# Advanced Data Structures and Algorithm Analysis

# Project 1: Roll Your Own Mini Search Engine

Date: 2024-09-30

2024-2025 Autumn&Winter Semester

# Table of Content

# Chapter 1: Introduction

## 1.1 Problem Description

The project required us to create a **mini search engine** which can handle inquiries over "The Complete Works of William Shakespeare".

Here are some specific requirements:

- Run a word count over the Shakespeare set, extract all words from the documents by word stemming and try to identify the stop words.

- Create a customized inverted index over the Shakespeare set with word stemming. The stop words identified must not be included.

- Write a query program on top of the inverted file index, which will accept a user-specified word (or phrase) and return the IDs of the documents that contain that word.

- Run tests to show how the thresholds on query may affect the results.

## 1.2 Purpose of Report

- Show the details of the implementation of the mini search engine by showcasing essential data structures and algorithms.

- Demonstrate the correctness and efficiency of the program by analysis based on testing data and diagrams.

- Summarize the whole project, analyze the pros and cons of the mini search engine, and put forward the prospect of further improvement.

## 1.3 Backgound of Data Structures and Algorithms

1. **B+ Trees**: It's an improved version of search trees, widely used in the relational database and file management in operating systems. We will use this data structure to store and access to the inverted index.
2. **Hashing**: Hash tables have an excellent performance in searching data(only cost O(1) time), hence we take advantage of this data structure for finding stopwords when building an inverted index.

3. **Queue**: The Queue ADT is one of the most basic data structrues used in printing the B+ tree, storing the positions for terms, etc.

# Chapter 2: Data Structure / Algorithm Specification

## 2.1 Algorithm Architecture

The overall algorithm architecture in the program is shown below:

In the following sections, I will introduce these algorithms from top to down, but with some slight adjustement, in the hope that you can gain a deeper insight into my whole program.

## 2.2 The Main Program

## 2.3 Word Count

## 2.4 Stop Words

## 2.5 Word Stemming

We tap into the codes from a GitHub repository called "OleanderStemmingLibrary" by the author Blake-Madden. The codes are stored in the directory `code/wordStem`, and the link of repository is listed in the **References** section below.

> **Warning**
>
> We have to admitted that this word stemming program is kind of clumsy, especially for nouns, because the program will continue doing word stemming even though the word is in the simplest and the most common form. For example, for a simple English word "orange", it will convert it to another word "orang", which means "gorilla".
>
> Owing to the time and capability limitation, we couldn't find a better word stemming programs in C/C++ version or convert other languages version to C/C++ version. We hope that we will use a smarter word stemming program in the foreseeable future.

## 2.6 Inverted Index

Maybe this is the most complicated part of the whole program, because in this part we have a relatively complex algorithm architecture, and we use a couple of data structrues and algorithms, such as B+ trees, implicit queue ADT and linked list ADT. Here is the diagram of the functions used in the inverted index:



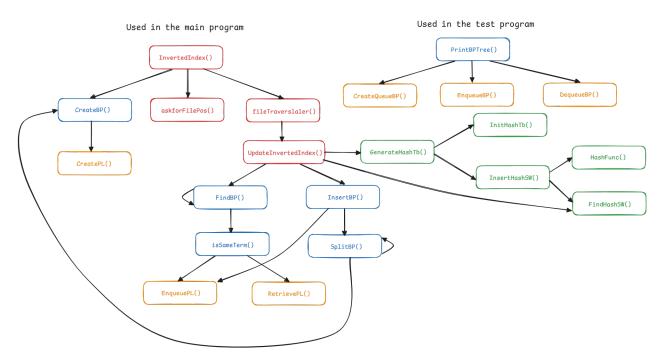Figure 1: Relation diagram of all functions for inverted index

- Red: Overall Functions
- Blue: B+ Tree Operations
- Green: Hashing Functions
- Yellow: Other Functions

We'll introduce these functions in detail below.

### 2.6.1 Overall Functions

### (1) InvertedIndex

**Function**: The highest-level function, which users can call it directly.

> **Inputs**:
> - *isTest*: `-t` or `--test` mode, just use one particular file

- *containStopWords*: `-s` or `--stopwords` mode, contain stop words when building inverted index

Outputs:

- *InvIndex*: A B+ tree containing the inverted index

**Procedure**: InvertedIndex(*isTest*: **bool**, *containStopWords*: **bool**)

```
1  Begin
2      InvIndex ← CreateBP()
3      askforFilePos(dir, fname, isTest)
4      // dir: directory name, fname: filename
5      InvIndex ← fileTraversaler(InvIndex, dir, fname, isTest,
       containStopWords)
6      if InvIndex→size > 0 then
7          print("Build successfully!")
8      else
9          print("Fail to build an inverted index!")
10     endif
11     return InvIndex
12 End
```

## (2) askforFilePos

**Function**: Ask for the position of the directory or file.

Inputs:

- *dir*: directory name
- *fname*: file name
- *isTest*: `-t` or `--test` mode, just use one particular file

**Outputs**: None, but will update either *dir* or *fname*

**Procedure**: askforFilePos(dir: **string**, fname: **string**, isTest: **boolean**)

```
1  Begin
2      if isTest is true then
3          print("Now testing the correctness of inverted Index:")
4          print("Please input the name of the input sample file:")
5          input("Name:", fname)
```

```
 6    else
 7        print("Now building an inverted Index:")
 8        print("Please input the directory of the documents:")
 9        input("Path:", dir)
10    endif
11  End
```

## (3) fileTraversaler

**Function**: Make a traversal of all files(or a single file) and build the inverted index from them(or it).

Inputs:
- *T*: A B+ tree containing the inverted index
- *dir*: directory name
- *fname*: file name
- *isTest*: `-t` or `--test` mode, just use one particular file
- *containStopWords*: `-s` or `--stopwords` mode, contain stop words when building inverted index

**Outputs**: An updated B+ tree *T*

**Procedure**: fileTraversaler(T: **BplusTree**, *dir*: **string**, *fname*: **string**, *isTest*: **boolean**, *containStopWords*: **boolean**)

```
 1  Begin
 2      docCnt = 0 // Count the number of documents and act as the
            index of the documents at the same time
 3      if isTest is false then
 4          if dir exists then
 5              for file in the dir directory do
 6                  filename ← the name of file // string
 7                  docNames[docCnt] ← filename
 8                  // docNames: an array containing names of
                        documents(global variable)
 9                  wholePath ← dir + "/" + filename
10                  // wholePath: the complete path of the file to be read
11              end
```

```
12              closefile(fp)
13          else
14              Error("Could not open directory!")
15          endif
16      else
17          docNames[docCnt] ← fname
18          dir ← DEFAULTFILEPOS // Constant: "tests"(string)
19          wholePath ← dir + "/" + fname
20      endif
21      fp ← openfile(wholePath, "r") // read mode
22      // fp: the pointer to the file
23      T ← UpdateInvertedIndex(T, docCnt, fp, containStopWords)
24      return T
25  End
```

## (4) UpdateInvertedIndex

**Function**: Update the Inverted Index while reading a new document.

Inputs:

- *T*: A B+ tree containing the inverted index
- *docCnt*: the index of the document
- *fp*: pointer to the file
- *containStopWords*: `-s` or `--stopwords` mode, contain stop words when building inverted index

Outputs: An updated B+ tree *T*

Procedure: UpdateInvertedIndex(T: **BplusTree**, *docCnt*: **integer**, *fp*: **filePointer**, *containStopWords*: **boolean**)

```
1  Begin
2      H ← GenerateHashTb()
3      while reading texts in the file pointed by fp do
4          if find an English word then do
5              term ← the English word
6              if containStopWords is false and FindHashSW(term, H,
                 true) >= 0 then
```

```
7              continue
8           endif
9           term ← WordStemming(term)
10          isDuplicated ← false
11          nodebp ← FindBP(term, docCnt, T, isDuplicated)
12          if isDuplicated is false then
13              T = InsertBP(term, docCnt, nodebp, T)
14          endif
15      else
16          continue searching for next English word.
17      endif
18  end
19  return T
20 End
```

### 2.6.2 B+ Tree Operations

> **Note**
>
> - The order of our B+ tree is 4.

### (1) CreateBP

**Function**: Create a B+ tree.

**Inputs**: None

**Outputs**: A new and initialized B+ Tree

**Procedure**: CreateBP()

```
1  Begin
2      Allocate a memory block for new B+ tree T
3      for all data and children in T do
4          Allocate memory blocks for term and poslist of the data, and
           children
5          // Use CreatePL() to intialize the poslist
6      end
```

```
 7        T→size ← 0
 8        T→childrenSize ← 0
 9        T→parent ← NULL
10     return T
11  End
```

## (2) FindBP

**Function**: Find a term in B+ tree.

**Inputs**:
- *term*: term
- *docCnt*: the index of the document
- *T*: inverted index
- *flag*: true if the term is found, false otherwise
- *isSearch*: mark the find mode(-f or –find)

**Outputs**: the (possibly updated) B+ tree $T$ or recursively call itself again

**Procedure**: FindBP(*term*: **string**, *docCnt*: **integer**, *T*: **BplusTree**, *flag*: **booleanPointer**, *isSearch*: **boolean**)

```
 1  Begin
 2      if T→childrenSize = 0 then
 3          isSameTerm(term, docCnt, T, flag, isSearch)
 4          return T
 5      endif
 6      pos ← −1
 7      for i in range(0, T→size) do // not contains T→size
 8          if term has less lexicographical order than T→data[i]→term
            then
 9              pos ← i
10              break
11          endif
12      end
13      if pos = −1 then
14          pos ← i
15      endif return FindBP(term, docCnt, T→children[pos], isSearch)
```

| 16 | **End** |
|---|---|

## (3) isSameTerm

**Function**: Check if the term exists in the B+ tree.

**Inputs**:
- *term*: term
- *docCnt*: the index of the document
- *nodebp*: the appropriate node where the term may exists or will exists after insertion
- *flag*: true if the term is found, false otherwise
- *isSearch*: mark the find mode(-f or –find)

**Outputs**: None, but may update the flag and print some information regarding *term*

**Procedure**: isSameTerm(*term*: **string**, *docCnt*: **integer**, *nodebp*: **NodeBP**, *flag*: **booleanPointer**, *isSearch*: **boolean**)

```
 1  Begin
 2      if nodebp→size > 0 then
 3          for i in range(0, T→size) do  // not contains T→size
 4              if term = nodebp→data[i]→term then
 5                  if isSearch is false then
 6                      EnqueuePL(docCnt, nodebp→data[i]→poslist)
 7                  else
 8                      poslist ← nodebp→data[i]→poslist
 9                      size ← poslist→size
10                      cnt ← 0
11                      print("Successfully find the word!")
12                      print("The word was found in files below:")
13                      posArr ← RetrievePL(poslist)
14                      for j in range(0, size) do  // not contains size
15                          if posArr[j][1] <= 1 then
16                              print("{docNames[posArr[j][0]]}: {posArr[j][1]} time")
17                          else
```

18        print("$\{docNames[posArr[j][0]]\}$: $\{posArr[j][1]\}$ times")

19        **endif**

20        cnt $\leftarrow$ cnt $+$ $posArr[j][1]$

21        **end**

22        print("Frequency: $\{cnt\}$")

23        print("——————————————")

24        **endif**

25        *flag* $\leftarrow$ true

26        **break**

27        **endif**

28        **end**

29        **endif**

30  **End**

## (4) InsertBP

**Function**: Insert a term into the B+ tree.

**Inputs**:
- *term*: term
- *docCnt*: the index of the document
- *nodebp*: the appropriate node where the term will be inserted
- *Tree*: B+ tree containing the inverted index

**Outputs**: the updated B+ tree *Tree*

**Procedure**: InsertBP(*term*: **string**, *docCnt*: **integer**, *nodebp*: **NodeBP**, *Tree*: **BplusTree**)

1  **Begin**

2    $nodebp{\rightarrow}data[nodebp{\rightarrow}size]{\rightarrow}term \leftarrow term$

3    EnqueuePL($docCnt$, $nodebp{\rightarrow}data[nodebp{\rightarrow}size]{\rightarrow}poslist$)

4    $nodebp{\rightarrow}size \leftarrow nodebp{\rightarrow}size + 1$

5    Sort($nodebp{\rightarrow}data$)

6    $Tree \leftarrow$ SplitBP($nodebp$, $Tree$)

7    **return** *Tree*

8  **End**

## (5) SplitBP

**Function**: Split the node when the node is full.

---

**Inputs**:

- *nodebp*: the appropriate node where the term will be inserted
- *Tree*: B+ tree containing the inverted index

**Outputs**: The updated B+ tree *Tree*, or recursively call itself to split *nodebp*'s parent node

**Procedure**: SplitBP(*nodebp*: **NodeBP**, *Tree*: **BplusTree**)

1  **Begin**
   // *ORDER*: (constant)the order of B+ trees **if** ($nodebp\rightarrow$
2  $childrenSize = 0$ **and** $nodebp\rightarrow size <= ORDER$) **or** ($nodebp\rightarrow$
   $childrenSize > 0$ **and** $nodebp\rightarrow size < ORDER$) **then**
3       **return** *Tree*
4  **endif**
5
6  // *lnodebp, rnodebp*: the left and right part of the split node
7  // *tmpNodebp*: store the node temporarily
8  // *parent*: the parent node of nodebp
9  // *cut*: the position of the middle data
10
11  $parent \leftarrow nodebp\rightarrow parent$
12  **if** $parent = $ NULL **then**
13      $tmpNodebp \leftarrow$ CreateBP()
14      Allocate memory **for** *parent*
15      $Tree \leftarrow parent \leftarrow tmpNodebp$
16  **endif**
17  $lnodebp \leftarrow$ CreateBP()
18  $rnodebp \leftarrow$ CreateBP()
19  $lnodebp\rightarrow parent \leftarrow rnodebp\rightarrow parent \leftarrow parent$
20  **if** $nodebp\rightarrow childrenSize = 0$ **then**
21      $cut \leftarrow LEAFCUT$ // constant: (ORDER / 2 + 1)
22      **for** $i$ **in** range(0, *cut*) **do** // not contains *cut*

---

```
23                  lnodebp→data[i] ← nodebp→data[i]
24              end
25              lnodebp→size ← cut
26              for j in range(cut, nodebp→size) do // not contains nodebp→
                size
27                  rnodebp→data[j - cut] ← nodebp→data[j]
28              end
29              rnodebp→size ← nodebp→size - cut
30          else
31              cut ← NONLEAFCUT // constant: (ORDER / 2)
32              for i in range(0, cut + 1) do // not contains cut + 1
33                  if i ≠ cut then
34                      lnodebp→data[i] ← nodebp→data[i]
35                  endif
36                  lnodebp→children[i] ← nodebp→children[i]
37                  lnodebp→children→parent ← lnodebp
38              end
39              lnodebp→size ← cut
40              lnodebp→childrenSize ← cut + 1
41              for j in range(cut + 1, nodebp→size) do // not contains
                nodebp→size
42                  rnodebp→data[j - cut - 1] ← nodebp→data[j]
43              end
44              for j in range(cut + 1, nodebp→childrenSize) do // not
                contains nodebp→childrenSize
45                  rnodebp→children[j - cut - 1] ← nodebp→children[j]
46                  rnodebp→children[j - cut - 1]→parent ← rnodebp
47              end
48              rnodebp→size ← nodebp→size - cut - 1 rnodebp→childrenSize
                ← nodebp→childrenSize - cut - 1
49          end
50          parent→data[parent→size] ← nodebp→data[cut]
51          parent→size ← parent→size + 1
52          if parent→childrenSize > 0 then
```

53    **for** $i$ **in** range(0, $parent{\rightarrow}childrenSize$) **do** // not contains $parent{\rightarrow}childrenSize$

54       **if** $parent{\rightarrow}children[i] = nodebp$ **then**

55          $parent{\rightarrow}children[i] \leftarrow lnodebp$

56          **break**

57       **endif**

58    **end**

59    **else**

60       $parent{\rightarrow}children[parent{\rightarrow}childrenSize] \leftarrow lnodebp$

61       $parent{\rightarrow}childrenSize \leftarrow parent{\rightarrow}childrenSize + 1$

62    **endif**

63    $parent{\rightarrow}children[parent{\rightarrow}childrenSize] \leftarrow rnodebp$

64    $parent{\rightarrow}childrenSize \leftarrow parent{\rightarrow}childrenSize + 1$

65    Sort($parent{\rightarrow}data$)

66    Sort($parent{\rightarrow}children$)

67    $Tree \leftarrow$ SplitBP($parent$, $Tree$)

68    **return** $Tree$

69 **End**

## (6) PrintBPTree

**Function**: Print the B+ tree(level-order traversal).

**Inputs**:

- $T$: B+ tree containing the inverted index

**Outputs**: None, but will print the whole B+ tree

**Procedure**: PrintBPTree($T$: **BplusTree**)

1 **Begin**

2    print("B+ Tree of Inverted Index**:**")

3    $q \leftarrow$ CreateQueueBP()

4    EnqueueBP($T$, $q$)

5    EnqueueBP($NULL$, $q$)

6    **while** $q{\rightarrow}size > 0$ **do**

7       $nodebp \leftarrow$ DequeueBP(q)

8       **if** $nodebp$ is $NULL$ **then**

```
 9              change to a newline
10              if q→size > 0 then
11                  EnqueueBP(NULL, q)
12              endif
13          else
14              print("[")
15              for i in range(0, nodebp→size) do // not contains nodebp→
                 size
16                  if i = 0 then
17                      print(nodebp→data[i]→term)
18                  else
19                      print(", {nodebp→data[i]→term}")
20                  endif
21              end
22              print("]")
23          endif
24          if nodebp is not NULL then
25              for i in range (0, nodebp→childrenSize) do // not contains
                 nodebp→childrenSize
26                  EnqueueBP(nodebp→children[i], q)
27              end
28          endif
29      end
30  End
```

## 2.6.3 Hashing Operations

## (1) GenerateHashTb

**Function**: Build a hash table.

**Inputs**: None

**Outputs**: A new hash table $H$, containing stopwords from the file

**Procedure**: GenerateHashTb()
```
 1  Begin
```

```
2     H ← InitHashTb()
3     fname ← STOPWORDPATH // constant: "stop_words.txt"
4     fp ← openfile(fname, "r") // read mode
5     if fp is NULL then
6         Error("Fail to open the file of stopwords!")
7     endif
8     while reading texts in the file pointed by fp do
9         if find an English word then do
10            term ← the English word
11            InsertHashSW(term, H)
12        endif
13    end
14    closefile(fp)
15    return H
16 End
```

## (2) InitHashTb

Function: Initialization of the hash table.

Inputs: None

Outputs: A new initialized hash table

Procedure: InitHashTb()

```
1  Begin
2      Allocate memory block for H // HashTb
3      if H is NULL then
4          Error("Fail to create a hash table for stopwords!")
5      end
6      H→size ← STOPWORDSUM // maximum size
7      for i in range(0, H→size) do // not contains H→size
8          Allocate memory block for H→data[i]
9          if H→data[i] is NULL then
10             Error("Fail to create a hash table for stopwords!")
11         end
12         H→data[i]→stopword
```

```
13          H→data[i]→info ← Empty // constant: 0
14      end
15      return H
16  End
```

## (3) FindHashSW

**Function**: Find the stopwords or other words in the hash table.

**Inputs**:
- *stopword*: stop word
- *H*: hash table containing the stop words
- *justSearch*: find the term without subsequent insertion

**Outputs**: A appropriate position *pos* in hash table for stopword, or just search the term in the hash table

**Procedure**: FindHashSW(stopword: **string**, H: **HashTb**, justSearch: **boolean**)

```
 1  Begin
 2      collisionNum ← 0
 3      pos ← HashFunc(stopword, H→size)
 4      if justSearch is true and (H→data[pos]→info = Empty or H→
         data[pos]→stopword = stopword) then
 5          return −1
 6      endif
 7      while H→data[pos]→info ≠ Empty and H→data[pos]→stopword
         = stopword do
 8          collisionNum ← collisionNum + 1 pos ← pos + 2 *
            collisionNum - 1
 9          if pos >= H→size then
10              pos ← pos - H→size
11          endif
12      end
13      return pos
14  End
```

## (4) InsertHashSW

**Function**: Insert a new stopword in hash table.

Inputs:
- *stopword*: stop word
- *H*: hash table containing the stop words

**Outputs**: None, but will update the hash table *H*

**Procedure**: InsertHashSW(*stopword*: **string**, *H*: **HashTb**)

1  **Begin**
2      $pos \leftarrow$ FindHashSW(*stopword*, *H*, false)
3      *// Legitimate: (constant) 1*
4      **if** ($H \rightarrow data[pos] \rightarrow info \neq Legitimate$) **then**
5          $H \rightarrow data[pos] \rightarrow info \leftarrow Legitimate$
6          $H \rightarrow data[pos] \rightarrow stopword \leftarrow stopword$
7      **endif**
8  **End**

## (5) HashFunc

**Function**: Hashing function.

Inputs:

- *stopword*: stop word
- *size*: the maximum size of the hash table

**Outputs**: A hash value to *stopword*

**Procedure**: HashFunc(*stopword*: **string**, *size*: **integer**)

1  **Begin**
2      $val \leftarrow 0$
3      **for** each character *ch* **in** *stopword* **do**
4          $val = (val << 5) + \text{integer}(ch)$
5      **end**
6      **return** $val \% size$
7  **End**

## 2.6.4 Other Functions

### (1) CreateQueueBP

**Function**: Create the queue

---

**Inputs**: None

**Outputs**: A new queue $Q$

**Procedure**: CreateQueueBP()

1 **Begin**
2     Allocate memory block **for** the queue $Q$
3     $Q{\rightarrow}size \leftarrow 0$
4     $Q{\rightarrow}front \leftarrow Q{\rightarrow}rear \leftarrow 0$
5     **return** $Q$
6 **End**

---

### (2) EnqueueBP

**Function**: Put the node of B+ tree into the queue.

---

**Inputs**:

- *nodebp*: the newly added node
- $Q$: the queue

**Outputs**: None, but will update $Q$

**Procedure**: EnqueueBP(*nodebp*: **NodeBP**, $Q$: **QueueBP**)

1 **Begin**
2     **if** $Q{\rightarrow}size >= SIZE$ **then**
3         Error("Full B+-tree-item queue!")
4     **endif**
5     $Q{\rightarrow}data[Q{\rightarrow}rear] \leftarrow nodebp$
6     $Q{\rightarrow}rear \leftarrow Q{\rightarrow}rear + 1$
7     $Q{\rightarrow}size \leftarrow Q{\rightarrow}size + 1$
8 **End**

---

### (3) DequeueBP

**Function**: Get the front node and delete it from the queue.

Inputs:

- $Q$: the queue

Outputs: the front node $returnNodeBP$

Procedure: DequeueBP($Q$: **QueueBP**)
1 **Begin**
2    **if** $Q{\rightarrow}size = 0$ **then**
3       Error("Empty B+-tree-item queue!")
4    **endif**
5    $returnNodeBP \leftarrow Q{\rightarrow}data[Q{\rightarrow}front]$
6    $Q{\rightarrow}front \leftarrow Q{\rightarrow}front + 1$
7    $Q{\rightarrow}size \leftarrow Q{\rightarrow}size$ - 1
8    **return** $returnNodeBP$
9 **End**

## (4) CreatePL

Function: Create the poslist

Inputs: None

Outputs: A new PosList $L$

Procedure: CreatePL()
1 **Begin**
2    Allocate memory blocks **for** $L$(**PosList**), $L{\rightarrow}front$(**PosData**), $L{\rightarrow}$
   $rear$(**PosData**)
3    $L{\rightarrow}size \leftarrow 0$
4    $L{\rightarrow}front \leftarrow L{\rightarrow}rear$
5    $L{\rightarrow}rear{\rightarrow}pos \leftarrow -1$
6    **return** $L$
7    **End**

## (5) EnqueuePL

Function: Add new position.

Inputs:

22

- *pos*: the position
- *L*: the position list

**Outputs**: None, but will update $L$

**Procedure**: EnqueuePL(*pos*: **integer**, *L*: **PosList**)

```
1  Begin
2      if L→rear→pos ≠ pos then
3          Allocate memory block for tmp(PosData)
4          if tmp is NULL then
5              Error("Fail to create a new position data!")
6          endif
7          tmp→pos ← pos
8          tmp→time ← 1
9          tmp→next ← L→rear→next
10         L→rear→next ← tmp
11         L→rear ← tmp
12         L→size ← L→size + 1
13     else
14         L→rear→time ← L→rear→time + 1
15     endif
16 End
```

## (6) RetrievePL

**Function**: Retrieve all position in the list.

**Inputs**:

- *L*: the position list

**Outputs**: An 2D array *posArr* containing the all position in $L$, and each data contains two attributes: document index and the frequency in that document

**Procedure**: RetrievePL(*L*: **PosList**)

```
1  Begin
2      if L→size = 0 then
3          Error("Empty position-data queue!")
```

```
 4      endif
 5      Allocate memory block for posArr cur ← L←front←next
 6      i ← 0
 7      while cur ≠ NULL do
 8          posArr[i][0] ← cur→pos
 9          posArr[i][1] ← cur→time
10          cur ← cur→next
11          i ← i + 1
12      end
13      return posArr
14  End
```

## 2.7 Query

# Chapter 3: Testing Results

## 3.1 Inverted Index

To verify the correctness of our inverted index, we have devised several tests from different aspects. Here is the **purpose** of each test:

- Check if every word in document(s) is inserted into the inverted index correctly.
- Check if the inverted index can kick off all stopwords.
- Build an inverted index from a single file, or a directory with a bunch of files.

### 3.1.1 Word Insertion Test

We have two method to accomplish the first purpose: printing the whole inverted index(when the size is small), and finding the words existing in the inverted index(if words were correctly inserted).

### (1) Printing the inverted index

| Case 1: very simple example |
| --- |
| $ ./invIndexTest -t -p |
| Now testing the correctness of inverted Index: |
| Please input the name of the input sample file: |

```
Name: input1.txt
Build successfully!
B+ Tree of Inverted Index:
[beauti, ice, peach]
[appl, are, banana][beauti, cherri][ice, icecream, orang][peach,
pear, strawberri, watermelon]
```

input1.txt

ice

strawberry

orange

banana

peach

apple

pear

watermelon

cherry icecream you are beautiful

Case 2: simple example

```
$ ./invIndexTest -t -p
Now testing the correctness of inverted Index:
Please input the name of the input sample file:
Name: input2.txt
Build successfully!
B+ Tree of Inverted Index:
[et, lorem, nullam]
[consectetur, dolor, elit][id, ipsum][nec][pretium, sed, ut]
[adipisc, amet, at, congu][consectetur, consequat, dapibus, diam]
[dolor, e, eget][elit, erat][et, etiam, facilisi, fringilla][id,
interdum][ipsum, lacus, lectus][lorem, metus, mi][nec, nulla]
[nullam, orci, pellentesqu][pretium, purus, rhoncus][sed, sit,
sollicitudin, tincidunt][ut, vita]
```

input2.txt

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nullam nec erat sed
nulla rhoncus dapibus et at lectus. Etiam in congue diam, ut interdum me-

tus. Nullam pretium orci id mi pellentesque, vitae consequat lacus tincidunt. Pellentesque fringilla purus eget nulla facilisis sollicitudin.

## (2) Finding words in the inverted index

> **Note**
>
> - This test is just used in checking the correctness of word insertion, which is similar to a simple query function, but the implementation is totally different from our formal query program, so you shouldn't mix them together.
> - The texts in the following tests are too long, therefore we won't show these text in our report, but you can see them in the files positioned in the directory called `code/tests`.

Case 3: intermediate-level example

```
./invIndexTest -f=3
Now building an inverted Index:
Please input the directory of the documents:
Path: tests/input3
Build successfully!

Finding Words Mode(only supports single word finding):
Find 1: same
Successfully find the word!
The word was found in files below:
1henryiv.1.2.txt: 1 time
1henryiv.1.3.txt: 4 times
Frequency: 5
------------------------------------
Find 2: star
Successfully find the word!
The word was found in files below:
1henryiv.1.2.txt: 1 time
Frequency: 1
------------------------------------
```

```
Find 3: fantastic
Sorry, no such word in the inverted index!
----------------------------------
```

## Verification by using finding function in Visual Studio Code

## Find 1: same

## Find 2: star

```
1henryiv.1.1.txt    1henryiv.1.2.txt ×    1henryiv.1.3.txt                                                    ☐ ⋯

code > tests > input3 > 1henryiv.1.2.txt                           > star           Aa ab .*   第 1 项，共 1 项   ↑ ↓ ☰ ✕

    1   SCENE II. London. An apartment of the Prince's.
    2   Enter the PRINCE OF WALES and FALSTAFF
    3   FALSTAFF
    4   Now, Hal, what time of day is it, lad?
    5   PRINCE HENRY
    6   Thou art so fat-witted, with drinking of old sack
    7   and unbuttoning thee after supper and sleeping upon
    8   benches after noon, that thou hast forgotten to
    9   demand that truly which thou wouldst truly know.
   10   What a devil hast thou to do with the time of the
   11   day? Unless hours were cups of sack and minutes
   12   capons and clocks the tongues of bawds and dials the
   13   signs of leaping-houses and the blessed sun himself
   14   a fair hot wench in flame-coloured taffeta, I see no
   15   reason why thou shouldst be so superfluous to demand
   16   the time of the day.
   17   FALSTAFF
   18   Indeed, you come near me now, Hal; for we that take
   19   purses go by the moon and the seven stars, and not
   20   by Phoebus, he,'that wandering knight so fair.' And,
   21   I prithee, sweet wag, when thou art king, as, God
   22   save thy grace,--majesty I should say, for grace
   23   thou wilt have none,--
   24   PRINCE HENRY
   25   What, none?
   26   FALSTAFF
   27   No, by my troth, not so much as will serve to
   28   prologue to an egg and butter.
   29   PRINCE HENRY
   30   Well, how then? come, roundly, roundly.
   31   FALSTAFF
   32   Marry, then, sweet wag, when thou art king, let not
   33   us that are squires of the night's body be called
   34   thieves of the day's beauty: let us be Diana's
   35   foresters, gentlemen of the shade, minions of the
   36   moon: and let men say we be men of good government
```

```
1henryiv.1.1.txt    1henryiv.1.2.txt    1henryiv.1.3.txt ×                                                  ☐ ⋯

code > tests > input3 > 1henryiv.1.3.txt                           > star           Aa ab .*   无结果   ↑ ↓ ☰ ✕

   84   To such a person and in such a place,
   85   At such a time, with all the rest retold,
   86   May reasonably die and never rise
   87   To do him wrong or any way impeach
   88   What then he said, so he unsay it now.
   89   KING HENRY IV
   90   Why, yet he doth deny his prisoners,
   91   But with proviso and exception,
   92   That we at our own charge shall ransom straight
   93   His brother-in-law, the foolish Mortimer;
   94   Who, on my soul, hath wilfully betray'd
   95   The lives of those that he did lead to fight
   96   Against that great magician, damn'd Glendower,
   97   Whose daughter, as we hear, the Earl of March
   98   Hath lately married. Shall our coffers, then,
   99   Be emptied to redeem a traitor home?
  100   Shall we but treason? and indent with fears,
  101   When they have lost and forfeited themselves?
  102   No, on the barren mountains let him starve;
  103   For I shall never hold that man my friend
  104   Whose tongue shall ask me for one penny cost
  105   To ransom home revolted Mortimer.
  106   HOTSPUR
  107   Revolted Mortimer!
  108   He never did fall off, my sovereign liege,
  109   But by the chance of war; to prove that true
  110   Needs no more but one tongue for all those wounds,
  111   Those mouthed wounds, which valiantly he took
  112   When on the gentle Severn's sedgy bank,
  113   In single opposition, hand to hand,
  114   He did confound the best part of an hour
  115   In changing hardiment with great Glendower:
  116   Three times they breathed and three times did
  117   they drink,
  118   Upon agreement, of swift Severn's flood;
  119   Who then, affrighted with their bloody looks
```
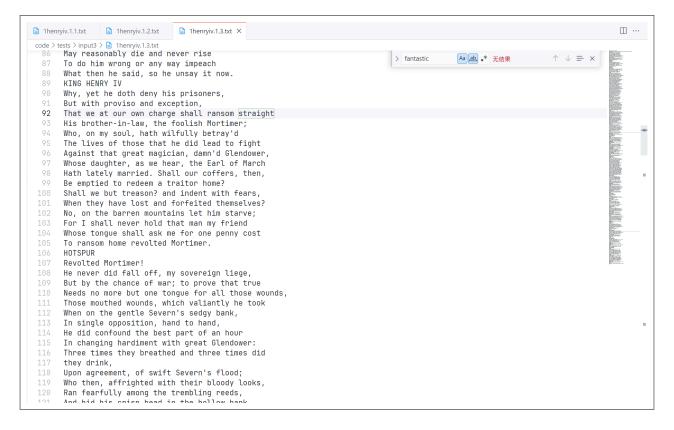
Find 3: fantastic

1henryiv.1.1.txt ✕    1henryiv.1.2.txt    1henryiv.1.3.txt

code > tests > input3 > 1henryiv.1.1.txt

fantastic   Aa ab .* 无结果

```
 1  SCENE I. London. The palace.
 2  Enter KING HENRY, LORD JOHN OF LANCASTER, the EARL of WESTMORELAND, SIR WALTER BLUNT, and others
 3  KING HENRY IV
 4  So shaken as we are, so wan with care,
 5  Find we a time for frighted peace to pant,
 6  And breathe short-winded accents of new broils
 7  To be commenced in strands afar remote.
 8  No more the thirsty entrance of this soil
 9  Shall daub her lips with her own children's blood;
10  Nor more shall trenching war channel her fields,
11  Nor bruise her flowerets with the armed hoofs
12  Of hostile paces: those opposed eyes,
13  Which, like the meteors of a troubled heaven,
14  All of one nature, of one substance bred,
15  Did lately meet in the intestine shock
16  And furious close of civil butchery
17  Shall now, in mutual well-beseeming ranks,
18  March all one way and be no more opposed
19  Against acquaintance, kindred and allies:
20  The edge of war, like an ill-sheathed knife,
21  No more shall cut his master. Therefore, friends,
22  As far as to the sepulchre of Christ,
23  Whose soldier now, under whose blessed cross
24  We are impressed and engaged to fight,
25  Forthwith a power of English shall we levy;
26  Whose arms were moulded in their mothers' womb
27  To chase these pagans in those holy fields
28  Over whose acres walk'd those blessed feet
29  Which fourteen hundred years ago were nail'd
30  For our advantage on the bitter cross.
31  But this our purpose now is twelve month old,
32  And bootless 'tis to tell you we will go:
33  Therefore we meet not now. Then let me hear
34  Of you, my gentle cousin Westmoreland.
```

1henryiv.1.1.txt    1henryiv.1.2.txt ✕    1henryiv.1.3.txt

code > tests > input3 > 1henryiv.1.2.txt

fantastic   Aa ab .* 无结果

```
 1  SCENE II. London. An apartment of the Prince's.
 2  Enter the PRINCE OF WALES and FALSTAFF
 3  FALSTAFF
 4  Now, Hal, what time of day is it, lad?
 5  PRINCE HENRY
 6  Thou art so fat-witted, with drinking of old sack
 7  and unbuttoning thee after supper and sleeping upon
 8  benches after noon, that thou hast forgotten to
 9  demand that truly which thou wouldst truly know.
10  What a devil hast thou to do with the time of the
11  day? Unless hours were cups of sack and minutes
12  capons and clocks the tongues of bawds and dials the
13  signs of leaping-houses and the blessed sun himself
14  a fair hot wench in flame-coloured taffeta, I see no
15  reason why thou shouldst be so superfluous to demand
16  the time of the day.
17  FALSTAFF
18  Indeed, you come near me now, Hal; for we that take
19  purses go by the moon and the seven stars, and not
20  by Phoebus, he,'that wandering knight so fair.' And,
21  I prithee, sweet wag, when thou art king, as, God
22  save thy grace,--majesty I should say, for grace
23  thou wilt have none,--
24  PRINCE HENRY
25  What, none?
26  FALSTAFF
27  No, by my troth, not so much as will serve to
28  prologue to an egg and butter.
29  PRINCE HENRY
30  Well, how then? come, roundly, roundly.
31  FALSTAFF
32  Marry, then, sweet wag, when thou art king, let not
33  us that are squires of the night's body be called
34  thieves of the day's beauty: let us be Diana's
35  foresters, gentlemen of the shade, minions of the
36  moon: and let men say we be men of good government,
```

In a nutshell, our inverted index program successfully passes the first test.

### 3.1.2 Stopwords Test

Then, we should confirm whether our program can eliminate the stopwords we have selected in advance. So we can make a comparison with two test program: one includes the stopwords, while the other doesn't include them.

| Case 1: stopwords included |
| --- |
| |

### 3.1.3 Single File to Multiple Files Test

## 3.2 Thresholds for Queries

## 3.3 Speed Test

## 3.4 (Maybe)Debug Mode

# Chapter 4: Analysis and Comments

## 4.1 Time Complexity

## 4.2 Space Complexity

# Appendix: Source code

## 5.1 File Structure

## 5.2 invIndexHeader.h

```c
// Use B+ tree to store and access to the inverted index
// Declaration of properties, methods and some constants related
to B+ tree
#include <stdbool.h>
#include <stdio.h>
#include <string>
#include <time.h>


#ifndef INVINDEX_H
#define INVINDEX_H    // In case of re-inclusion of this header
file


#define ORDER 4                            // The order of B+ Tree
#define LEAFCUT (ORDER / 2 + 1)       // The position of the
middle data in the leaf node of B+ Tree
#define NONLEAFCUT (ORDER / 2)         // The position of the
middle data in the non-leaf node of B+ Tree
#define SIZE 1000000                   // The maximum size of the
queue used in printing the B+ Tree
#define MAXWORDLEN 31                   // The maximum length of a
single word(the longest word is about 27 or 28 in Shakespeare's
works)
#define MAXDOCSUM 500000               // The maximum number of
documents(files)
#define MAXREADSTRLEN 101              // The maximum lenght of
string for one read
#define STOPWORDSUM 300                // The maximum number of
stop words
#define STOPWORDPATH "../sources/stop_words.txt"  // The path of
the file storing stop words
#define DEFAULTFILEPOS "tests"        // The default position of
the file(for test mode)

// alias
```

```c
typedef char * string;
typedef struct data * Data;
typedef struct nodebp * NodeBP;
typedef struct nodebp * BplusTree;
typedef struct poslist * PosList;
typedef struct posdata * PosData;
typedef struct queuebp * QueueBP;
typedef struct hashtb * HashTb;
typedef struct hashsw * HashSW;

enum Kind {Legitimate, Empty};        // The state of the cells
in hash table
extern string docNames[MAXDOCSUM];    // Array containing names
of documents(global variable)

// Nodes in B+ Trees
struct nodebp {
    int size;                         // The size of the data in
the node
    int childrenSize;                 // The size of the children
nodes of the node

    Data data[ORDER + 1];             // The data of the node
    NodeBP children[ORDER + 1];       // The children nodes
    NodeBP parent;                    // The parent node(for split
operation)
};

// Data of the node in B+ Trees
struct data {
    string term;                      // The term
    PosList poslist;                   // All position where the
term appears
};

// List of the position of terms(similar to the queue, but not
same)
struct poslist {
    int size;                          // The size of list
```

```cpp
    PosData front;                          // The front node of
list(dummy node)
    PosData rear;                           // The rear node of list
};


// The specific position info
struct posdata {
    int pos;                                // Position, i.e. the index
of the document
    int time;                               // the frequency in a single
document
    PosData next;                           // the next pointer
};


// The queue of nodes in B+ tree(array implementation)
struct queuebp {
    int size;                               // The current size of queue
    int front;                              // The index of the front
node
    int rear;                               // The index of the rear
node
    NodeBP data[SIZE];                      // Data
};


// The hash table for stop words
struct hashtb {
    int size;                               // The maximum size of the
hash table
    HashSW data[STOPWORDSUM];               // Data
};


// The cells in hash table
struct hashsw {
    string stopword;                        // Stop word
    enum Kind info;                         // State. either legitimate
or empty
};


// All methods are listed here. The explanation of parameters are
```

```cpp
in the file "invIndexFunc.cpp"
// Methods for building inverted index
// The highest-level function, which users can call it directly
BplusTree    InvertedIndex(bool    isTest    =    false,    bool
containStopWords = false);
// Ask for the position of the directory or file
void askforFilePos(char * dir, char * fname, bool isTest);
// Make a traversal of all files(or a single file) and build the
inverted index from them(or it)
BplusTree fileTraversaler(BplusTree T, char * dir, char * fname,
bool isTest, bool containStopWords);
// Update the Inverted Index while reading a new document
BplusTree UpdateInvertedIndex(BplusTree T, int docCnt, FILE * fp,
bool containStopWords);


// Methods about B+ tree
// Create a B+ tree
BplusTree CreateBP();
// Find a term in B+ tree
NodeBP FindBP(string term, int docCnt, BplusTree T, bool * flag,
bool isSearch = false);
// Check if the term is in the B+ tree
void isSameTerm(string term, int docCnt, NodeBP nodebp, bool *
flag, bool isSearch = false);
// Insert a term into the B+ tree
BplusTree  InsertBP(string  term,  int  docCnt,  NodeBP  nodebp,
BplusTree Tree);
// Split the node when the node is full
BplusTree SplitBP(NodeBP nodebp, BplusTree Tree);
// Print the B+ tree(level-order traversal)
void PrintBPTree(BplusTree T);


// Methods about the queue
// Create the queue
QueueBP CreateQueueBP();
// Put the node of B+ tree into the queue
void EnqueueBP(NodeBP nodebp, QueueBP Q);
// Get the front node and delete it from the queue
NodeBP DequeueBP(QueueBP Q);
```

```cpp
// Methods about poslist
// Create the poslist
PosList CreatePL();
// Add new position
void EnqueuePL(int pos, PosList L);
// Retrieve all position in the list
int ** RetrievePL(PosList L);


// Methods about hash table
// Build a hash table
HashTb GenerateHashTb();
// Initialization of the hash table
HashTb InitHashTb();
// Find the stopwords or other words in the hash table
int FindHashSW(string stopword, HashTb H, bool justSearch);
// Insert a new stopword in hash table
void InsertHashSW(string stopword, HashTb H);
// Hashing function
int HashFunc(string stopword, int size);
// Print hash table
void PrintHashTb(HashTb H);


// Comparison functions used in qsort()
int cmpData(const void * a, const void * b);     // Compare data
of the node in B+ tree
int cmpNodeBP(const void * a, const void * b);   // Compare the
node by their data


// wstring ⟷ char *, for word stemming
std::wstring chararrToWstring(char * st);
char * wstringToChararr(std::wstring wst);


// Word Stmming wrapper
void WordStem(string term);


// Print the ticks and duration, for -tr or --time function
void PrintTime(clock_t start, clock_t end);
```

```
#endif
```

# 5.3 invIndexFunc.cpp

```cpp
// Implementation of methods related to B+ tree in invIndex.h
#include "invIndexHeader.h"
#include "wordStem/english_stem.h"
#include <algorithm>
#include <codecvt>
#include <filesystem>
#include <locale>
#include <string>
// To avoid unexpected import problems
extern "C" {
    #include <ctype.h>
    #include <stdbool.h>
    #include <stdio.h>
    #include <stdlib.h>
    #include <string.h>
}
// file system namespace
namespace fs = std::filesystem;

string docNames[MAXDOCSUM];          // Array containing names of
documents(global variable)
HashTb H;                            // Hash table storing the stop
words

// The highest-level function, which users can call it directly
// isTest: -t or --test mode, just use one particular file
// containStopWords: -s or --stopwords mode, contain stop words
when building inverted index
BplusTree InvertedIndex(bool isTest, bool containStopWords) {
    char dir[MAXREADSTRLEN];               // Directory name
    char fname[MAXREADSTRLEN];              // File name
    BplusTree InvIndex = CreateBP();     // Inverted index, stored
in B+ tree

    askforFilePos(dir, fname, isTest);   // Ask for the position
```

```c
of file or directory
    // Make a traversal in the directory(or a single file) and
build the inverted index from it
    InvIndex = fileTraversaler(InvIndex, dir, fname, isTest,
containStopWords);
    if (InvIndex→size) {              // If the inverted index
contains the data, it indicates the success of building
        printf("Build successfully!\n");
    } else {                          // Otherwise, it fails
        printf("Fail to build an inverted index!\n");
    }

    return InvIndex;                          // Return the final
inverted index
}

// Ask for the position of the directory or file
// dir: directory name
// fname: file name
// isTest: -t or --test mode, just use one particular file
void askforFilePos(char * dir, char * fname, bool isTest) {
    if (isTest) {    // If we choose the test mode in the test
file,
        printf("Now testing the correctness of inverted Index:
\n");
        printf("Please input the name of the input sample file:
\nName: ");
        scanf("%s", fname);  // then we should input the filename
    } else {          // Otherwise(in the main program or other
default situations),
        printf("Now building an inverted Index:\n");
        printf("Please input the directory of the documents:
\nPath: ");
        scanf("%s", dir);    // we should input the name of the
directory
    }
}

// Make a traversal of all files(or a single file) and build the
```

```cpp
inverted index from them(or it)
// T: B+ tree containing the inverted index
// dir: directory name
// fname: file name
// isTest: -t or --test mode, just use one particular file
// containStopWords: -s or --stopwords mode, contain stop words
when building inverted index
BplusTree fileTraversaler(BplusTree T, char * dir, char * fname,
bool isTest, bool containStopWords) {
    int docCnt = 0;       // Count the number of documents and act
as the index of the documents at the same time
    char * wholePath;     // The whole path name
    FILE * fp = NULL;     // File pointer

    H = GenerateHashTb();                     // Build a hash table
for stop words
    if (containStopWords) {                   // If in stopwords
mode, print the hash table
        PrintHashTb(H);
    }


    wholePath = new char[MAXREADSTRLEN];
    if (!isTest) {        // If we choose the test mode in the
test file,
        fs::path dirPath(dir);
        if (fs::exists(dirPath) && fs::is_directory(dirPath))
{   // Make a traversal in the directory
            for (const auto& entry : fs::directory_iterator(dirPath))
{
                if (fs::is_regular_file(entry)) {    // entry: a
single file
                                    std::string  filename  =
entry.path().filename().string();  // Get the file name
                    docNames[docCnt] = new char[filename.length()
+ 1];
                    strcpy(docNames[docCnt], filename.c_str());   //
Store the filename
                    strcpy(wholePath, (dirPath.string() + "/" +
filename).c_str());  // Get the whole path name
```

```cpp
                        // Open the file
                        fp = fopen(wholePath, "r");
                        if (!fp) {  // Error handler
                            printf("Couldn't open the file!\n");
                            exit(1);
                        }
                        // Update the inverted index
                        T = UpdateInvertedIndex(T, docCnt++, fp,
containStopWords);
                    }
                }
                if (fp)  // Don't forget close the file pointer
                    fclose(fp);
            } else {  // Input wrong directory
                perror("Could not open directory");

            }
        } else {
            strcpy(dir, DEFAULTFILEPOS);       // The file is in the
default position
            std::string sdir(dir);
            std::string sfname(fname);
            strcpy(wholePath, (sdir + "/" + sfname).c_str());  // Get
the whole path name

            docNames[docCnt] = (string)malloc(sizeof(char) *
(strlen(fname) + 1));
            strcpy(docNames[docCnt], fname);                    //
Store the filename
            // Open the file
            fp = fopen(wholePath, "r");
            if (!fp) {  // Error handler
                printf("Couldn't open the file!\n");
                exit(1);
            }
            // Update the inverted index
            T = UpdateInvertedIndex(T, docCnt++, fp,
containStopWords);
            fclose(fp);     // Don't forget close the file pointer
```

```c
    }

    return T;
}

// Update the Inverted Index while reading a new document
// T: B+ tree containing the inverted index
// docCnt: the index of the document
// fp: file pointer
// containStopWords: -s or --stopwords mode, contain stop words
when building inverted index
BplusTree UpdateInvertedIndex(BplusTree T, int docCnt, FILE * fp,
bool containStopWords) {
    int i;
    int pre, cur;                           // Mark the start and
the end of one word
    char tmp[MAXREADSTRLEN];                // Memory space storing
the reading data temporarily
    string term;                            // Term(or word)
    bool isDuplicated;                      // A flag, record
whether the term exists in the B+ tree
    NodeBP nodebp;                          // Node in B+ tree


    while (fgets(tmp, MAXREADSTRLEN - 1, fp) != NULL) {   //
Continue reading the file, until arrive at the end of file
        pre = cur = 0;                      // Initialization
        for (i = 0; i < strlen(tmp); i++) {   // Retrieve all
characters in the tmp string
            if (!isalpha(tmp[i])) {         // Maybe it's time
to record a word
                cur = i;
                if (cur > pre) {            // Legitimate
situation
                    term = (char *)malloc(sizeof(char) * (cur -
pre + 1));
                    strncpy(term, tmp + pre, cur - pre);
                    term[cur - pre] = '\0';    // Don't forget
this step
```

```c
                    // If we consider the stop words(default) and
assure the term is a stop word,
                    if (!containStopWords && FindHashSW(term, H,
true) >= 0) {
                        pre = cur + 1;
                        continue;  // then we should ignore it
                    }

                    // Word stemming
                    WordStem(term);

                    isDuplicated = false;
                nodebp = FindBP(term, docCnt, T, &isDuplicated);  //
Find the appropriate position for the term
                    // If isDuplicated is true, then the time of
the term will +1 in function isSameTerm()
                    if (!isDuplicated) {  // If it's a new term,
insert it!
                        T = InsertBP(term, docCnt, nodebp, T);
                    }
                }

                pre = cur + 1;     // Move the start position for
possible new word
            }
        }

        // Handle the last possible word in the tmp string
        if (!cur || pre > cur && pre != i) {
            cur = i;
            term = (char *)malloc(sizeof(char) * (cur - pre +
1));
            strncpy(term, tmp + pre, cur - pre);
            term[cur - pre] = '\0';

            // Word stemming
            WordStem(term);
```

```c
            isDuplicated = false;
            nodebp = FindBP(term, docCnt, T, &isDuplicated); //
Find the appropriate position for the term
            // If isDuplicated is true, then the time of the term
will +1 in function isSameTerm()
            if (!isDuplicated) {  // If it's a new term, insert
it!
                T = InsertBP(term, docCnt, nodebp, T);
            }
        }
    }

    return T;
}

// Create a B+ tree
BplusTree CreateBP() {
    BplusTree T = (BplusTree)malloc(sizeof(struct nodebp));  //
Allocate the memory space for new B+ tree
    if (T == NULL) {     // Allocation failure
        printf("Failed to create a B+ Tree!\n");
        return T;
    }

    int i;

    // Memory allocation and initialization of data and children
    for (i = 0; i ≤ ORDER; i++) {
        T→data[i] = (Data)malloc(sizeof(struct data));
            T→data[i]→term = (string)malloc(sizeof(char) *
MAXWORDLEN);
        T→data[i]→poslist = CreatePL();
        T→children[i] = (NodeBP)malloc(sizeof(struct nodebp));
    }
    // Initialization of other fields
    T→size = 0;
    T→childrenSize = 0;
    T→parent = NULL;
```

```
    return T;
}


// Find a term in B+ tree
// term: term
// docCnt: the index of the document
// T: inverted index
// flag: true if the term is found, false otherwise
// isSearch: mark the find mode(-f or --find)
NodeBP FindBP(string term, int docCnt, BplusTree T, bool * flag,
bool isSearch) {
    int i;

    if (!T) {    // If the tree is empty, return the tree(actually,
it's impossible in our program)
        return T;
    } else if (!T→childrenSize) {  // If we arrive at the leaf
node, search its data
        isSameTerm(term, docCnt, T, flag, isSearch);
        return T;
    }

    int pos = -1;  // The index of the appopriate non-leaf node
    for (i = 0; i < T→size; i++) {
        if (strcmp(term, T→data[i]→term) < 0) { // Find the
first node which have term with higher lexicographic number
            pos = i;
            break;
        }
    }
    if (pos == -1) {  // If no position found in above loop,
choose the last node
        pos = i;
    }

        return FindBP(term, docCnt, T→children[pos], flag,
isSearch);  // Continue finding in the children node
}
```

```c
// Check if the term exists in the B+ tree
// term: term
// docCnt: the index of the document
// nodebp: the appropriate node where the term may exists or will
exists after insertion
// flag: true if the term is found, false otherwise
// isSearch: mark the find mode(-f or --find)
void isSameTerm(string term, int docCnt, NodeBP nodebp, bool *
flag, bool isSearch) {
    int i;

    if (nodebp→size) {      // If it's not an empty node, start
searching
        for (i = 0; i < nodebp→size; i++) {
            if (!strcmp(term, nodebp→data[i]→term)) {      // If
the term exists in the inverted index
                if (!isSearch) {      // If it's not in the find
mode
                    EnqueuePL(docCnt, nodebp→data[i]→poslist);  //
Update the poslist of the term
                } else {     // Otherwise, print all info of the
term
                    PosList poslist = nodebp→data[i]→poslist;     //
Position list
                    int size = poslist→size;                 //
The number of all documents where the term appears
                    int cnt = 0;                             //
Record the total frequency of the term

                    printf("Successfully find the word!\n");     //
Some banners
                        printf("The word was found in files below:
\n");

                    int j;
                    int ** posArr = (int **)malloc(sizeof(int *)
* size);   // Allocation of a 2D array
                        for (j = 0; j < size; j++) {
                            posArr[i] = (int *)malloc(sizeof(int) *
```

```c
2);
                }

            posArr = RetrievePL(poslist);                    //
Put the poslist in a 2D array

                for (j = 0; j < size; j++) {   // Print the
name of documents and their frequency respectively
                    if (posArr[j][1] <= 1)   // Singular
                    printf("%s: %d time\n", docNames[posArr[j]
[0]], posArr[j][1]);
                    else    // Plural
                    printf("%s: %d times\n", docNames[posArr[j]
[0]], posArr[j][1]);
                    cnt += posArr[j][1];
                }
            printf("Frequency: %d\n", cnt);                   //
The total frequency
                printf("---------------------------------
\n");
            }

            *flag = true;       // mark the flag, indicating
we find the term
            break;
        }
      }
    }
}


// Insert a term into the B+ tree
// term: term
// docCnt: the index of the document
// nodebp: the appropriate node where the term will be inserted
// Tree: B+ tree containing the inverted index
BplusTree InsertBP(string term, int docCnt, NodeBP nodebp,
BplusTree Tree) {
    int i;
```

```c
    strcpy(nodebp→data[nodebp→size]→term, term);
    EnqueuePL(docCnt, nodebp→data[nodebp→size++]→poslist);    //
Add the data info
      qsort(nodebp→data, nodebp→size, sizeof(nodebp→data[0]),
cmpData);    // Sort the data in time

    Tree = SplitBP(nodebp, Tree);     // Split the node
    return Tree;
}


// Split the node when the node is full
// nodebp: the appropriate node where the term will be inserted
// Tree: B+ tree containing the inverted index
BplusTree SplitBP(NodeBP nodebp, BplusTree Tree) {
    if (!nodebp→childrenSize && nodebp→size ≤ ORDER     // If
the node is not full
        || nodebp→childrenSize && nodebp→size < ORDER) {    //
(consider both leaf node and non-leaf node),
        return Tree;                                        // do
nothing!
    }

     // lnodebp, rnodebp: the left and right part of the split
node
    // tmpNodebp: store the node temporarily
    // parent: the parent node of nodebp
    NodeBP lnodebp, rnodebp, tmpNodebp, parent;
    int cut;    // The position of the middle data
    int i, j;

    parent = nodebp→parent;

    if (!parent) {  // If the node has no parent(i.e. this node
is the root),
        tmpNodebp = CreateBP();  // create a new node as the
parent(and also the root of the tree)
        parent = (NodeBP)malloc(sizeof(struct nodebp));
        Tree = parent = tmpNodebp;
    }
```

```
    lnodebp = CreateBP();
    rnodebp = CreateBP();
    lnodebp→parent = rnodebp→parent = parent;   // Connect the
two parts with the parent node


    if (!nodebp→childrenSize) {     // If the node is the leaf
node
        cut = LEAFCUT;

        for (i = 0; i < cut; i++) {     // Assign the data in the
left part of original node to lnodebp
            lnodebp→data[i] = nodebp→data[i];
        }
        lnodebp→size = cut;

        for (j = cut; j < nodebp→size; j++) {    // Assign the
data in the right part of original node to rnodebp
            rnodebp→data[j - cut] = nodebp→data[j];
        }
        rnodebp→size = nodebp→size - cut;

    } else {                         // If the node is the non-leaf
node
        cut = NONLEAFCUT;

        for (i = 0; i ≤ cut; i++) {     // Assign the data and
children in the left part of original node to lnodebp
            if (i ≠ cut)
                lnodebp→data[i] = nodebp→data[i];
            lnodebp→children[i] = nodebp→children[i];
            lnodebp→children[i]→parent = lnodebp;
        }
        lnodebp→size = cut;
        lnodebp→childrenSize = cut + 1;

        // Assign the data and children in the right part of
original node to rnodebp
        for (j = cut + 1; j < nodebp→size; j++) {
```

48

```c
            rnodebp→data[j - cut - 1] = nodebp→data[j];
        }
        for (j = cut + 1; j < nodebp→childrenSize; j++) {
            rnodebp→children[j - cut - 1] = nodebp→children[j];
            rnodebp→children[j - cut - 1]→parent = rnodebp;
        }
        rnodebp→size = nodebp→size - cut - 1;
        rnodebp→childrenSize = nodebp→childrenSize - cut - 1;
    }


    // Assign the middle data in the original node to its parent
    parent→data[parent→size++] = nodebp→data[cut];
    if (parent→childrenSize) {    // If the parent has children(not
be created newly)
        for (i = 0; i < parent→childrenSize; i++) {
            if (parent→children[i] == nodebp) {    // Replace
the original node with lnodebp
                parent→children[i] = lnodebp;
                break;
            }
        }
    } else {    // newly created parent
        parent→children[parent→childrenSize++] = lnodebp;  //
Insert the lnodebp
    }
    parent→children[parent→childrenSize++] = rnodebp;       //
Insert the rnodebp

    // Sort the data and children of the parent
    qsort(parent→data, parent→size, sizeof(parent→data[0]),
cmpData);
    qsort(parent→children, parent→childrenSize, sizeof(parent-
>children[0]), cmpNodeBP);

    free(nodebp); // Free the memory of the original node

    Tree = SplitBP(parent, Tree);  // Continue spliting the upper
node
```

```c
        return Tree;
}


// Print the B+ tree(level-order traversal)
// T: B+ tree containing the inverted index
void PrintBPTree(BplusTree T) {
    int i;
    NodeBP nodebp;        // The node obtained from the queue
    QueueBP q;            // The queue containing the nodes from
B+ tree

    printf("B+ Tree of Inverted Index:\n");

    q = CreateQueueBP();  // Create an empty queue
    EnqueueBP(T, q);        // Put the root of the tree into the
queue first
    EnqueueBP(NULL, q);    // Put the NULL pointer, for creation
of newline

    while (q→size) {      // If the queue isn't empty, repeat the
following steps
        nodebp = DequeueBP(q);  // Get the front node
        if (!nodebp) {          // If it's an NULL pointer, it's
time to add a newline
            printf("\n");
            if (q→size) {       // If the queue isn't empty,
continue add a new NULL pointer
                EnqueueBP(NULL, q);
            }
        } else {
            printf("[");        // Print the node's data(just the
term)
            for (i = 0; i < nodebp→size; i++) {
                if (!i) {
                    printf("%s", nodebp→data[i]→term);
                } else {
                    printf(", %s", nodebp→data[i]→term);
                }
            }
```

```c
                printf("]");
            }

            if (nodebp) {           // If nodebp isn't a NULL pointer,
then put its children into the queue
                for (i = 0; i < nodebp→childrenSize; i++) {
                    EnqueueBP(nodebp→children[i], q);
                }
            }
        }
    }
}


// Create the queue
QueueBP CreateQueueBP() {
    QueueBP Q = (QueueBP)malloc(sizeof(struct queuebp));
    Q→size = 0;
    Q→front = Q→rear = 0;

    return Q;
}


// Put the node of B+ tree into the queue
// nodebp: the newly added node
// Q: the queue
void EnqueueBP(NodeBP nodebp, QueueBP Q) {
    if (Q→size ≥ SIZE) {  // If the queue is full, enqueue
operation fails
        printf("Full B+-tree-item queue!\n");
        exit(1);
    }
    Q→data[Q→rear++] = nodebp;  // Add new node
    Q→size++;
}


// Get the front node and delete it from the queue
// Q: the queue
NodeBP DequeueBP(QueueBP Q) {
    if (!Q→size) {              // If the queue is empty, dequeue
operation fails
```

```c
        printf("Empty B+-tree-item queue!\n");
        exit(1);
    }
    NodeBP returnNodeBP = Q→data[Q→front++];  // Get the front
node
    Q→size--;      // Delete the node from queue
    return returnNodeBP;
}

// Create the poslist
PosList CreatePL() {
    PosList L;

    L = (PosList)malloc(sizeof(struct poslist));
    L→size = 0;
    L→front = (PosData)malloc(sizeof(struct posdata));
    L→rear = (PosData)malloc(sizeof(struct posdata));
    L→front = L→rear;
    L→rear→pos = -1;   // Distinguish from other nodes

    return L;
}

// Add new position
// pos: the position
// L: the position list
void EnqueuePL(int pos, PosList L) {
    if (L→rear→pos ≠ pos) {  // If it's a new position
        PosData tmp = (PosData)malloc(sizeof(struct posdata));
        if (!tmp) {
            printf("Fail to create a new position data!\n");
            exit(1);
        }  // Insert the new one in the position list
        tmp→pos = pos;
        tmp→time = 1;
        tmp→next = L→rear→next;
        L→rear→next = tmp;
        L→rear = tmp;
        L→size++;
```

```c
    } else {      // Otherwise, just increment the frequency
        L→rear→time++;
    }


}


// Retrieve all position in the list
// L: the position list
int ** RetrievePL(PosList L) {
    if (!L→size) {        // If the list is empty, retrieve
operation fails
        printf("Empty position-data queue!\n");
        exit(1);
    }

    int i = 0, j;
    int ** posArr = (int **)malloc(sizeof(int *) * L→size);
    for (j = 0; j < L→size; j++) {  // Memory Allocation for 2D
array
        posArr[j] = (int *)malloc(sizeof(int) * 2);
    }

    PosData cur = L→front→next;

    while (cur ≠ NULL) {  // Make a traversal in the position
list
        posArr[i][0] = cur→pos;  // Get the specific info of the
position
        posArr[i][1] = cur→time;
        cur = cur→next;
        i++;
    }

    return posArr;
}

// Build a hash table
HashTb GenerateHashTb() {
    int i;
```

```c
    int pre, cur;                          // Mark the start and the end
of one word
    HashTb H;                              // The hash table containing
the stop words
    FILE * fp;                             // File pointer
    char fname[MAXWORDLEN];         // File name
     char tmp[MAXREADSTRLEN];          // Memory space storing the
reading data temporarily
    char * term;                           // Term(or word)

    H = InitHashTb();                 // Initialization

        strcpy(fname,  STOPWORDPATH);       //  Use   default
path(stop_words.txt)
    fp = fopen(fname, "r");        // Open the file
    if (!fp) {
        printf("Fail to open the file of stopwords!\n");
        exit(1);
    }

      while (fgets(tmp, MAXREADSTRLEN - 1, fp) ≠ NULL) {   //
Continue reading the file, until arrive at the end of file
        pre = cur = 0;                          // Initialization
         for (i = 0; i < strlen(tmp); i++) {       // Retrieve
all characters in the tmp string
            if (!isalpha(tmp[i])) {              // Maybe it's
time to record a word
                cur = i;
                if (cur > pre) {                 // Legitimate
situation
                    term = (char *)malloc(sizeof(char) * (cur -
pre + 1));
                    strncpy(term, tmp + pre, cur - pre);
                    term[cur - pre] = '\0';
                    InsertHashSW(term, H);        // Insert the
new term
                }

                pre = cur + 1;
```

```c
            }
        }

        // Handle the last possible word in the tmp string
        if (!cur || pre > cur && pre != i) {
            cur = i;
            term = (char *)malloc(sizeof(char) * (cur - pre + 1));
            strncpy(term, tmp + pre, cur - pre);
            term[cur - pre] = '\0';
            InsertHashSW(term, H);              // Insert the new term
        }
    }
    fclose(fp);
    return H;
}

// Initialization of the hash table
HashTb InitHashTb() {
    HashTb H;                // Hash table
    int i;

    H = (HashTb)malloc(sizeof(struct hashtb));    // Memory allocation for the whole table
    if (H == NULL) {
        printf("Fail to create a hash table for stopwords!\n");
        exit(1);
    }

    H->size = STOPWORDSUM;      // maxixum size

    for (i = 0; i < H->size; i++) {
        H->data[i] = (HashSW)malloc(sizeof(hashsw));
        if (H->data[i] == NULL) {      // Memory allocation for cells
            printf("Fail to create a hash table for stopwords!\n");
            exit(1);
```

```c
        }
            H→data[i]→stopword = (string)malloc(sizeof(char) *
MAXWORDLEN);
        H→data[i]→info = Empty;
    }

    return H;
}


// Find the stopwords or other words in the hash table
// stopword: stop word
// H: hash table containing the stop words
// justSearch: find the term without subsequent insertion
int FindHashSW(string stopword, HashTb H, bool justSearch) {
    int pos;                                // Appropraite position
    int collisionNum = 0;                    // collision number,
for quadratic probe
    pos = HashFunc(stopword, H→size);    // Use hashing function
first

    // If we just search a term in the hash table and assure it's
not a stop word, then return
    if (justSearch && (H→data[pos]→info == Empty || strcmp(H-
>data[pos]→stopword, stopword))) {
        return -1;
    }

    // Collision occurs!
    while (H→data[pos]→info ≠ Empty && strcmp(H→data[pos]-
>stopword, stopword)) {
        pos += 2 * ++collisionNum - 1;  // Quadratic probe
        if (pos ≥ H→size)
            pos -= H→size;
    }
    return pos;
}


// Insert a new stopword in hash table
// stopword: stop word
```

```c
// H: hash table containing the stop words
void InsertHashSW(string stopword, HashTb H) {
    int pos;
    pos = FindHashSW(stopword, H, false);    // Find the correct
position
    if (H→data[pos]→info ≠ Legitimate)    // Insert the stop
word
    {
        H→data[pos]→info = Legitimate;
        strcpy(H→data[pos]→stopword, stopword);
    }
}


// Hashing function
// stopword: stop word
// size: the maximum size of the hash table
int HashFunc(string stopword, int size) {
    unsigned int val = 0;
    while (*stopword ≠ '\0')
        val = (val << 5) + *stopword++;    // Generate the hash
value from every character in the string
    return val % size;
}


// Print hash table
void PrintHashTb(HashTb H) {
    int i;

    printf("Stopwords in hash table:\n");
    for (i = 0; i < H→size; i++) {
        if (H→data[i]→info ≠ Empty) {
            printf("%d: %s\n", i, H→data[i]→stopword);
        }
    }
    printf("\n");
}


// Comparison functions used in qsort()
int cmpData(const void * a, const void * b) {
```

```cpp
    const Data dataA = *(const Data*)a;
    const Data dataB = *(const Data*)b;

    return strcmp(dataA→term, dataB→term);
}

int cmpNodeBP(const void * a, const void * b) {
    const NodeBP nodebpA = *(const NodeBP*)a;
    const NodeBP nodebpB = *(const NodeBP*)b;

        return  strcmp(nodebpA→data[0]→term,  nodebpB→data[0]-
>term);
}

// wstring ⟷ char *, for word stemming
std::wstring chararrToWstring(char * st) {
    std::string tmp(st);
    std::wstring_convert<std::codecvt_utf8<wchar_t>> converter;
    std::wstring wstr = converter.from_bytes(tmp);

    return wstr;
}

char * wstringToChararr(std::wstring wst) {
    std::wstring_convert<std::codecvt_utf8<wchar_t>> converter;
    std::string tmp = converter.to_bytes(wst);
    char * st = new char[tmp.size() + 1];
    strcpy(st, tmp.c_str());

    return st;
}

// Word Stmming wrapper
void WordStem(string term) {
    std::wstring term_wstr;                       // the wstring form of
the term
        stemming::english_stem<> StemEnglish;   // Word stemming
function(a little clumsy)
```

```cpp
    term_wstr = chararrToWstring(term);
            transform(term_wstr.begin(),   term_wstr.end(),
term_wstr.begin(), ::tolower);
    StemEnglish(term_wstr);
    term = wstringToChararr(term_wstr);
}


// Print the ticks and duration, for -tr or --time function
void PrintTime(clock_t start, clock_t end) {
    clock_t tick;        // ticks
    double duration;     // duration(unit: seconds)

    tick = end - start;
    duration = ((double)(tick)) / CLOCKS_PER_SEC;
    printf("Ticks: %lu\n", (long)tick);     // Print the info
    printf("Duration: %.2fs\n", duration);
}
```

## 5.4 invIndexTest.cpp

```cpp
#include "invIndexHeader.h"
#include "wordStem/english_stem.h"
#include <stdbool.h>
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <string>
#include <time.h>


int main(int argc, char * argv[]) {          // Use command line
parameters
    int i;
    bool isTest = false;                      // Whether open test
mode(-t or --test)
    bool Print = false;                       // Whether open print
mode(-p or --print)
    bool isFound;                             // Whether the term is
found in the inverted index
    bool containStopWords = false;            // Whether open stopword
```

```cpp
mode(-s or --stopwords)
    bool timeRecord = false;                    // Whether open time
record mode(-tr or --time)
    int findCnt = 0;                            // The time of finding
a single word for find mode(-f=n or --find=n)
    char * pos;                                 // Get the number after
`=` in parameters
    char tmp[MAXWORDLEN];                       // Store the input
string temporarily
    char * word;                                // The word to be
searched
    double duration;                            // Duration of running
a function
    std::wstring word_wstr;                     // wstring form of
word
      stemming::english_stem<> StemEnglish;   // Word stemming
function
    BplusTree InvIndex;                         // Inverted Index
    clock_t start, end, tick;                   // Record the start
and the end of the clock

    for (i = 1; i < argc; i++) {                // Read the parameters
        if (!strcmp(argv[i], "--test") || !strcmp(argv[i], "-t"))
{   // Test mode
            isTest = true;
        } else if (strstr(argv[i], "--find") || strstr(argv[i],
"-f")) {   // Find mode
            if ((pos = strchr(argv[i], '='))) {
                findCnt = atoi(pos + 1);    // Get the number
behind `=`
                if (!findCnt) {
                    printf("Wrong Number!\n");
                    exit(1);
                }
            } else {
                findCnt = 1;    // Use default number
            }
        } else if (!strcmp(argv[i], "--print") || !strcmp(argv[i],
"-p")) {   // Print mode
```

```c
            Print = true;
             } else if (!strcmp(argv[i], "--stopwords") || !
strcmp(argv[i], "-s")) {   // Stopword mode
            containStopWords = true;
        } else if (!strcmp(argv[i], "--time") || !strcmp(argv[i],
"-tr")) {   // Time record mode
            timeRecord = true;
        } else {    // Error
            printf("Wrong Parameter!\n");
            exit(1);
        }
    }

    if (!timeRecord) {    // No time record
        InvIndex = InvertedIndex(isTest, containStopWords);
    } else {   // Time record
        start = clock();
        InvIndex = InvertedIndex(isTest, containStopWords);
        end = clock();
        PrintTime(start, end);
    }

    if (Print) {   // Print the B+ tree
        PrintBPTree(InvIndex);
    }

    if (InvIndex→size && findCnt) {   // Search the single word
        word = (char *)malloc(sizeof(char) * MAXWORDLEN);
         printf("\nFinding Words Mode(only supports single word
finding):\n");
        for (i = 0; i < findCnt; i++) {   // For every search
            isFound = false;
            printf("Find %d: ", i + 1);
            scanf("%s", tmp);
            strcpy(word, tmp);

            // Word stemming
            word_wstr = chararrToWstring(word);
            StemEnglish(word_wstr);
```

```
                word = wstringToChararr(word_wstr);

                // Find the word
                if (!timeRecord) {   // No time record
                    FindBP(word, -1, InvIndex, &isFound, findCnt);
                } else {   // Time record
                    start = clock();
                    FindBP(word, -1, InvIndex, &isFound, findCnt);
                    end = clock();
                    PrintTime(start, end);
                }

                // If not found, then give relevant information
                if (!isFound) {
                    printf("Sorry, no such word in the inverted index!
\n");

                    printf("--------------------------------\n");
                }
            }
        }

    return 0;
}
```

# References

[1] Blake-Madden, OleanderStemmingLibrary, https://github.com/Blake-Madden/OleanderStemmingLibrary

# Declaration

*We hereby declare that all the work done in this project titled "Roll Your Own Mini Search Engine" is of our independent effort as a group.*