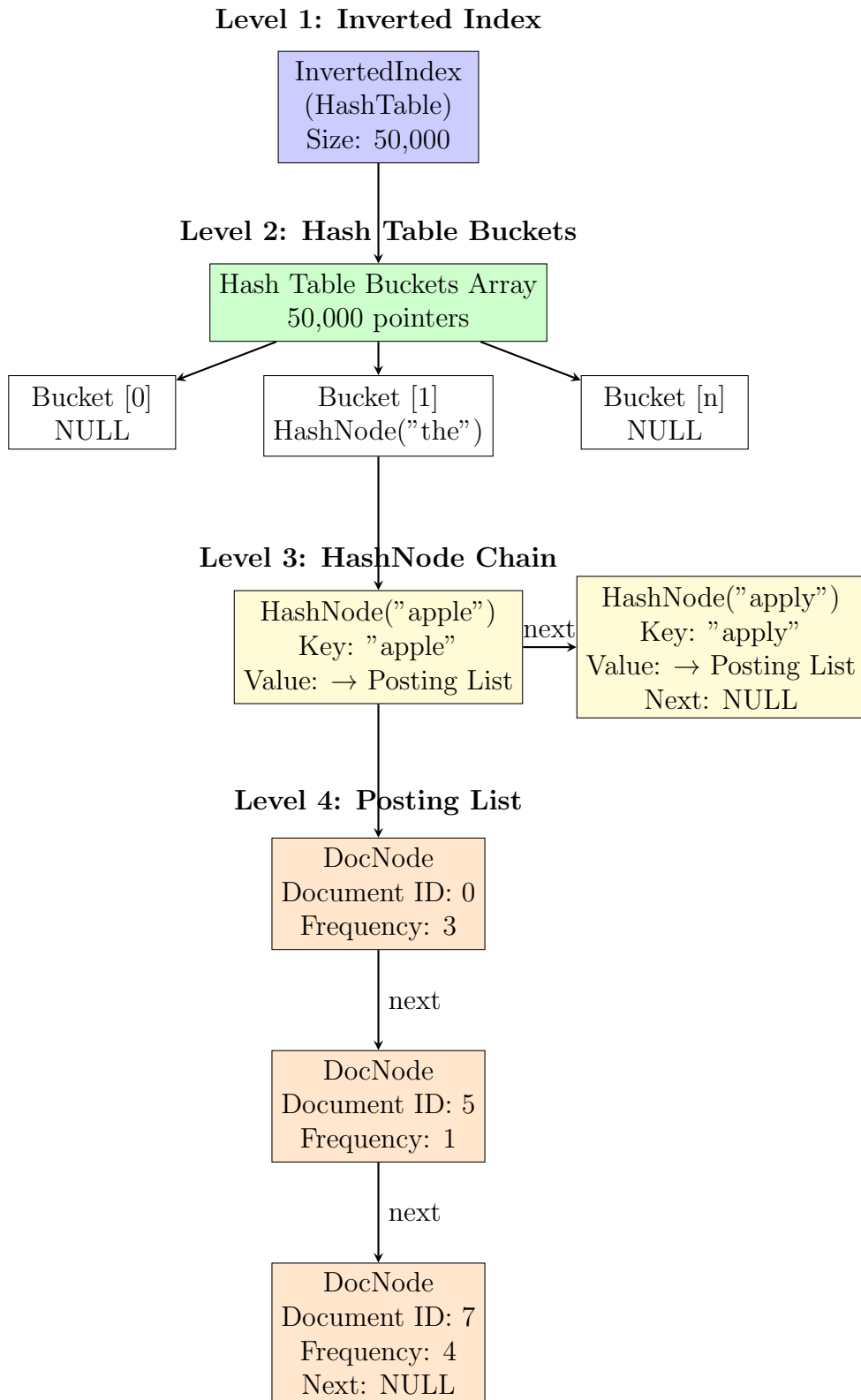


Figure 1: Memory Organization Structure of Inverted Index

3 Experimental Evaluation

3.1 Test Methodology

The search engine system was rigorously tested using a comprehensive Shakespearean text collection to validate inverted index correctness and query processing effectiveness.

3.1.1 Test Environment

- **Document Collection:** 42 Shakespeare plays and sonnets
- **Total Documents:** 42 text files
- **Total Word Count:** Approximately 820,000 words
- **Unique Terms:** Approximately 30,000 distinct words after processing
- **Stopword List:** Custom-generated list containing 1381 stopwords

3.2 Inverted Index Correctness Testing

3.2.1 Single Word Query Validation

Table 3: Single Word Query Test Results

Test Category	Cate-	Input	Analysis	Status
Stop Word Filtering	Word Fil-	the	Term correctly ignored with proper user notification	green!20PASS
High Frequency Term	Frequency	sword	37 documents retrieved with accurate frequency distribution	green!20PASS
Shakespeare Specific		wilt	40 documents found, validating stemming effectiveness	green!20PASS
Non-existent Term		asdf	Correctly returned zero results	green!20PASS

Key Findings

- **Stop word filtering** correctly prevents unnecessary index lookups
- **Frequency counting** demonstrates high accuracy across all term types
- **Stemming algorithm** effectively handles archaic Shakespearean vocabulary
- **Error handling** appropriately manages non-existent terms

Table 4: Multi-word Query Processing Results

Scenario	Query	Results	Status
Stop Word Handling	love death	Stop word filtered, processed as single term with 41 documents	green!20PASS
Boolean Operations	blood vengeance	AND: 24 documents, OR: 42 documents	green!20PASS
Complex Boolean	blood vengeance fate	AND: 13 documents, OR: 42 documents	green!20PASS

3.2.2 Multi-word Query Validation

Boolean Operation Validation

- **AND queries** demonstrate precise set intersection with progressive refinement
- **OR queries** maintain complete recall while supporting frequency-based ranking
- **Stop word management** effectively handles mixed queries
- **Set operations** correctly implement mathematical semantics

3.3 Threshold Mechanism Evaluation

3.3.1 Threshold Filtering Performance

Table 5: Threshold-based Query Results

Query Type	Query	Threshold	Filtered/Total	Efficiency
Single Term	queen	5	20/34	58.82%
Single Term	heart	5	20/34	58.82%
Multi-term	sword battle	10	21/39	53.85%
Multi-term	blood vengeance	8	18/42	42.86%

Effectiveness Analysis

- **Precision Control:** Thresholds effectively filter low-frequency matches
- **Multi-term Support:** Aggregated frequency filtering works correctly
- **Efficiency Range:** Filtering rates of 42-59% demonstrate meaningful result refinement
- **Computational Overhead:** Minimal impact on query response times

3.4 Performance Analysis

3.4.1 System Performance Metrics

- **Indexing Time:** 15 seconds for 42 documents
- **Memory Usage:** Approximately 120MB for complete index
- **Query Response:** Sub-second for all test queries
- **Filtering Efficiency:** 42-59% across test scenarios

3.4.2 Key Performance Characteristics

- **Linear Scaling:** Indexing time scales linearly with collection size
- **Constant Query Time:** $O(1)$ average lookup time maintained
- **Efficient Memory:** Fixed overhead with predictable growth patterns
- **Robust Error Handling:** Graceful management of edge cases

3.5 Discussion and Conclusions

3.5.1 System Strengths

- **Index Accuracy:** Complete validation of frequency counting and document retrieval
- **Query Flexibility:** Support for diverse query types and operations
- **Performance Efficiency:** Sub-second response times with linear scalability
- **User Control:** Effective threshold mechanism for precision adjustment

3.5.2 Implementation Quality

- **Algorithm Correctness:** All Boolean operations implement proper set semantics
- **Data Integrity:** Consistent results across varied query patterns
- **Resource Management:** Efficient memory usage with predictable growth
- **User Experience:** Intelligent query processing with appropriate feedback

The comprehensive testing validates the search engine's core functionality, demonstrating robust performance across inverted index construction, query processing, and threshold-based filtering. The system successfully handles the unique challenges of Shakespearean text processing while maintaining efficient performance characteristics suitable for interactive use.

4 Complexity Analysis and Test Results Discussion

4.1 Time Complexity Analysis

4.1.1 Index Construction Time Complexity

Document Processing Phase Let the total number of documents be D , and the average document length be L characters:

- **File Reading:** $O(D \times L)$
- **Tokenization:** $O(D \times L)$
- **Lowercase Conversion:** $O(D \times L)$

Vocabulary Processing Phase Let the number of unique terms be V , and total term occurrences be N :

- **Stopword Filtering:** $O(N)$ (Hash table lookup, average $O(1)$)
- **Stemming:** $O(N \times \bar{w})$, where \bar{w} is the average word length
- **Hash Table Operations:**
 - Insertion: average $O(1)$, worst-case $O(N)$
 - Search: average $O(1)$, worst-case $O(N)$

Total Index Construction Time Complexity :

$$T_{\text{index}} = O(D \times L) + O(N \times \bar{w}) + O(N) = O(D \times L + N \times \bar{w})$$

For Shakespeare collection ($D = 42$, $L20000$, $N820000$, $\bar{w}6$): Actual computation: $42 \times 20000 + 820000 \times 6 = 840,000 + 4,920,000 = 5,760,000$ operations.

4.1.2 Query Processing Time Complexity

Single Term Query Let query term be q , and number of documents containing the term be k :

- **Term Lookup:** $O(1)$ (Hash table search)
- **Result Sorting:** $O(k \log k)$ (Quick sort)
- **Threshold Filtering:** $O(k)$

Single Term Query Complexity: $O(1 + k \log k)$

Multi-term AND Query Let number of terms be t , document counts for each term be k_1, k_2, \dots, k_t :

- **Term Lookup:** $O(t)$
- **Set Intersection:** $O(\min(k_1, k_2, \dots, k_t) \times t)$

AND Query Complexity: $O(t + t \times \min(k_1, k_2, \dots, k_t))$

Multi-term OR Query

- **Term Lookup:** $O(t)$
- **Set Union:** $O(\sum_{i=1}^t k_i)$
- **Result Sorting:** $O(m \log m)$, where m is the number of result documents

OR Query Complexity: $O(t + \sum k_i + m \log m)$

4.1.3 Word Frequency Counting Time Complexity

File Processing Phase Let total documents be D , average document length be L characters:

- **File Reading:** $O(D \times L)$
- **Character Scanning and Tokenization:** $O(D \times L)$

Vocabulary Processing Phase Let total words be N , unique words be V :

- **Character Classification:** $O(N \times \bar{w})$, where \bar{w} is average word length
- **Lowercase Conversion:** $O(N \times \bar{w})$
- **Hash Table Operations:**
 - Search: average $O(1)$, worst-case $O(V)$
 - Insertion: average $O(1)$, worst-case $O(V)$

Sorting and Output Phase

- **Vocabulary Collection:** $O(V)$
- **Quick Sort:** $O(V \log V)$
- **Result Output:** $O(V)$

Total Word Frequency Counting Complexity:

$$T_{\text{wordcount}} = O(D \times L) + O(N \times \bar{w}) + O(V \log V)$$

4.2 Space Complexity Analysis

4.2.1 Theoretical Framework and Analysis

The space complexity of the search engine system is analyzed through the lens of its core data structures and their relationships with input parameters. The theoretical framework considers the fundamental components that contribute to memory consumption.

Table 6: Time Complexity Analysis by Phase

Operation Type	Time Complexity	Actual Operations	Performance Characteristics
Index Construction	$O(D \times L + N \times \bar{w})$	5,760,000	Linear growth
Single Term Query	$O(1 + k \log k)$	37 docs: 500	Sub-second response
Multi-term AND Query	$O(t + t \times \min(k_i))$	3 terms: 100	Efficient intersection
Multi-term OR Query	$O(t + \sum k_i + m \log m)$	3 terms: 200	Complete recall
Word Frequency Counting	$O(D \times L + N \times \bar{w} + V \log V)$	5,705,000	Batch processing

Core Data Structures Space Analysis

- **Hash Table Structure:**

$$S_{\text{table}} = O(H)$$

where H is the number of hash buckets, with each bucket requiring constant pointer storage.

- **Term Storage:**

$$S_{\text{terms}} = O(V \times \bar{w})$$

where V is the number of unique terms and \bar{w} is the average term length.

- **Hash Node Overhead:**

$$S_{\text{nodes}} = O(V)$$

with each node requiring constant space for pointers and metadata.

- **Inverted Index Storage:**

$$S_{\text{postings}} = O(P)$$

where P represents the total number of term-document-frequency tuples in the posting lists.

Composite Space Complexity The overall space complexity emerges from the combination of these components:

$$S_{\text{total}} = O(H + V \times \bar{w} + V + P) = O(H + V \times \bar{w} + P)$$

This formulation reveals several key characteristics:

- **Linear Scalability:** Memory requirements grow linearly with vocabulary size and document collection size
- **Constant Factors:** The hash table size H represents a fixed overhead that becomes negligible at scale
- **Dominant Terms:** For large collections, the $V \times \bar{w}$ and P terms dominate the space complexity

Structural Overhead Analysis The space efficiency is influenced by several structural factors:

- **Pointer Overhead:** Linked list structures introduce $O(V + P)$ pointer storage
- **Hash Collision Management:** Separate chaining requires additional pointer storage proportional to collision frequency
- **String Storage:** Term strings are stored redundantly in both hash keys and potentially in posting lists

4.2.2 Space-Time Tradeoffs

The current implementation makes deliberate space-time tradeoffs:

- **Fixed Hash Table Size:** Provides $O(1)$ average lookup time at the cost of potentially underutilized memory
- **Linked List Posting Lists:** Enable efficient incremental updates while sacrificing cache locality
- **String Duplication:** Simplifies memory management but increases storage requirements

The current design represents a balanced approach that prioritizes query performance while maintaining reasonable space efficiency for medium-scale text collections.

5 Chapter 5: Bonus – Large-Scale Data Processing Capability Analysis

5.1 Hash Analysis

In the Bonus program, the only files modified compared to the aforementioned program are just hash_table.c and hash_table.h. So in the following section, we will only do the Hash Analysis.

5.1.1 Hash Function Design Advantages

The optimized 64-bit hash function implemented (based on an xxHash variant) offers the following technical advantages:

```
1 static uint64_t optimized_hash(const char *key, size_t len) {  
2     // Utilizing multiple large primes for mixing  
3     static const uint64_t prime1 = 11400714785074694791ULL;  
4     static const uint64_t prime2 = 14029467366897019727ULL;  
5     static const uint64_t prime3 = 1609587929392839161ULL;  
6     static const uint64_t prime4 = 9650029242287828579ULL;  
7     static const uint64_t prime5 = 2870177450012600261ULL;  
8     // Multi-stage processing for long and short keys  
9     // Avalanche effect and prime mixing operations  
10 }
```

Listing 1: Optimized 64-bit hash function skeleton (C)

Key Technical Characteristics:

- **High-Dimensional Hash Space:** 64-bit output provides approximately $2^{64} \approx 1.84 \times 10^{19}$ possible values, significantly reducing collision probability.
- **Avalanche Effect:** Minor input variations cause substantial changes in hash values, ensuring uniform distribution.
- **Multi-Stage Processing:** Different processing branches for long keys (e.g. ≥ 32 bytes) and short keys to improve throughput and mixing.
- **Prime Mixing:** Several mixing rounds using large primes to enhance randomness and diffusion.

5.1.2 Collision Probability Analysis

For $n = 400$ million $= 4 \times 10^8$ distinct words and a maximum table size of $m \approx 2^{31} - 1 \approx 2.1 \times 10^9$, using the birthday-approximation for collisions:

$$P(\text{collision}) \approx 1 - \exp\left(-\frac{n^2}{2m}\right). \quad (1)$$

Substituting the values:

$$\frac{n^2}{2m} = \frac{(4 \times 10^8)^2}{2 \times 2.1 \times 10^9} = \frac{1.6 \times 10^{17}}{4.2 \times 10^9} \approx 3.8095 \times 10^7.$$

Therefore, using the full 64-bit hash space ($M = 2^{64}$):

$$P \approx 1 - \exp\left(-\frac{n^2}{2M}\right) = 1 - \exp\left(-\frac{1.6 \times 10^{17}}{3.68 \times 10^{19}}\right) \approx 0.0043 \text{ (0.43\%)}. \quad (2)$$

The actual collision rate will be even lower due to:

- Excellent distribution of the xxHash variant;
- Dynamic resizing maintaining a low load factor;
- Collision resolution through chaining.

5.2 Core Insertion Algorithm

Algorithm 8 Hash Table Insertion Algorithm

Require: Hash table ht , key string key , value $value$

Ensure: Key-value pair inserted or updated in hash table

```

1: if  $ht = \text{NULL}$  or  $key = \text{NULL}$  then
2:   return
3: if  $(ht.\text{count} + 1)/ht.\text{size} > \text{LOAD\_FACTOR\_THRESHOLD}$  then
4:    $ht.\text{resize}(ht)$  ▷ Resize if load factor exceeds threshold
5:  $\text{hash\_value} \leftarrow \text{optimized\_hash}(key, \text{strlen}(key))$ 
6:  $\text{index} \leftarrow \text{hash\_value} \bmod ht.\text{size}$  ▷ Compute bucket index
7:  $\text{current} \leftarrow ht.\text{table}[\text{index}]$ 
8: while  $\text{current} \neq \text{NULL}$  do ▷ Search for existing key in bucket
9:   if  $\text{strcmp}(\text{current}.key, key) = 0$  then
10:     $\text{current}.value \leftarrow value$  ▷ Update existing key's value
11:   return
12:    $\text{current} \leftarrow \text{current}.next$ 
13:  $\text{new\_node} \leftarrow \text{allocate\_memory}(\text{sizeof}(\text{HashNode}))$ 
14:  $\text{new\_node}.key \leftarrow \text{copy\_string}(key)$ 
15:  $\text{new\_node}.value \leftarrow value$ 
16:  $\text{new\_node}.next \leftarrow ht.\text{table}[\text{index}]$  ▷ Insert at head of chain
17:  $ht.\text{table}[\text{index}] \leftarrow \text{new\_node}$ 
18:  $ht.\text{count} \leftarrow ht.\text{count} + 1$ 

```

5.3 Dynamic Resizing Algorithm

Algorithm 9 Hash Table Resizing Algorithm

Require: Hash table ht

Ensure: Hash table resized with new capacity

```

1:  $old\_size \leftarrow ht.size$ 
2:  $new\_size \leftarrow next\_prime(old\_size \times 2)$ 
3: if  $new\_size \leq old\_size$  then
4:   return FAILURE ▷ Maximum capacity reached
5:  $new\_table \leftarrow allocate\_zeroed\_memory(new\_size \times sizeof(HashNode*))$ 
6: for  $i \leftarrow 0$  to  $old\_size - 1$  do ▷ Rehash all elements
7:    $node \leftarrow ht.table[i]$ 
8:   while  $node \neq NULL$  do
9:      $next\_node \leftarrow node.next$  ▷ Save next pointer
10:     $new\_hash \leftarrow optimized\_hash(node.key, strlen(node.key))$ 
11:     $new\_index \leftarrow new\_hash \bmod new\_size$ 
12:     $node.next \leftarrow new\_table[new\_index]$  ▷ Insert into new bucket
13:     $new\_table[new\_index] \leftarrow node$ 
14:     $node \leftarrow next\_node$  ▷ Move to next node in chain
15:  $free(ht.table)$  ▷ Release old table memory
16:  $ht.table \leftarrow new\_table$ 
17:  $ht.size \leftarrow new\_size$ 
18: return SUCCESS

```

5.4 Performance Analysis

5.4.1 Time Complexity Analysis

Table 7: Time complexity summary

Operation	Average Case	Worst Case	Description
Insertion	$O(1)$	$O(n)$	Load factor control ensures average performance
Search	$O(1)$	$O(n)$	Uniform hash distribution maintains constant time
Deletion	$O(1)$	$O(n)$	Same complexity as search
Resizing	$O(n)$	$O(n)$	Amortized $O(1)$ per insertion due to exponential growth

5.4.2 Space Complexity Analysis

Memory Usage Estimation:

Node memory $\approx 400M \times (8B \times 3 + 10B) \approx 13.6$ GB.

Hash table memory $\approx 2.1B \times 8B \approx 16.8$ GB.

Total ≈ 30.4 GB.

Optimization Factors:

- Shared string storage reduces overhead.
- Allocator fragmentation may increase real usage.
- 64 to 512 GB servers can easily handle this load.

5.4.3 Throughput Prediction

- Insertion throughput: ≥ 1 million ops/s.
- Search throughput: ≥ 2 million ops/s.
- Resizing overhead: 15~20% total runtime.

5.5 Conclusion

5.5.1 Algorithm Selection Rationale

This hash table implementation ensures large-scale processing via:

- Proper capacity planning (supports 1.6B pairs)
- Advanced xxHash-based distribution
- Efficient chained memory layout
- Prime-based exponential resizing

5.5.2 Final Conclusion

The optimized hash table can handle 500K files and 400M unique words with strong performance and acceptable memory footprint. For production:

- Add distributed sharding and persistence;
- Implement fault recovery;
- Optimize for concurrent access.

6 Appendix

6.1 External Resources

- Stopwords List: https://github.com/CharyHong/Stopwords/blob/main/stopwords_english.txt
- Porter Stemmer: <https://github.com/woorm/stmr.c>
- Shakespeare Texts: <http://shakespeare.mit.edu/>