# AED - Algoritmos e Estruturas de Dados
# Hash Table implementation

Hash Tables implementation using singly linked lists and binary trees

| Dinis Cruz | Duarte Mortágua | Tiago Oliveira |
|---|---|---|
| 92080 | 92963 | 93456 |
| 33,3% | 33,3% | 33,3% |

Teacher: Tomás Oliveira e SIlva

Departamento de Eletrónica,

Telecomunicações e Informática

Universidade de Aveiro

december 27, 2019

# Contents

# List of Tables

# List of Figures

# 1. Introduction

## 1.1   A brief description of the problem

To start of, the main goal of this practical work was the implementation of an hash-table. The purpose of an hash-table implementation is to store data, so we can see an hash-table as a data structure. With the hast-table implementation the main goal was to store all the different words in a determined text and how many times every single word occurred in that text in each linked list entry corresponding each hash-table entry. Besides, we were told to store in each linked-list entry the first and last locations of the word and the maximum, minimum and medium distances between words. As we had to store a lot of information about each distended word, we thought that was clever to use structures as the representation of each word and the linked lists entries would be this structures. The language we used to implement this was C, although we are not very familiar with it we think we have done a good job using it. As soon as we started thinking about the implementation we decided

that the data that we were going to store would pass through an hash function first, however as we all know hashing methods always have collisions and as we want to have a efficient implementation we decided to use, in every hash-table entry, a linked list. Liked lists are basically separated elements that have a pointer to next element and so on, in this way every time we got an hashing collision we would put the next element in the linked list and the element that was already there with a pointer to it. As soon as we finish implementing the hash-table with linked lists we thought we could implement the hash-table but now with binary trees instead of linked lists to compare which one is the most efficient way.

# 2. Implementation methods

## 2.1 Singly linked lists implementation

As we described earlier, the linked list implementation is based in each element pointing to the next. Every time we have an hashing collision the linked list in that hashing position will grow. Nevertheless, every time we hash something to the hash-table we need to check if that position is still empty or not and if the word we are hashing is already in the hash-table or not. First of all, as soon as we have the hashing position that the word we are hashing will occupy, we go through the linked list in that position and in each element we check if it matches the word we are storing. If we find a match then we update all the data of that structure(which are what we used to represent each word) such as the number of times it occurred and the last, maximum, minimum and medium distances. On the other hand, while we are going trough the linked list if we reach a pointer that points do NULL then we know for sure that we are checking the last element of the linked list and as

we have not got any match this means that the word we are hashing still does not exists in the hash-table so we need to add him. As we already described the adding in the linked list is basically changing the pointer of the last element of the linked list to point to the new element so the new element becomes the last element of the that linked list.

## 2.2 Binary trees implementation

The binary trees implementation is based in having more efficiency in the search. Different than the linked list implementation, every structure in the hash table has a pointer to a binary tree. We thought that it would not make sense if we used unordered binary trees because that's what makes the search efficient. Of course every binary tree node has the same information as the linked list structure, such as the number of times the word occurred and the last, minimum, maximum and medium distances and a few more details we found imperative. On one hand, the collisions problem is solved exactly the same way the linked list implementation is, every time we have a collision, if the hashing position has already a pointer to the head of the binary tree, we go through all the binary tree and once we arrive the end of the tree if we found no match to the word then we add a new node to the binary tree.

On the other hand, the adding of a new node is different. As the binary tree is ordered, every time we have to add a new node to a tree, we compare the word we are adding with the head, if it is lower we go left, if it is greater we go right, and we go through the binary tree until we find the exact spot that fits the word.

About the resize, we did not implement it a long with the binary tree implementation because the efficiency is in the capacity of the search

being shorter, and so we think that a resize would not be a must need functionality (see 3.1.2).

# 3. Dynamic resizing

In order to reduce the number of collisions as the hash table gets filled, we implemented a dynamic resize function that doubles the size of the hash table and reallocates the words in the table. This action is triggered whenever the load factor is reached. The load factor is the number of linked lists (or binary trees) stored in the hash table divided by it's capacity, and in this case, we use a load factor of 0.5.

## 3.1   Approach description

Our hash table structure is not only characterized by the table itself, but by the elements "count" and "size" too. Whenever a word is read and occupies a new entry of the hash table, the count element is increased, so the count keeps track of the hash table slots that are not null. The size initially represents the number of slots of the hash table, equaling 2000. When $count/size >= 0.5$, the resize action is triggered.

### 3.1.1 Singly linked lists

We start the resize process by creating a new empty table of pointers to words (our new hash table). This new table will have the size equal to two times the old size. Then, we iterate thought the hash table entries. In each one, we detach the singly link list, keeping track of the head pointer in a new variable, "next". Then, as the size will be changed and so the hash-code too, we cant just move the entire linked list to a new entry of the new table, we have to detach the entire linked list and map each one of it's words to a new linked list in the new hash table. In order to do that, and having the "next" pointer, we iterate through all the words of each linked list, and since each word has an attribute "hash" (djb2 direct hash), we get the index of the word's new entry in the new hash table by getting the rest of the division of the hash by the size of the new hash table. After getting the index, instead of appending the word to the end of the new linked list, we attach it to the beginning. We do this by causing the new word to point to the first word in the list, and then replace the first word in the list with the new word.
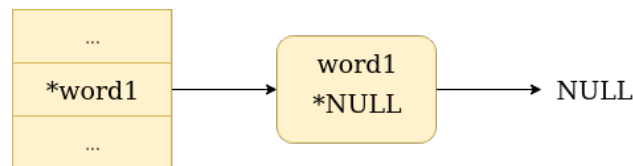


**Figure 3.1:** Adding word1 to empty linked list

**Figure 3.2:** Adding word2 to linked list containing word1

After doing this to all the linked lists in the old hash table, we free the memory of the old hash table, replace it by the new one and update the hash table size with the new size.

## 3.1.2  Binary trees

We decided not to implement dynamic resize when implementing the hash table with binary trees because:

- When talking about adding elements, resizing the table or not, the time complexity is $O(1)$ for both implementations, being that the dynamic resize would only have impact in time execution when searching for a word in the table.

- Each time we do the dynamic resize, in linked lists, we go through all the words in all the old lists one by one ($O(n)$) and add them to the new linked lists in the new table ($O(1)$), having a total complexity of ($O(n)$). On the other hand, in binary trees, we go through all the words in all the old trees one by one ($O(n)$) and then insert them on the correct position in the new binary trees ($O(n)$), having a total complexity of ($O(n^2)$). So, resizing when the

objective is adding words to the hash table isn't such a great idea when talking about binary trees implementation.

- Ignoring dynamic resize in binary trees doesn't affect the complexity that much when searching, because the fact of using ordered binary trees, although huge, decreases the search time of O(n) of the linked lists search to $O(logn)$, so it becomes fast enough to keep up with the search when dealing with linked lists (resized).

In sum, since the proposed problem was to add words to the hash table, we thought it would be less time expensive to implement an hash table with singly linked lists with dynamic resizing and an hash table with binary trees without dynamic resizing.

# 4. Tests

## 4.1  Singly linked lists tests

### 4.1.1  Output

To conclude, as we finished the linked lists implementation we started testing it so we could analyze if our implementation was doing what it was supposed to and also to prove that it was actually working. To help us doing the testing we used C asserts. First of all we made the program print the whole hash-table along with each words characteristics. We have also printed the number of words read, the number of words inserted in the hash-table, the size of the hash-table and the number hash-table slots that were occupied just for us to have an idea during the testing time, if we run the program with the same text file these numbers are supposed to be the same.

Next test was to verify if every single word read was in the correct

hash-table position, so after the hash-table was completed we read the text again and hashed every single word again and compare it to the word that was already in that position so we could assure the words were hashed correctly.

Finally the last test we have made was the top 10 words with more occurrences in the text file and the next step was to compare it with the binary tree (that will be approached in another topic). We also made a test only in linked lists that is commented because it takes to long to run (about 20min), which is to check if every single word occurs in the hash table once and only once. As we have so much words to read and search, this test has to check the whole hash-table the number of times as the number of words it would read so this explains the 20 min time running.

This tests were written along with the hash-table count in the end of the program so they would be showed as we run the program in the output.

We have also calculated some execution times, and concluded that the insertion of all the Sherlock Holmes book words (657438 words) takes about 0.11s (almost the same as the binary trees).

## 4.2 Binary trees solution

### 4.2.1 Output

The tests we have made about the binary tree implementation were exactly the same as in the linked list implementation, the only difference was in the way we go trough the hash-table as with binary trees it works different and also in the hash-table print, here we printed it in a graphical way so the words would actually be displayed in a tree format.

## 4.3   Solution comparison

To sum up, when the whole testing session was done we compared the execution times, concluding that the hash table implementation with binary trees (not resizable) takes near 0.14s and that the hash tables implementation with singly linked lists (resizable) takes near 0.11s to run.

# A. Code appendix

Below is the code of our ht_sll.c and ht_bt.c files respectively. All the comments made by us are in capital letter. Some comments made by the teacher were omitted in order to keep this appendix concise. These comments were no longer relevant since they were supposed to guide our code work in an early stage.

## A.1  Hash Tables with Singly Linked Lists (ht_sll.c)

```c
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <math.h>
5  #include <assert.h>
6  #include <time.h>
7
8  typedef struct file_data {
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  //Representa cada palavra distinta num determinado ficheiro
19  typedef struct word {
20      struct word* next;
21      char word[64];
22      unsigned long hash;
23      int first_location;
24      int last_location;
25      int max_dist;
26      int min_dist;
27      int medium_dist;
28      int count;
29  } word_t;
30
31  typedef struct hash_table {
32      unsigned int size;
33      unsigned int count;
34      word_t** table;
35  } hash_table_t;
36
37  int open_text_file(char* file_name, file_data_t* fd)
```

```
38  {
39      fd->fp = fopen(file_name, "r");
40      if (fd->fp == NULL){
41          printf("File does not exist.\n");
42          return -1;
43      }
44      fd->word_pos = -1;
45      fd->word_num = -1;
46      ;
47      fd->word[0] = '\0';
48      fd->current_pos = -1;
49      return 0;
50  }
51
52  void close_text_file(file_data_t* fd)
53  {
54      fclose(fd->fp);
55      fd->fp = NULL;
56  }
57
58  int read_word(file_data_t* fd)
59  {
60      int i, c;
61      // skip white spaces
62      do {
63          c = fgetc(fd->fp);
64          if (c == EOF)
65              return -1;
66          fd->current_pos++;
67
68      } while (c <= 32);
69      //record word
70      fd->word_pos = fd->current_pos;
71      fd->word_num++;
72      fd->word[0] = (char)c;
73      for (i = 1; i < (int)sizeof(fd->word) - 1; i++) {
74          c = fgetc(fd->fp);
75          if (c == EOF)
76              break;
77          // end of file
78          fd->current_pos++;
79          if (c <= 32)
80              break;
81          // terminate word
82          fd->word[i] = (char)c;
83      }
84      fd->word[i] = '\0';
85      return 0;
86  }
87
88  unsigned long hash(unsigned char* str)
89  {
90      unsigned long hash = 5381;
91      int c;
92
93      while (c = *str++)
94          hash = ((hash << 5) + hash) + c; /* hash * 33 + c */
95      return abs(hash);
96  }
97
98  int main(int argc, char* argv[])
99  {
100     file_data_t* fl;
101     fl = (file_data_t*)malloc(sizeof(file_data_t));
102
103     if (open_text_file(argv[1], fl) == -1) {
104         return EXIT_FAILURE;
105     }
106     double time_spent_total = 0;
107
108     clock_t begin_total = clock();
109
110     // HASHTABLE INITIALIZATION
111     hash_table_t* hash_table = NULL;
112     hash_table = malloc(sizeof(hash_table_t));
113     hash_table->table = malloc(2000 * sizeof(word_t*));
114     hash_table->size = 2000;
```

```
115        hash_table->count = 0;
116        for (int i = 0; i < hash_table->size; i++) {
117            hash_table->table[i] = NULL;
118        }
119
120        // POINTERS DECLARATION TO USE INSIDE WHILE CYCLE
121        word_t* head;
122        word_t* prev;
123        head = (word_t*)malloc(sizeof(word_t));
124        prev = (word_t*)malloc(sizeof(word_t));
125
126        int hashcode=0;
127        int word_counter=0;
128        int resize_counter=0;
129        double time_spent_resize = 0;
130
131        while (read_word(fl) != -1) {
132
133            //DYNAMIC RESIZE
134            if (hash_table->count >= hash_table->size / 2) {
135                clock_t begin_resize = clock();
136                word_t **table, *curr, *next;
137                size_t i, k;
138                next = malloc(sizeof(word_t));
139                curr = malloc(sizeof(word_t));
140                int new_size = hash_table->size * 2;
141                table = malloc(new_size * sizeof(word_t*));
142                if (!table)
143                    return -1; // OUT OF MEMORY
144
145                // INITIALIZE NEW TABLE TO EMPTY
146                for (i = 0; i < new_size; i++) {
147                    table[i] = NULL;
148                }
149
150                for (i = 0; i < hash_table->size; i++) {
151                    // DETACH THE SINGLY LINKED LIST
152                    next = hash_table->table[i];
153                    hash_table->table[i] = NULL;
154                    while (next) {
155                        // DETACH THE NEXT ELEMENT AS CURRENT
156                        curr = next;
157                        next = next->next;
158
159                        // K IS THE INDEX OF CURR IN THE NEW TABLE
160                        k = curr->hash % new_size; // o curr->hash   o resultado da word em curr ao passar pela fun   o
        hash()
161
162                        // PREPEND TO THE LINKED LIST IN TABLE[K]
163                        if (curr != table[k]) {
164                            curr->next = table[k];
165                            table[k] = curr;
166                        }
167                    }
168                }
169                // NO LONGER NEED NEXT AND CURR
170                free(next);
171                free(curr);
172
173                // NO LONGER NEED THE OLD HASH TABLE
174                free(hash_table->table);
175
176                // REPLACE THE OLD HASH TABLE WITH THE NEW ONE
177                hash_table->table = table;
178                hash_table->size = new_size;
179                clock_t end_resize = clock();
180                time_spent_resize = (double)(end_resize - begin_resize) / CLOCKS_PER_SEC;
181                resize_counter++;
182            }
183
184            word_counter++;
185            int flag = 0; //IF A WORD IS FOUND, THEN WE DONT NEED TO CREATE IT
186            hashcode = hash(fl->word) % hash_table->size;
187            head = hash_table->table[hashcode];
188            if (head == NULL) { // IF THERES NOTHING IN TABLE[HASHCODE], CREATE A NEW WORD THERE
189                word_t* new;
190                new = (word_t*)malloc(sizeof(word_t));
```

```
191              new->next = NULL;
192              new->hash = hash(fl->word);
193              new->first_location = fl->current_pos;
194              new->last_location = fl->current_pos;
195              new->max_dist = NULL;
196              new->min_dist = NULL;
197              new->medium_dist = 0;
198              new->count = 1;
199              strcpy(new->word, fl->word);
200              hash_table->table[hashcode] = new;
201              hash_table->count += 1;
202          }
203          else {
204              while (head != NULL) {
205                  if (strcmp(head->word, fl->word) == 0) { // IF MATCH IS FOUND
206                      flag = 1; // MATCH FOUND
207                      int temp = head->last_location;
208                      int dist = fl->current_pos - temp;
209                      head->last_location = fl->current_pos;
210                      if (dist > head->max_dist || head->max_dist == NULL) {
211                          head->max_dist = dist;
212                      }
213                      if (dist < head->min_dist || head->min_dist == NULL) {
214                          head->min_dist = dist;
215                      }
216                      head->medium_dist = head->medium_dist + (dist - head->medium_dist) / head->count;
217                      head->count++;
218                      break;
219                  }
220                  prev = head;
221                  head = head->next;
222              }
223              if (flag == 0) { // MATCH WAS FOUND? IF NOT, CREATE NEW WORD AND ATTACH IT TO THE LAST WORD (PREV)
224                  word_t* new;
225                  new = (word_t*)malloc(sizeof(word_t));
226                  new->next = NULL;
227                  new->hash = hash(fl->word);
228                  new->first_location = fl->current_pos;
229                  new->last_location = fl->current_pos;
230                  new->max_dist = NULL;
231                  new->min_dist = NULL;
232                  new->medium_dist = 0;
233                  new->count = 1;
234                  strcpy(new->word, fl->word);
235                  prev->next = new;
236              }
237          }
238      }
239      clock_t end_total = clock();
240      time_spent_total = (double)(end_total - begin_total) / CLOCKS_PER_SEC;
241
242      // HASHING DONE //
243
244      // TESTING //
245      int new_words = 0;
246      word_t* word;
247
248      // PRINT HASH TABLE
249      for (int k = 0; k < hash_table->size; k++) {
250          printf("%d: ", k);
251          if (hash_table->table[k] == NULL) {
252              printf("NULL\n");
253          }
254          else {
255              word = hash_table->table[k];
256              while (word->next != NULL) {
257                  printf("%s (FL: %d, LL: %d, MAXD: %d, MIND: %d, MEDD: %d, WC: %d) --> ", word->word, word->
    first_location, word->last_location, word->max_dist, word->min_dist, word->medium_dist, word->count);
258                  word = word->next;
259                  new_words++;
260              }
261              printf("%s (FL: %d, LL: %d, MAXD: %d, MIND: %d, MEDD: %d, WC: %d) --> NULL\n", word->word, word->
    first_location, word->last_location, word->max_dist, word->min_dist, word->medium_dist, word->count);
262              new_words++;
263          }
264      }
265
```

```
266        printf("================================\n");
267        printf("TABLE STATS\n");
268        printf("Words read: %d\n", word_counter);
269        printf("Hash table count: %d\n", hash_table->count);
270        printf("Hash table size: %d (resized %d times)\n", hash_table->size, resize_counter);
271        printf("Total duration: %5.4fs\n", time_spent_total);
272        printf("Duration of last resize: %5.4fs\n", time_spent_resize);
273        printf("Words inside table: %d\n", new_words);
274
275
276
277        // SEARCH TEST // --> ASSERT THAT WORD IS IN HASHTABLE
278
279        file_data_t* ft;
280        ft = (file_data_t*)malloc(sizeof(file_data_t));
281        word_t* head_test;
282        head_test = (word_t*)malloc(sizeof(word_t));
283        if (open_text_file("Teste.txt", ft) == -1) {
284            return EXIT_FAILURE;
285        }
286
287        int flag_search = 0;
288        while (read_word(ft) != -1) {
289            int hashcode_test=0;
290            hashcode_test = hash(fl->word) % hash_table->size;
291            head_test = hash_table->table[hashcode];
292            while (head_test) {
293                if (strcmp(head_test->word, fl->word) == 0){
294                    flag_search++;
295                }
296                head_test = head_test->next;
297            }
298
299        }
300
301        // // SEARCH TEST // --> ASSERT THAT A WORD APPEARS ONCE AND ONLY ONCE IN THE HASH TABLE (COMMENTED BACAUSE TAKES
              TOO MUCH TIME)
302
303        // file_data_t* ft;
304        // ft = (file_data_t*)malloc(sizeof(file_data_t));
305        // word_t* head_test;
306        // head_test = (word_t*)malloc(sizeof(word_t));
307        // if (open_text_file("Teste.txt", ft) == -1) {
308        //     return EXIT_FAILURE;
309        // }
310
311
312        // while (read_word(ft) != -1) {
313        //     int hashcode_test=0;
314        //     int flag_search = 0;
315
316        //     for (int k = 0; k < hash_table->size; k++){
317        //         head_test = hash_table->table[k];
318        //         while (head_test) {
319        //             if (strcmp(head_test->word, fl->word) == 0){
320        //                 flag_search++;
321        //             }
322        //             head_test = head_test->next;
323        //         }
324        //     }
325        //     assert(flag_search==1);
326        // }
327
328
329        // WORDS INSIDE TABLE TEST //
330
331        int test_counter=0;
332        word_t* head_test_2;
333        head_test_2 = (word_t*)malloc(sizeof(word_t));
334
335        for (int k = 0; k < hash_table->size; k++){
336            head_test_2 = hash_table->table[k];
337            while (head_test_2) {
338                test_counter++;
339                head_test_2 = head_test_2->next;
340            }
341        }
```

```
342
343        assert(new_words == test_counter);
344
345        // TOP 10 MORE FREQUENT WORDS //
346
347        char a[10][64];
348        int max_test = 0;
349        word_t* head_test_3;
350        char word_test[64];
351        int maxs[10];
352        head_test_3 = (word_t*)malloc(sizeof(word_t));
353        for (int j = 0; j < 10; j++){
354            max_test = 0;
355            for (int k = 0; k < hash_table->size; k++){
356                head_test_3 = hash_table->table[k];
357                while (head_test_3) {
358                    if (head_test_3->count > max_test){
359                        int flagg = 1;
360                        for (int l = 0; l < j; l++){
361                            if (strcmp(a[l],head_test_3->word) == 0){
362                                flagg = 0;
363                                break;
364                            }
365                        }
366                        if (flagg){
367                            max_test = head_test_3->count;
368                            memset(word_test, 0, 64);
369                            strcpy(word_test, head_test_3->word);
370                        }
371                    }
372                    head_test_3 = head_test_3->next;
373                }
374            }
375
376            strcpy(a[j], word_test);
377            maxs[j] = max_test;
378
379        }
380        printf("================================\n");
381        printf("TOP 10 MOST FREQUENT WORDS\n");
382        for (int k = 0; k < 10; k++){
383            printf("%-5s (%d)\n", a[k], maxs[k]);
384        }
385        return 0;
386
387 }
```

## A.2   Hash Tables with Binary Trees (ht_bt.c)

```
 1 #include <stdio.h>
 2 #include <stdlib.h>
 3 #include <string.h>
 4 #include <math.h>
 5 #include <assert.h>
 6 #include <time.h>
 7
 8 typedef struct file_data {
 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 //Representa cada palavra distinta num determinado ficheiro
19 typedef struct word {
20     struct word* left;
21     struct word* right;
22     char word[64];
23     int hash;
24     int first_location;
```

```
25      int last_location;
26      int max_dist;
27      int min_dist;
28      int medium_dist;
29      int count;
30  } word_t;
31
32  typedef struct hash_table {
33      unsigned int size;
34      unsigned int count;
35      word_t** table;
36  } hash_table_t;
37
38  int open_text_file(char* file_name, file_data_t* fd)
39  {
40      fd->fp = fopen(file_name, "r");
41      if (fd->fp == NULL){
42          printf("File does not exist.\n");
43          return -1;
44      }
45      fd->word_pos = -1;
46      fd->word_num = -1;
47      ;
48      fd->word[0] = '\0';
49      fd->current_pos = -1;
50      return 0;
51  }
52
53  void close_text_file(file_data_t* fd)
54  {
55      fclose(fd->fp);
56      fd->fp = NULL;
57  }
58
59  int read_word(file_data_t* fd)
60  {
61      int i, c;
62      // skip white spaces
63      do {
64          c = fgetc(fd->fp);
65          if (c == EOF)
66              return -1;
67          fd->current_pos++;
68
69      } while (c <= 32);
70      //record word
71      fd->word_pos = fd->current_pos;
72      fd->word_num++;
73      fd->word[0] = (char)c;
74      for (i = 1; i < (int)sizeof(fd->word) - 1; i++) {
75          c = fgetc(fd->fp);
76          if (c == EOF)
77              break;
78          // end of file
79          fd->current_pos++;
80          if (c <= 32)
81              break;
82          // terminate word
83          fd->word[i] = (char)c;
84      }
85      fd->word[i] = '\0';
86      return 0;
87  }
88
89  unsigned long
90  hash(unsigned char* str)
91  {
92      unsigned long hash = 5381;
93      int c;
94
95      while (c = *str++)
96          hash = ((hash << 5) + hash) + c; /* hash * 33 + c */
97
98      return abs(hash);
99  }
100 void print2DUtil(word_t* root, int space)
101 {
```

```
102        // Base case
103        if (root == NULL) {
104            printf("\n");
105            for (int i = 10; i < space + 10; i++)
106                printf(" ");
107            printf("%15s\n", NULL);
108            return;
109        }
110        // Increase distance between levels
111        space += 10;
112
113        // Process right child first
114        print2DUtil(root->right, space);
115
116        // Print current node after space
117        // count
118        printf("\n");
119        for (int i = 10; i < space; i++)
120            printf(" ");
121        printf("%15s\n", root->word);
122
123        // Process left child
124        print2DUtil(root->left, space);
125    }
126
127    void visit(word_t* word)
128    {
129        printf("%s (FL: %d, LL: %d, MAXD: %d, MIND: %d, MEDD: %d, WC: %d)\n", word->word, word->first_location, word->
                last_location, word->max_dist, word->min_dist, word->medium_dist, word->count);
130    }
131
132    word_t* search_recursive(word_t* link, char* data)
133    {
134        if (link == NULL || strcmp(link->word, data) == 0)
135            return link;
136        if (strcmp(link->word, data) > 0)
137            return search_recursive(link->left, data);
138        else
139            return search_recursive(link->right, data);
140    }
141
142    void insert_non_recursive(word_t** link, word_t** insert, char* data)
143    {
144        word_t* parent = NULL;
145        while (*link != NULL) {
146            parent = *link;
147            link = (strcmp((*link)->word, data) > 0) ? &((*link)->left) : &((*link)->right);
148        }
149        *link = *insert;
150    }
151
152    void traverse_in_order_recursive(word_t* link)
153    {
154        if (link != NULL) {
155            traverse_in_order_recursive(link->left);
156            //printf("--------------------------\n");
157            visit(link);
158            //printf("--------------------------\n");
159            traverse_in_order_recursive(link->right);
160        }
161    }
162    void most_used_words(hash_table_t* hash_table)
163    {
164        word_t* head_test;
165        word_t* checker;
166        head_test = (word_t*)malloc(sizeof(word_t));
167        checker = (word_t*)malloc(sizeof(word_t));
168        int max_prev = __INT_MAX__;
169        int max_in_cicle = 0;
170        int occ;
171        char palavra[64] = "word";
172        for (int c = 0; c < 10; c++) {
173            max_in_cicle = 0;
174            for (int f = 0; f < hash_table->size; f++) {
175
176                checker = hash_table->table[f];
177                if (checker->count < max_prev) {
```

```
178                if (checker->count > max_in_cicle) {
179                    strcpy(palavra, checker->word);
180                    max_in_cicle = checker->count;
181                    occ = checker->count;
182                }
183            }
184        }
185        printf("%-5s (%d)\n", palavra, occ);
186        max_prev = max_in_cicle;
187    }
188 }
189
190 int main(int argc, char* argv[])
191 {
192     file_data_t* fl;
193     fl = (file_data_t*)malloc(sizeof(file_data_t));
194
195     if (open_text_file(argv[1], fl) == -1) {
196         return EXIT_FAILURE;
197     }
198
199     double time_spent_total = 0;
200
201     clock_t begin_total = clock();
202
203     // HASHTABLE INITIALIZATION
204     hash_table_t* hash_table = NULL;
205     hash_table = malloc(sizeof(hash_table_t));
206     hash_table->table = malloc(2000 * sizeof(word_t*));
207     hash_table->size = 2000;
208     hash_table->count = 0;
209     for (int i = 0; i < 2000; i++) {
210         hash_table->table[i] = NULL;
211     }
212
213     // HEAD DECLARATION TO USE INSIDE WHILE CYCLE
214     word_t* head;
215     head = (word_t*)malloc(sizeof(word_t));
216     int counterrr = 0;
217     int hashcode = 0;
218     int word_counter = 0;
219     while (read_word(fl) != -1) {
220         word_counter++;
221         hashcode = hash(fl->word) % hash_table->size;
222         head = hash_table->table[hashcode];
223         if (head == NULL) {
224             word_t* new;
225             new = (word_t*)malloc(sizeof(word_t));
226             new->left = NULL;
227             new->right = NULL;
228             new->hash = hash(fl->word);
229             new->first_location = fl->current_pos;
230             new->last_location = fl->current_pos;
231             new->max_dist = NULL;
232             new->min_dist = NULL;
233             new->medium_dist = 0;
234             new->count = 1;
235             strcpy(new->word, fl->word);
236             hash_table->table[hashcode] = new;
237             hash_table->count += 1;
238             counterrr++;
239         }
240         else {
241             word_t* this_one;
242             this_one = search_recursive(head, fl->word);
243             if (this_one != NULL) {
244                 int temp = this_one->last_location;
245                 int dist = fl->current_pos - temp;
246                 this_one->last_location = fl->current_pos;
247                 if (dist > this_one->max_dist || this_one->max_dist == NULL) {
248                     this_one->max_dist = dist;
249                 }
250                 if (dist < this_one->min_dist || this_one->min_dist == NULL) {
251                     this_one->min_dist = dist;
252                 }
253                 this_one->medium_dist = this_one->medium_dist + (dist - this_one->medium_dist) / this_one->count;
254                 this_one->count++;
```

```
255                }
256            else {
257                word_t* new1;
258                new1 = (word_t*)malloc(sizeof(word_t));
259                new1->left = NULL;
260                new1->right = NULL;
261                new1->hash = hash(fl->word);
262                new1->first_location = fl->current_pos;
263                new1->last_location = fl->current_pos;
264                new1->max_dist = NULL;
265                new1->min_dist = NULL;
266                new1->medium_dist = 0;
267                new1->count = 1;
268                strcpy(new1->word, fl->word);
269                insert_non_recursive(&head, &new1, new1->word);
270                counterrr++;
271            }
272        }
273    }
274
275    clock_t end_total = clock();
276    time_spent_total = (double)(end_total - begin_total) / CLOCKS_PER_SEC;
277
278    word_t* word;
279    for (int k = 0; k < hash_table->size; k++) {
280        printf("==================\n");
281        if (hash_table->table[k] == NULL) {
282            printf("NULL\n");
283        }
284        else {
285            word = hash_table->table[k];
286            traverse_in_order_recursive(word);
287            //print2DUtil(word,0);
288        }
289        printf("==================\n");
290        printf("LINE: %d ABOVE\n", k);
291    }
292    printf("==================\n");
293    printf("TABLE STATS\n");
294    printf("Words read: %d\n", word_counter);
295    printf("Hash elements count: %d\n", hash_table->count);
296    printf("Hash elements size: %d\n", hash_table->size);
297    printf("Total duration: %5.4fs\n", time_spent_total);
298    printf("Number of words inside hash table: %d\n", counterrr);
299
300    file_data_t* ft;
301    ft = (file_data_t*)malloc(sizeof(file_data_t));
302    word_t* head_test;
303    word_t* checker;
304    head_test = (word_t*)malloc(sizeof(word_t));
305    checker = (word_t*)malloc(sizeof(word_t));
306    if (open_text_file("Teste.txt", ft) == -1) {
307        return EXIT_FAILURE;
308    }
309    int count_words = 0;
310    int hashcode_test = 0;
311    while (read_word(ft) != -1) {
312        hashcode_test = 0;
313        int flag_test = 0;
314
315        hashcode_test = hash(fl->word) % hash_table->size;
316        head_test = hash_table->table[hashcode_test];
317        checker = search_recursive(head_test, fl->word);
318        if (checker != NULL) {
319            flag_test = 1;
320            count_words = count_words + checker->count;
321        }
322        assert(flag_test);
323    }
324    assert(count_words == word_counter);
325    printf("==================\n");
326    printf("TOP 10 MOST FREQUENT WORDS\n");
327    most_used_words(hash_table);
328    return 0;
329 }
```