

# Algorithms and Data Structures

## Occurrences Counter

### Students:

Isadora Loredó – 91322-50%

Margarida Martins – 93169-50%

31<sup>th</sup> of December 2019

# Index

The Assignment	2
Implementation	3
<b>manipulate_file.h</b>	3
hash_size	3
file_data_t	3
word_stats	3
open_text_file()	4
close_text_file()	4
read_word()	5
hash_function()	5
hash_data	6
new_hash_data()	6
hash_resize()	6
hash_data_bt	7
new_hash_data_bt()	7
hash_resize_bt()	8
traverse_tree()	8
insert_bt()	9
<b>Linked List</b>	10
<b>Binary Tree</b>	11
Word statistics of text file	12
<b>Number of occurrences</b>	12
<b>First and last appearances</b>	13
<b>Smallest, medium and largest distance</b>	14
<b>Time of execution and resize</b>	15
Appendix	17
<b>Occurrences_tree.c</b>	17
<b>Occurrences.c</b>	19
manipulate_file.h	21

# The Assignment

Count the number of occurrences of each word of a text file (book Sherlock Holmes – A Study In Scarlet), records the location of the first and last occurrences of each distinct word and records the smallest, largest, and average distances between consecutive occurrences of the same distinct word.

The program must use a **hash table** (separate chaining), the hash table size should grow dynamically, and each hash table entry should point to either a **linked list** or an ordered **binary tree**.

# Implementation

## manipulate\_file.h

The header file 'manipulate\_file.h' contains functions declaration and data structures that are used by both 'c' programs, with **linked list** and with **binary tree**, through the directive **#include**.

### hash\_size

Global variable used to define the size of the **hash table**.

```
unsigned int hash_size= 2000u;
```

### file\_data\_t

The new data type struct 'file\_data\_t' defines the stats of the text file read and stores the current word.

- **long** word\_pos: position of the current word being read in the file;
- **long** word\_num: number of words read in the file;
- **char** word[]: array of characters meaning the current word being read in the file;
- **FILE \*fp**: file pointer pointing to the file being accessed;
- **long** current\_pos: current position in the file for the read word;

```
// FILE STRUCT
typedef struct file_data_t {
    // public data
    long word_pos; // zero-based
    long word_num; // zero-based
    char word[64];
    // private data
    FILE *fp;
    long current_pos; // zero-based
} file_data_t;
```

### word\_stats

The new data type struct 'word\_stats' defines the stats of each word read by the text file.

- **int** number\_occurrences: number of occurrences of the word in the file;
- **int** first\_appearance: first appearance of the word in the file;
- **int** last\_appearance: last appearance of the word in the file;
- **int** s\_distance: smallest distance between consecutive occurrences of the same distinct word;
- **int** l\_distance: largest distance between consecutive occurrences of the same distinct word;
- **int** m\_distance: medium distance between consecutive occurrences of the same distinct word, so the average distance can be calculated;

- `int t_distance`: total distance between consecutive occurrences of the same distinct word;

```
// STRUCT THAT DEFINES THE STATS FOR THE WORDS IN THE HASH-TABLE
typedef struct word_stats {
    int number_occurrences;
    int first_appearance;
    int last_appearance;
    int s_distance; // smallest distance
    int l_distance; // largest distance
    int m_distance; // medium distance
    int t_distance; // total distance
} word_stats;
```

## open\_text\_file()

The function 'open\_text\_file' is used to open the file to be read and initialize the variables of the data type 'file\_data\_t'.

Arguments:

- `char *file_name`: file pointer pointing to the text file to be accessed;
- `file_data_t *fd`: file stats pointer pointing to the struct to be initialized;

Return:

- `int value`: returns a integer value, '0' if success or '-1' if file stats pointer is null;

```
// OPEN FILE
int open_text_file(char *file_name, file_data_t *fd) {
    fd->fp = fopen(file_name, "r");
    if (fd->fp == NULL)
        return -1;
    fd->word_pos = -1;
    fd->word_num = -1;
    fd->word[0] = '\0';
    fd->current_pos = -1;
    return 0;
}
```

## close\_text\_file()

The function 'close\_text\_file' is used to close the file read.

Arguments:

- `file_data_t *fd`: file stats pointer pointing to the struct initialized;

```
// CLOSE FILE
void close_text_file(file_data_t *fd) {
    fclose(fd->fp);
    fd->fp = NULL;
}
```

## read\_word()

The function 'read\_word' is called to read the next word of the accessed file and save update its stats. It reads each character of the file using white spaces, special characters and punctuations to identify the beginning and end of a word. The characters of the read words are passed to lower case so same words are not detected as different.

Arguments:

- `file_data_t *fd`: file stats pointer pointing to the struct of the file being accessed;

Return:

- `int value`: returns a integer value, '-1' if end of file reached or 0' if end of file still not reached;

```
// READ WORD
int read_word(file_data_t *fd) {
    int i,c;
    // skip white spaces
    do {
        c = fgetc(fd->fp);
        if(c == EOF)
            return -1;
        fd->current_pos++;
    } while(!((c >= 48 && c<58) || (c>=65 && c<=90) || (c>=97 && c<=122) || (c>=192)));
    // record word
    fd->word_pos = fd->current_pos;
    fd->word_num++;
    fd->word[0] = (char)c;
    for(i = 1; i < (int)sizeof(fd->word) - 1; i++) {
        c = fgetc(fd->fp);
        if(c == EOF)
            break; // end of file
        fd->current_pos++;
        if(!((c >= 48 && c<58) || (c>=65 && c<=90) || (c>=97 && c<=122) || (c>=192)))
            break; // terminate word
        fd->word[i] = (char)c;
    }
    fd->word[i] = '\0';
    for(int j = 0; fd->word[j]; j++){
        fd->word[j] = tolower(fd->word[j]);
    }
    return 0;
}
```

## hash\_function()

The function 'hash\_function' maps each possible key (word of the file) to an integer. The integer will then be used to access the **hash table** array, as the index to the array.

Arguments:

- `const char *str`: the key (word) to be mapped;
- `unsigned int s`: size of the hash table;

Return:

- `unsigned int value`: calculated value as the index to the hash table;

```
// HASH-FUNCTION
unsigned int hash_function(const char *str, unsigned int s) {
    unsigned int h;
    for(h = 0; *str != '\0'; str++)
        h = 157u * h + (0xFFu & (unsigned int)*str); // arithmetic overflow may occur here (just ignore it!)
    return h % s; // due to the unsigned int data type, it is guaranteed that 0 <= h % s < s
}
```

## hash\_data

The new data type struct 'hash\_data' defines the nodes of data to be inserted in the **hash table** to the **linked list** implementation.

- `struct hash_data *next`: pointer to the next node in the linked list;
- `char key[]`: word read in the file;
- `struct word_stats word`: stats of the word (key);

```
// HASH-TABLE NODE LINKED LIST
typedef struct hash_data {
    struct hash_data *next;
    char key[64];
    struct word_stats word;
} hash_data;
```

## new\_hash\_data()

The function 'new\_hash\_data' initializes a new hash\_data type struct, new node, allocating the necessary memory size.

Return:

- `hash_data *hd`: pointer to the start location to the initialized node;

```
// ALLOCATE NEW HASH-DATA
hash_data *new_hash_data(void) {
    hash_data *hd = (hash_data *)malloc(sizeof(hash_data));
    if(hd == NULL) {
        fprintf(stderr, "Out of memory\n");
        exit(1);
    }
    return hd;
}
```

## hash\_resize()

The function 'hash\_resize' is called when the **hash table** reaches a certain number of nodes to be inserted and needs to grow in order to avoid an unwanted too large size of the **linked lists**. It allocates the necessary memory to the hash table with its new size then initializes its null nodes. Finally goes through the old hash table and its linked lists and recalculates the index of each node's key to relocate them in the new resized hash table.

Arguments:

- `hash_data **hash_table`: pointer pointing to the address that points to the beginning of the **hash table**;
- `unsigned int inc`: value of the increment to be added to the hash size;

Return:

- `hash_data **hash_data`: pointer pointing to the address that points to the beginning of the **hash table** with the increased size;

```
//RESIZE LINKED LIST
struct hash_data ** hash_resize(struct hash_data **hash_table, unsigned int inc){
    //printf("Resizing.....\n");
    hash_data *next;
    hash_size+=inc;
    struct hash_data **hash_table_new= malloc((hash_size)*sizeof(struct hash_data));
    int new_idx;
    for(int m = 0;m < hash_size;m++) hash_table_new[m] = NULL;
    for(int l=0;l<hash_size-inc;l++){
        while(hash_table[l]!=NULL){
            new_idx= hash_function(hash_table[l]->key, hash_size);
            next=hash_table[l]->next;
            hash_table[l]->next=hash_table_new[new_idx];
            hash_table_new[new_idx]=hash_table[l];
            hash_table[l]=next;
        }
    }
    free(hash_table);
    return hash_table_new;
}
```

## hash\_data\_bt

The new data type struct 'hash\_data\_bt' defines the nodes of data to be inserted in the **hash table** to the **binary tree** implementation.

- `struct hash_data *left`: pointer to the node on the left side (the left branch);
- `struct hash_data *right`: pointer to the node on the right side (the right branch);
- `char key[64]`: word read in the file;
- `struct word_stats word`: stats of the word (key);

```
// HASH-TABLE NODE BINARY TREE
typedef struct hash_data_bt {
    struct hash_data_bt *left;
    struct hash_data_bt *right;
    char key[64];
    struct word_stats word;
} hash_data_bt;
```

## new\_hash\_data\_bt()

The function 'new\_hash\_data\_bt' initializes a new hash\_data\_bt type struct, new node, allocating the necessary memory size.

Return:

- `hash_data *hd`: pointer to the start location to the initialized node;



```
// ALLOCATE NEW HASH-DATA
hash_data_bt *new_hash_data_bt(void) {
    hash_data_bt *hd = (hash_data_bt *)malloc(sizeof(hash_data_bt));
    if(hd == NULL) {
        fprintf(stderr, "Out of memory\n");
        exit(1);
    }
    return hd;
}
```

## hash\_resize\_bt()

The function 'hash\_resize\_bt' is called when the **hash table** reaches a certain number of nodes to be inserted and needs to grow in order to avoid an unwanted too large size of the **binary trees**. It allocates the necessary memory to the **hash table** with its new size then initializes its null nodes. Finally goes through the old **hash table** and its **binary trees** and recalculates the index of each node's key to relocate them in the new resized hash table.

Arguments:

- **hash\_data\_bt \*\*hash\_table**: pointer pointing to the address that points to the beginning of the **hash table**;
- **unsigned int inc**: value of the increment to be added to the hash size;

Return:

- **hash\_data\_bt \*\*hash\_data**: pointer pointing to the address that points to the beginning of the **hash table** with the increased size;

```
//RESIZE BINARY TREE
struct hash_data_bt ** hash_resize_bt(struct hash_data_bt **hash_table, unsigned int inc){
    //printf("Resizing.....\n");
    hash_data_bt *next;
    struct hash_data_bt **hash_table_new= malloc((hash_size+inc)*sizeof(struct hash_data_bt));
    int new_idx;
    hash_size+=inc;
    for(int m = 0; m < hash_size;m++) hash_table_new[m] = NULL;
    for(int l=0;l<hash_size-inc;l++){
        traverse_tree(hash_table[l],hash_table_new);
    }
    free(hash_table);
    return hash_table_new;
}

// CLOSE FILE
void close_text_file(file_data_t *fd) {
    fclose(fd->fp);
    fd->fp = NULL;
}
```

## traverse\_tree()

The function 'traverse\_tree' is called when is necessary to perform visits in the **binary trees**. The visits are made using depth-first search recursively. It performs the traverse search, then it calculates the new index of the given node in the resized hash table and if the calculated index has still no allocated node it allocates it as the root of a new tree, otherwise it calls the function to insert the node in the tree.

Arguments:

- `hash_data_bt *hb`: pointer pointing to the node to be inserted in the resized **hash table**;
- `hash_data_bt **hash_table`: pointer pointing to the address that points to the beginning of the resized **hash table**;

```
// TRAVERSE_TREE
void traverse_tree(struct hash_data_bt *hb, struct hash_data_bt **hash_table){
    if(hb!=NULL){
        int new_idx= hash_function(hb->key, hash_size);
        if(hash_table[new_idx]==NULL){
            hash_table[new_idx]=new_hash_data_bt();
            strcpy(hash_table[new_idx]->key,hb->key);
            hash_table[new_idx]->left=NULL;
            hash_table[new_idx]->right=NULL;
            hash_table[new_idx]->word=hb->word;
        }
        else{
            hash_data_bt *temp =new_hash_data_bt();
            temp->left=NULL;
            temp->right=NULL;
            temp->word=hb->word;
            strcpy(temp->key,hb->key);
            insert_bt(hash_table[new_idx],temp);
        }
        traverse_tree(hb->left, hash_table);
        traverse_tree(hb->right, hash_table);
    }
}
```

## insert\_bt()

The function 'insert\_bt' is called when is necessary to insert a new node in the **binary tree** at the appropriate location, and it is done recursively. It receives the root node and the node to be inserted as arguments, it compares the keys (word) of the arguments so it can insert the node or keep going forward until it finds the proper leaf to insert it as it's child.

Arguments:

- `hash_data_bt *root`: pointer pointing to the node to be inserted in the resized **hash table**;
- `hash_data_bt *hd`: pointer pointing to the address that points to the beginning of the resized **hash table**;

```
// INSERT
void insert_bt(hash_data_bt *root, hash_data_bt *hd ){
    if(strcmp(hd->key,root->key)<0){
        if(root->left!=NULL)
            insert_bt(root->left,hd);
        else
            root->left=hd;
    }
    if(strcmp(hd->key,root->key)>0){
        if(root->right!=NULL)
            insert_bt(root->right,hd);
        else
            root->right=hd;
    }
}
```

## Linked List

The first solution implemented uses a hash table that grows dynamically and each hash table entry points to a **linked list** when a collision occurs.

The structure [\*\\*hd\*](#) is the node to be stored to the hash table, before it does its index is calculated using a [hash function](#). The word's stats are stored, location of the first and last occurrences and the smallest, largest and average distances between consecutive occurrences. It also counts the number of occurrences of each distinct word in the text file and the number of times the hash table is resized.

```
int main(int argc, char **argv){
    file_data_t fd;
    int idx;
    int j = 0;
    int count = 0;
    int i = 0;
    int r_count = 0;
    struct hash_data **hash_table = malloc(hash_size * sizeof(struct hash_data));
    for (int k = 0; k < hash_size; k++){
        hash_table[k] = NULL;
    }
    open_text_file("SherlockHolmes.txt", &fd);
    while (read_word(&fd) != -1){
        idx = hash_function(fd.word, hash_size);
        hash_data *hd;
        for (hd = hash_table[idx]; hd != NULL && strcmp(fd.word, hd->key) != 0; hd = hd->next)
            i++;
        if (hd == NULL){
            if (i == 0){
                hd = new_hash_data();
                strcpy(hd->key, fd.word);
                hd->word.number_occurrences = 1;
                hd->word.first_appearance = count;
                hd->word.last_appearance = count;
                hd->next = NULL;
                hash_table[idx] = hd;
            } else{
                hd = new_hash_data();
                hd->next = hash_table[idx];
                strcpy(hd->key, fd.word);
                hd->word.number_occurrences = 1;
                hd->word.first_appearance = count;
                hd->word.last_appearance = count;
                hash_table[idx] = hd;
            }
        }
        if (j > hash_size * 5){
            hash_table = hash_resize(hash_table, 1000u);
            r_count++;
        }
        j++;
    } else{
        hd->word.number_occurrences += 1;
        if (hd->word.number_occurrences == 2){
            hd->word.s_distance = count - hd->word.first_appearance;
            hd->word.l_distance = count - hd->word.first_appearance;
            hd->word.t_distance += hd->word.s_distance;
        } else{
            if (hd->word.s_distance > count - hd->word.last_appearance)
                hd->word.s_distance = count - hd->word.last_appearance;
            if (hd->word.l_distance < count - hd->word.last_appearance)
                hd->word.l_distance = count - hd->word.last_appearance;
            hd->word.t_distance += count - hd->word.last_appearance;
        }
        hd->word.m_distance = hd->word.t_distance / hd->word.number_occurrences;
        hd->word.last_appearance = count;
    }
    i = 0;
    count++;
}
close_text_file(&fd);
free(hash_table);
return 1;
```

## Binary Tree

The second solution implemented uses a **hash table** that grows dynamically and each **hash table** entry points to an ordered **binary tree** when a collision occurs.

The structure `*hd` is the node to be stored to the hash table, before it does its index is calculated using a [hash function](#). The word's stats are stored, location of the first and last occurrences and the smallest, largest and average distances between consecutive occurrences. It also counts the number of occurrences of each distinct word in the text file and the number of times the hash table is resized.

```
int main(int argc, char **argv){
    file_data_t fd;
    int idx, j = 0, i = 0, count = 0, r_count = 0;
    struct hash_data_bt **hash_table = malloc(hash_size * sizeof(struct hash_data));
    for (int k = 0; k < hash_size; k++) hash_table[k] = NULL;
    open_text_file("SherlockHolmes.txt", &fd);
    while (read_word(&fd) != -1){
        idx = hash_function(fd.word, hash_size);
        hash_data_bt *hd = hash_table[idx];
        while (hd != NULL && strcmp(fd.word, hd->key) != 0){
            if (strcmp(fd.word, hd->key) < 0){
                hd = hd->left;
            } else{
                hd = hd->right;
            }
        }
        i++;
    }
    if (hd == NULL){
        if (i == 0){
            hd = new_hash_data_bt();
            strcpy(hd->key, fd.word);
            hd->word.number_occurrences = 1;
            hd->word.first_appearance = count;
            hd->word.last_appearance = count;
            hd->right = NULL;
            hd->left = NULL;
            hash_table[idx] = hd;
        } else{
            hd = new_hash_data_bt();
            hd->word.number_occurrences = 1;
            hd->word.first_appearance = count;
            hd->word.last_appearance = count;
            hd->right = NULL;
            hd->left = NULL;
            strcpy(hd->key, fd.word);
            insert_bt(hash_table[idx], hd);
        }
        if (j > hash_size * 5){
            hash_table = hash_resize_bt(hash_table, 10000);
            r_count++;
        }
        j++;
    } else{
        hd->word.number_occurrences += 1;
        if (hd->word.number_occurrences == 2){
            hd->word.s_distance = count - hd->word.first_appearance;
            hd->word.l_distance = count - hd->word.first_appearance;
            hd->word.t_distance += hd->word.s_distance;
        } else{
            if (hd->word.s_distance > count - hd->word.last_appearance)
                hd->word.s_distance = count - hd->word.last_appearance;
            if (hd->word.l_distance < count - hd->word.last_appearance)
                hd->word.l_distance = count - hd->word.last_appearance;
            hd->word.t_distance += count - hd->word.last_appearance;
        }
        hd->word.m_distance = hd->word.t_distance / hd->word.number_occurrences;
        hd->word.last_appearance = count;
    }
    i = 0;
    count++;
}
close_text_file(&fd);
free(hash_table);
return 1;
```

# Word statistics of text file

In order analyze the statistics of the words present in the book “Sherlock Holmes – A Study In Scarlet” the words stats present in the **hash table** (outcome of running the “.c” files) were stored in multiples “.txt” files.

## Number of occurrences

According to the results obtained the book has 19192 different words in a total of 668874.

The word with the biggest number of appearances is the word “the” showing up 36233 times throughout the text (5,4% of all word occurrences).

On the graph below (figure 1) we can observe that the number of words decreases very fast with the number of occurrences. There are 6729 words who only appear once 2799 who appear twice and 1655 thrice, together they make up more than half of the words present in the text (58,3%).

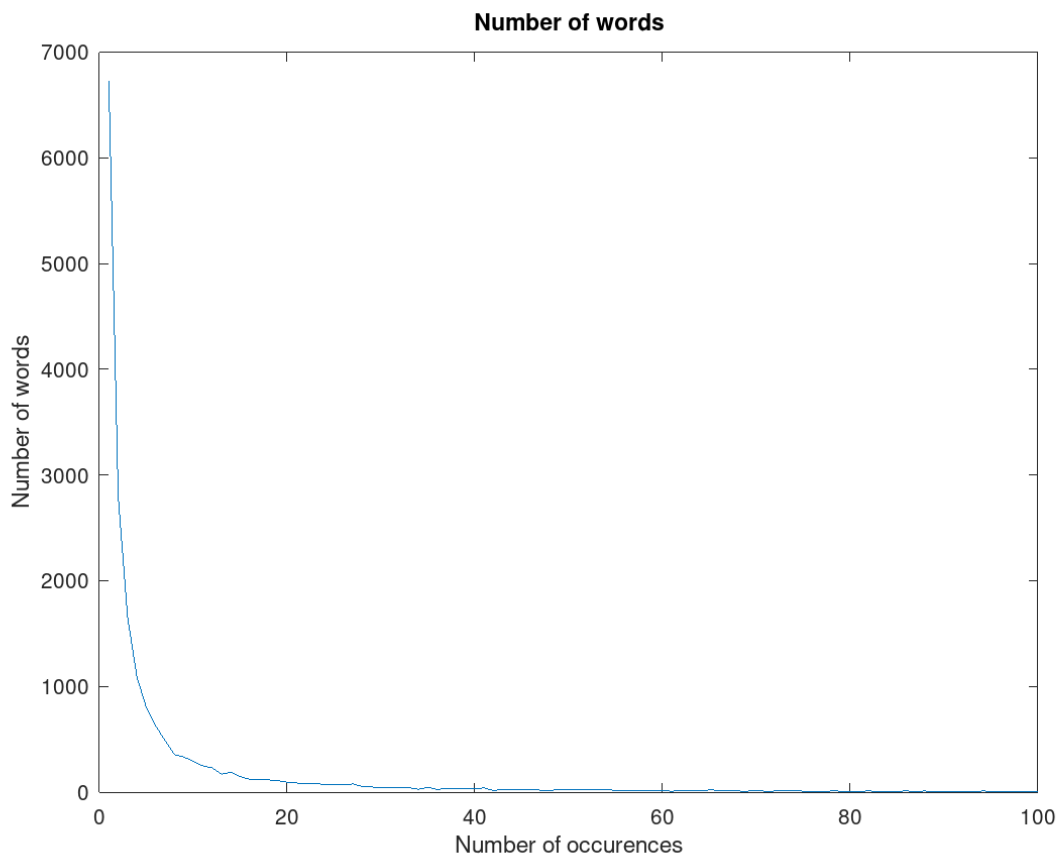


Figure 1- Number of words who occur between 1 and 100 times throughout the text.

## First and last appearances

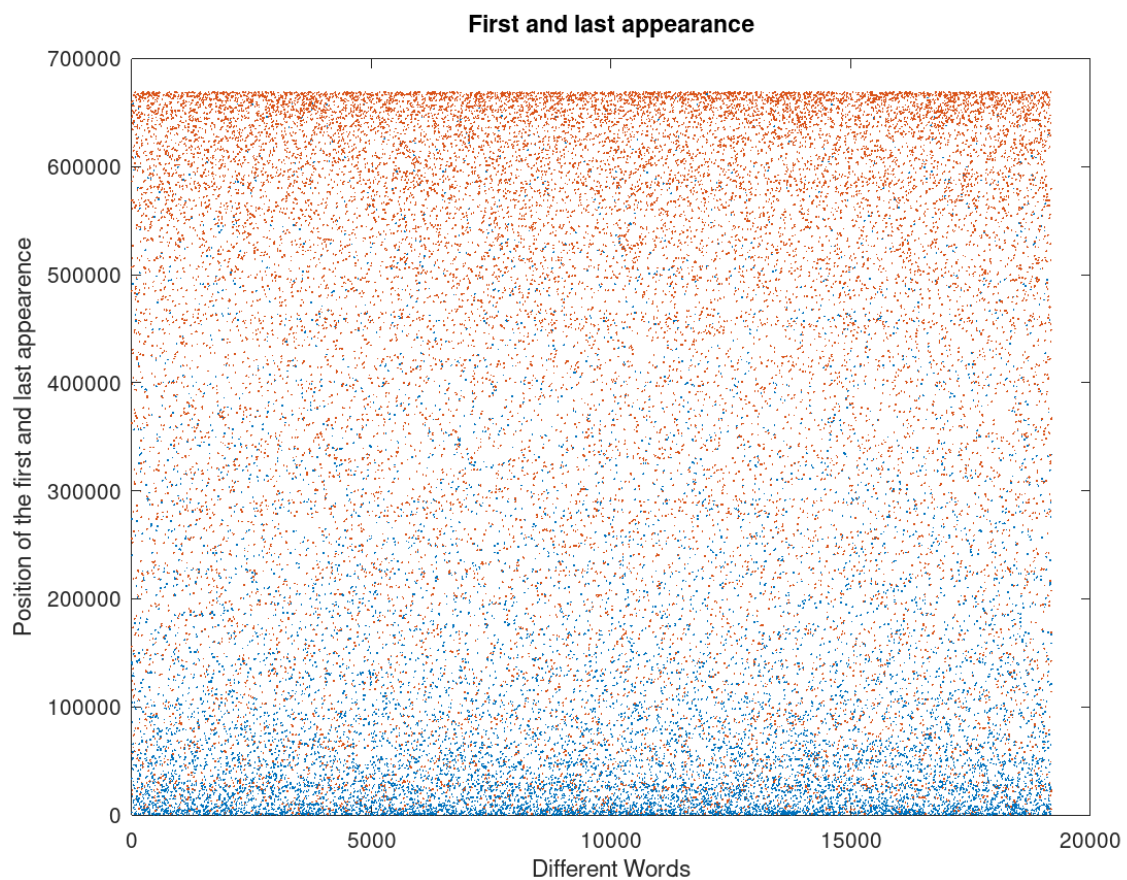


Figure 2- First (blue) and last (orange) appearances of different words in the text.

As expected, the position of the first appearance of a word is concentrated on the first 50000 words while the last appearance is concentrated on the last 50000 too. Still there are a considerable big number of words which appear, for the first or the last time, in the middle of the text, this is justified given the high amount of words who occur only once or twice in the book.

## Smallest, medium and largest distance

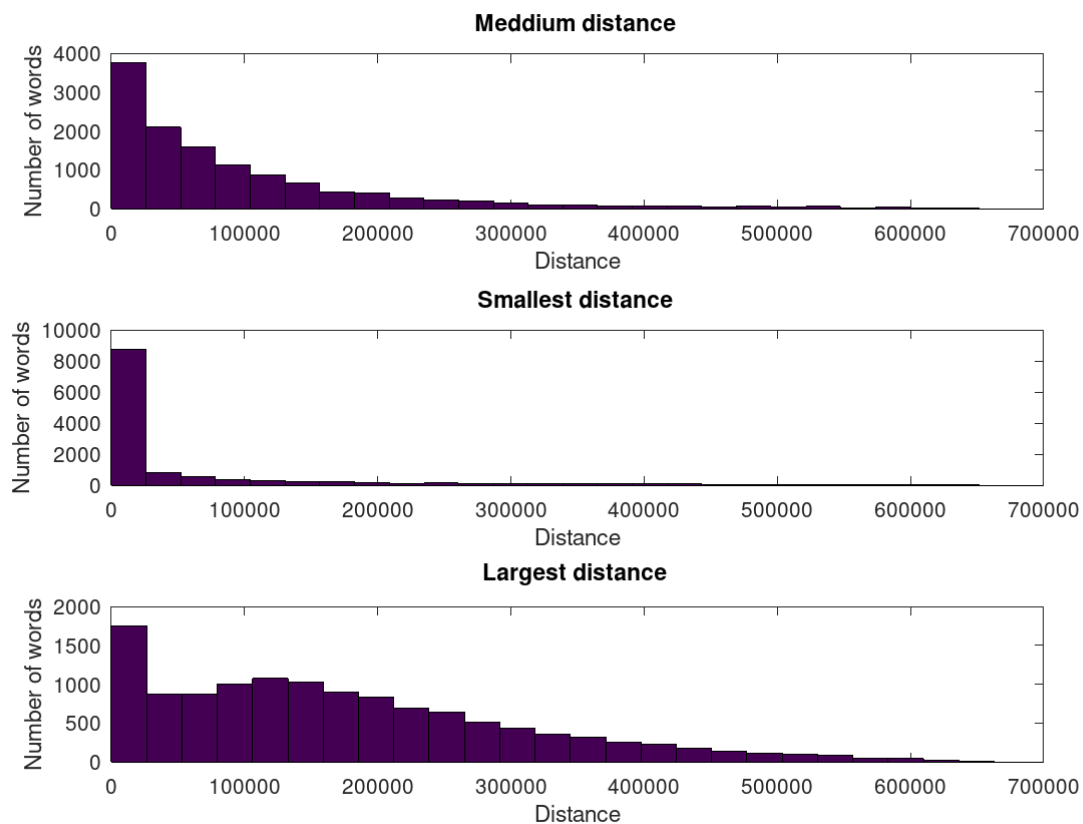


Figure 3- Medium, smallest and largest distance between equal words.

Almost every word has less than 100,000 words as the smallest distance, however the values for the largest distance are more scattered and interestingly a considerably big number of words (~1,750) has the largest distance below 25,000 words. Almost 4,000 words have a medium distance of less than 25,000.

## Time of execution and resize

For both the **linked list** and the **binary tree** approach it was established that the **hash tables** would resize when the number of different words was bigger than 5 times the size of the **hash table** (this value could be bigger for the **binary tree** approach given the properties of the structure).

Setting the increment of the **hash table** size after each resize at 1000, the time of execution of the programs was measured changing the initial size of the **hash table**. The results are shown below (figure 4).

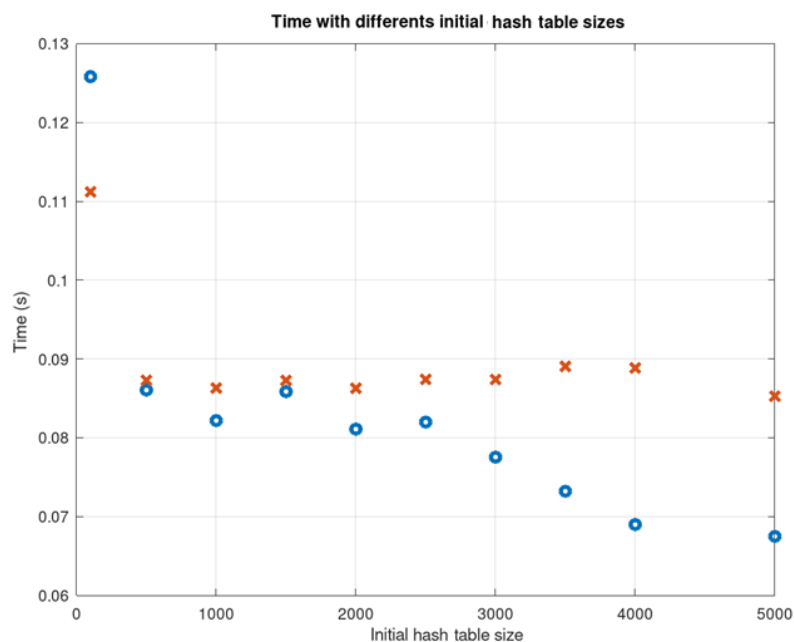


Figure 4 – Time of execution with different hash table sizes, for the binary tree approach (blue) and for the linked list approach (orange).

Hash table size	Resizes
100	4
500	4
1000	3
1500	3
2000	2
2500	2
3000	1
3500	1
4000	0
5000	0

Table 1 – Number of resizes done for each initial hash table size.

Apart from the **hash table** size 100, in the **linked list** approach the time of execution is quite constant as the **hash table** size increases. This doesn't happen with the **binary tree** approach where we can observe that the lowest times of execution occur when there is no resize at all (4000 and 5000).



Setting the initial size of the **hash table** at 2000, the time of execution of the programs was measured changing the increment on the **hash table** size when resizing. The results are shown below (figure 5).

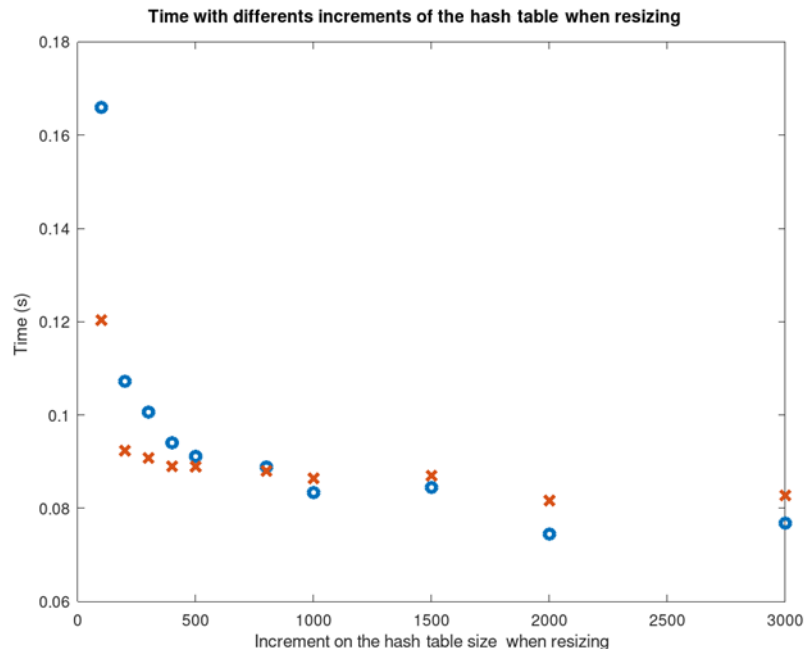


Figure 5- Time of execution with different hash table increments, for the binary tree approach (blue) and for the linked list approach (orange).

Hash table increment	Resizes
100	19
200	10
300	7
400	5
500	4
800	3
1000	2
1500	2
2000	1
3000	1

Table 2 – Number of resizes done for each hash table size increment.

Both approaches show a decrease of the time of execution when the increment is bigger (less resizes), however the **binary tree** approach decreases much faster having a time of execution lower than the **linked list** approach when the increment is bigger than 800.

Considering the results of the two graphs above we can conclude that if the number of elements to insert on a **hash table** is known (or belongs to a narrow interval), it is more efficient to use a **binary tree** approach, because searching an element will usually be faster, given that the tree is ordered. Meanwhile if the number of elements is unknown and it is necessary to resize the **hash table** multiple times the **linked list** approach is better because the cost of inserting a member should be taken account and insertion at the head is an  $O(1)$  operation, way faster than inserting an element on an **binary tree**.

# Appendix

## Occurrences\_tree.c

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <stdbool.h>
5  #include "manipulate_file.h"
6  #include <time.h>
7
8  static double elapsed_time(void)
9  {
10     static struct timespec last_time,current_time;
11
12     last_time = current_time;
13     if(clock_gettime(CLOCK_PROCESS_CPUTIME_ID,&current_time) != 0)
14         return -1.0; // clock_gettime() failed!!!
15     return ((double)current_time.tv_sec - (double)last_time.tv_sec)
16         + 1.0e-9 * ((double)current_time.tv_nsec - (double)last_time.tv_nsec);
17 }
18
19
20 int main(int argc,char **argv){
21     FILE *out= fopen("outBinaryTree_size.txt", "w");
22     for(int x=1;x<argc;x++){
23         file_data_t fd;
24         int idx;
25         int j=0;
26         int i=0;
27         int count=0;
28         int r_count=0;
29         (void)elapsed_time();
30         hash_size=(unsigned int)atoi(argv[x]);
31         struct hash_data_bt **hash_table= malloc(hash_size*sizeof(struct hash_data));
32         for(int k = 0;k < hash_size;k++) hash_table[k] = NULL;
33         open_text_file("SherlockHolmes.txt", &fd);
34         while(read_word(&fd) != -1) {
35             idx = hash_function(fd.word, hash_size);
36             hash_data_bt *hd=hash_table[idx];
37             while (hd != NULL && strcmp(fd.word,hd->key) != 0)
38             {
39                 if(strcmp(fd.word,hd->key)<0){
40                     hd=hd->left;
41                 }
42                 else
43                 {
44                     hd=hd->right;
45                 }
46
47                 i++;
48             }
49
50             if(hd == NULL ){
51                 if(i==0){
52                     hd=new_hash_data_bt();
53                     strcpy(hd->key, fd.word);
54                     hd->word.number_occurrences=1;
55                     hd->word.first_appearance=count;
56                     hd->word.last_appearance=count;
57                     hd->right=NULL;
58                     hd->left=NULL;
59                     hash_table[idx]=hd;
60
61                 }
```

```

62         else{
63             hd=new_hash_data_bt();
64             hd->word.number_occurrences=1;
65             hd->word.first_appearance=count;
66             hd->word.last_appearance=count;
67             hd->right=NULL;
68             hd->left=NULL;
69             strcpy(hd->key, fd.word);
70             insert_bt(hash_table[idx],hd);
71         }
72         if (j>hash_size*5){
73
74             hash_table=hash_resize_bt(hash_table,1000u);
75
76             r_count++;
77         }
78         j++;
79     }
80     else {
81         hd->word.number_occurrences++;
82         if(hd->word.number_occurrences==2){
83             hd->word.s_distance=count - hd->word.first_appearance;
84             hd->word.l_distance=count - hd->word.first_appearance;
85             hd->word.t_distance=hd->word.s_distance;
86         } else{
87             if(hd->word.s_distance > count-hd->word.last_appearance)
88                 hd->word.s_distance=count-hd->word.last_appearance;
89             if(hd->word.l_distance < count-hd->word.last_appearance)
90                 hd->word.l_distance=count-hd->word.last_appearance;
91             hd->word.t_distance+=count-hd->word.last_appearance;
92         }
93         hd->word.m_distance=hd->word.t_distance/(hd->word.number_occurrences -1);
94         hd->word.last_appearance=count;
95     }
96     i=0;
97     count++;
98     /*
99     printf("hash: %s %d \n", hash_table[idx]->key,hash_table[idx]->word.number_occurrences);
100     if(hash_table[idx]->left!=NULL){
101         printf("%s\n" , hash_table[idx]->left->key);
102     }
103     */
104 }
105 double cpu_time = elapsed_time();
106 close_text_file(&fd);
107 fprintf(out, "%d %.10f %d\n",(unsigned int)atoi(argv[x]), cpu_time,r_count);
108 free(hash_table);
109 }
110 fclose(out);
111 return 1;
112 }

```

## Occurrences.c

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <stdbool.h>
5  #include "manipulate_file.h"
6  #include <time.h>
7
8  static double elapsed_time(void)
9  {
10     static struct timespec last_time,current_time;
11
12     last_time = current_time;
13     if(clock_gettime(CLOCK_PROCESS_CPUTIME_ID,&current_time) != 0)
14         return -1.0; // clock_gettime() failed!!!
15     return ((double)current_time.tv_sec - (double)last_time.tv_sec)
16           + 1.0e-9 * ((double)current_time.tv_nsec - (double)last_time.tv_nsec);
17 }
18
19
20 int main(int argc,char **argv){
21     //FILE *out= fopen("outlinked_inc.txt", "w");
22     FILE *out= fopen("large_distance.txt", "w");
23     for(int x=1;x<argc;x++){
24         //printf("%d\n",(unsigned int)atoi(argv[x]));
25         file_data_t fd;
26         int idx;
27         int j=0;
28         int count=0;
29         int i=0;
30         int r_count=0;
31         (void)elapsed_time();
32         struct hash_data **hash_table= malloc(hash_size*sizeof(struct hash_data));
33         for(int k = 0;k < hash_size;k++) hash_table[k] = NULL;
34         open_text_file("SherlockHolmes.txt", &fd);
35         while(read_word(&fd) != -1) {
36             idx = hash_function(fd.word, hash_size);
37             hash_data *hd;
38             for(hd = hash_table[idx];hd != NULL && strcmp(fd.word,hd->key) != 0;hd = hd->next)
39                 i++;
40             if(hd == NULL ){
41                 if(i==0){
42                     hd=new_hash_data();
43                     strcpy(hd->key, fd.word);
44                     hd->word.number_occurrences=1;
45                     hd->word.first_appearance=count;
46                     hd->word.last_appearance=count;
47                     hd->next=NULL;
48                     hash_table[idx]=hd;
49
50                 }
51                 else{
52                     hd=new_hash_data();
53                     hd->next=hash_table[idx];
54                     strcpy(hd->key, fd.word);
55                     hd->word.number_occurrences=1;
56                     hd->word.first_appearance=count;
57                     hd->word.last_appearance=count;
58                     hash_table[idx]=hd;
59                 }
60             }
```

```

60         if (j>hash_size*5){
61
62             hash_table=hash_resize(hash_table,(unsigned int)atoi(argv[x]));
63             r_count++;
64         }
65         j++;
66     }
67     else {
68         hd->word.number_occurrences+=1;
69         if(hd->word.number_occurrences==2){
70             hd->word.s_distance=count - hd->word.first_appearance;
71             hd->word.l_distance=count - hd->word.first_appearance;
72             hd->word.t_distance=hd->word.s_distance;
73         } else{
74             if(hd->word.s_distance > count-hd->word.last_appearance)
75                 hd->word.s_distance=count-hd->word.last_appearance;
76             if(hd->word.l_distance < count-hd->word.last_appearance)
77                 hd->word.l_distance=count-hd->word.last_appearance;
78             hd->word.t_distance+=count-hd->word.last_appearance;
79         }
80         hd->word.m_distance=hd->word.t_distance/(hd->word.number_occurrences -1);
81         hd->word.last_appearance=count;
82     }
83     i=0;
84     count++;
85     //printf("hash: %s %d \n", hash_table[idx]->key,hash_table[idx]->word.number_occurrences);
86     //if(hash_table[idx]->next!=NULL){
87     //    printf("%s\n" , hash_table[idx]->next->key);
88     //}
89 }
90 double cpu_time = elapsed_time();
91 close_text_file(&fd);
92 //fprintf(out, "%d %.10f %d\n",(unsigned int)atoi(argv[x]), cpu_time,r_count);
93 for(int l=0;l<hash_size;l++){
94     while(hash_table[l]!=NULL){
95         fprintf(out,"%d\n",hash_table[l]->word.l_distance);
96         hash_table[l]=hash_table[l]->next;
97     }
98 }
99 }
100
101 free(hash_table);
102 hash_size=2000u;
103 }
104 fclose(out);
105 return 1;
106 }

```

## manipulate\_file.h

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <ctype.h>
5  #include <unistd.h>
6
7  // FILE STRUCT
8  typedef struct file_data_t {
9      // public data
10     long word_pos; // zero-based
11     long word_num; // zero-based
12     char word[64];
13     // private data
14     FILE *fp;
15     long current_pos; // zero-based
16 } file_data_t;
17
18 // HASH-FUNCTION
19 unsigned int hash_function(const char *str, unsigned int s) {
20     unsigned int h;
21     for(h = 0u; *str != '\0'; str++)
22         h = 157u * h + (0xFFu & (unsigned int)*str); // arithmetic overflow may occur here (just ignore it!)
23     return h % s; // due to the unsigned int data type, it is guaranteed that 0 <= h % s < s
24 }
25
26 // STRUCT THAT DEFINES THE STATS FOR THE WORDS IN THE HASH-TABLE
27 typedef struct word_stats {
28     int number_occurrences;
29     int first_appearance;
30     int last_appearance;
31     int s_distance; // smallest distance
32     int l_distance; // largest distance
33     int m_distance; // medium distance
34     int t_distance; // total distance
35 } word_stats;
36
37 // HASH-TABLE NODE LINKED LIST
38 typedef struct hash_data {
39     struct hash_data *next;
40     char key[64];
41     struct word_stats word;
42 } hash_data;
43
44 // HASH-TABLE NODE BINARY TREE
45 typedef struct hash_data_bt {
46     struct hash_data_bt *left;
47     struct hash_data_bt *right;
48     char key[64];
49     struct word_stats word;
50 } hash_data_bt;
```

```

52 // ALLOCATE NEW HASH-DATA
53 hash_data *new_hash_data(void) {
54     hash_data *hd = (hash_data *)malloc(sizeof(hash_data));
55     if(hd == NULL) {
56         fprintf(stderr, "Out of memory\n");
57         exit(1);
58     }
59     return hd;
60 }
61 // ALLOCATE NEW HASH-DATA
62 hash_data_bt *new_hash_data_bt(void) {
63     hash_data_bt *hd = (hash_data_bt *)malloc(sizeof(hash_data_bt));
64     if(hd == NULL) {
65         fprintf(stderr, "Out of memory\n");
66         exit(1);
67     }
68     return hd;
69 }
70
71 unsigned int hash_size= 2000u;
72
73 // INSERT
74 void insert_bt(hash_data_bt *root, hash_data_bt *hd ){
75     if(strcmp(hd->key, root->key)<0)
76     {
77         if(root->left!=NULL)
78             insert_bt(root->left,hd);
79         else
80             root->left=hd;
81     }
82
83     if(strcmp(hd->key, root->key)>0)
84     {
85         if(root->right!=NULL)
86             insert_bt(root->right,hd);
87         else
88             root->right=hd;
89     }
90 }
91
92 // OPEN FILE
93 int open_text_file(char *file_name, file_data_t *fd) {
94     fd->fp = fopen(file_name, "r");
95     if(fd->fp == NULL)
96         return -1;
97     fd->word_pos = -1;
98     fd->word_num = -1;
99     fd->word[0] = '\0';
100    fd->current_pos = -1;
101    return 0;
102 }
103
104
105 //RESIZE LINKED LIST
106 struct hash_data ** hash_resize(struct hash_data **hash_table, unsigned int inc){
107     //printf("Resizing.....\n");
108     hash_data *next;
109     hash_size+=inc;
110     struct hash_data **hash_table_new= malloc((hash_size)*sizeof(struct hash_data));
111     int new_idx;

```

```

112     for(int m = 0; m < hash_size; m++) hash_table_new[m] = NULL;
113     for(int l=0; l<hash_size-inc; l++){
114         while(hash_table[l]!=NULL){
115             new_idx= hash_function(hash_table[l]->key, hash_size);
116             next=hash_table[l]->next;
117             hash_table[l]->next=hash_table_new[new_idx];
118             hash_table_new[new_idx]=hash_table[l];
119             hash_table[l]=next;
120         }
121     }
122     }
123     free(hash_table);
124     return hash_table_new;
125 }
126 //TRAVERSE_TREE
127 void traverse_tree(struct hash_data_bt *hb, struct hash_data_bt **hash_table){
128     if(hb!=NULL){
129         int new_idx= hash_function(hb->key, hash_size);
130         if(hash_table[new_idx]==NULL){
131             hash_table[new_idx]=new_hash_data_bt();
132             strcpy(hash_table[new_idx]->key, hb->key);
133             hash_table[new_idx]->left=NULL;
134             hash_table[new_idx]->right=NULL;
135             hash_table[new_idx]->word=hb->word;
136         }
137         else{
138             hash_data_bt *temp =new_hash_data_bt();
139             temp->left=NULL;
140             temp->right=NULL;
141             temp->word=hb->word;
142             strcpy(temp->key, hb->key);
143             insert_bt(hash_table[new_idx], temp);
144         }
145         traverse_tree(hb->left, hash_table);
146         traverse_tree(hb->right, hash_table);
147     }
148 }
149 //RESIZE BINARY TREE
150 struct hash_data_bt ** hash_resize_bt(struct hash_data_bt **hash_table, unsigned int inc){
151     //printf("Resizing.....\n");
152     hash_data_bt *next;
153     struct hash_data_bt **hash_table_new= malloc((hash_size+inc)*sizeof(struct hash_data_bt));
154     int new_idx;
155     hash_size+=inc;
156     for(int m = 0; m < hash_size; m++) hash_table_new[m] = NULL;
157     for(int l=0; l<hash_size-inc; l++){
158         traverse_tree(hash_table[l], hash_table_new);
159     }
160     free(hash_table);
161     return hash_table_new;
162 }
163 // CLOSE FILE
164 void close_text_file(file_data_t *fd) {
165     fclose(fd->fp);
166     fd->fp = NULL;
167 }

```



```

168
169 // READ WORD
170 int read_word(file_data_t *fd) {
171     int i,c;
172     // skip white spaces
173     do {
174         c = fgetc(fd->fp);
175         if(c == EOF)
176             return -1;
177         fd->current_pos++;
178     } while(!((c >= 48 && c<58) || (c>=65 && c<=90) || (c>=97 && c<=122) || (c>=192)));
179     // record word
180     fd->word_pos = fd->current_pos;
181     fd->word_num++;
182     fd->word[0] = (char)c;
183     for(i = 1; i < (int)sizeof(fd->word) - 1; i++) {
184         c = fgetc(fd->fp);
185         if(c == EOF)
186             break; // end of file
187         fd->current_pos++;
188         if(!((c >= 48 && c<58) || (c>=65 && c<=90) || (c>=97 && c<=122) || (c>=192)))
189             break; // terminate word
190         fd->word[i] = (char)c;
191     }
192     fd->word[i] = '\0';
193     for(int j = 0; fd->word[j]; j++){
194         fd->word[j] = tolower(fd->word[j]);
195     }
196     return 0;
197 }
198

```