

Huffman Codes

Group 18

Date:2020-04-30

CONTENT

Chapter 1 Introduction

- 1.1 Background
- 1.2 Problem Description

Chapter 2 Algorithm Specification

- 2.1 Sketch and Data Structure
- 2.2 Function Description
 - 2.2.1 Minimum-Heap
 - 2.2.1.1 Insert
 - 2.2.1.1 Delete_Min
 - 2.2.2 Construct a Huffman tree
 - 2.2.3 Compare

Chapter 3 Testing Results

- 3.1 Test case 1
 - Purpose of this case
 - Input(sample)
 - Output
 - Result
- 3.2 Test case 2
 - Purpose of this case
 - Input
 - Output
 - Result
- 3.3 Test case 3
 - Purpose of this case
 - Input
 - Output
 - Result
- 3.4 Test case 4
 - Purpose of this case
 - Input
 - Output
 - Result
- 3.5 Test case 5
 - Purpose of this case
 - Input
 - Output
 - Result
- 3.6 Test case 6
 - Purpose of this case
 - Input
 - Output
 - Result
- 3.7 Test case 7
 - Purpose of this case
 - Input
 - Output
 - Result

Chapter 4 Analysis and Comments

- 4.1 Analysis on Time Complexity
- 4.2 Analysis on Space Complexity
- 4.3 Comments

Appendix

- Appendix I Source Code in C

Appendix II Testing Set Generating Script in C++

Appendix III Declaration and Signatures

Declaration

Signatures

Chapter 1 Introduction

1.1 Background

Huffman coding is a type of variable word length coding (VLC). Huffman proposed the coding method in 1952. This method completely constructs the codeword with the shortest average length of different prexes based on the occurrence probability of characters. It is sometimes called the best coding. And it is generally called Huffman coding. The following cites a theorem, which ensures that the code length is allocated according to the occurrence probability of characters, so that the average code length is the shortest.

There are diverse variable word length coding for one array of characters, while obviously not all of them are correct. Only some specic VLC can be called as correct Huffman codes and the detailed strategies to judge them are described in this report.

1.2 Problem Description

In this project, for a specific array of characters with defined frequency, we need to find whether the given codes are correct or not. Note that the optimal solution is not necessarily generated by Huffman algorithm. Any prefix code with code length being optimal is considered correct. So we need to generate a variety of codes strategies for our tests and use our program to check their correctness.

Chapter 2 Algorithm Specification

2.1 Sketch and Data Structure

Here is our program sketch:

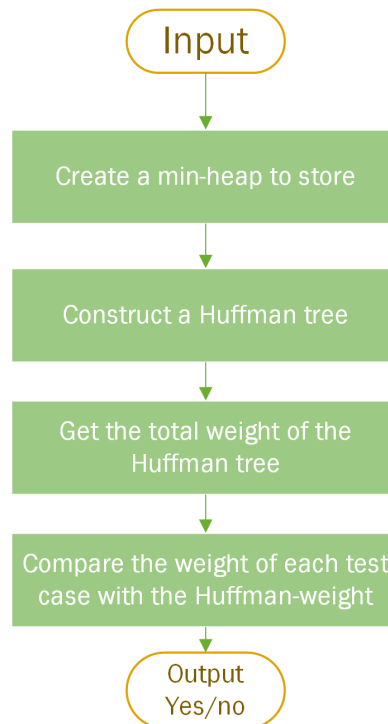


Figure1: Program Sketch

2.2 Function Description

There are several functions used in the program, complementing our main algorithm. We hereby list some functions which are acting as auxiliary in the program to make our report more completed.

2.2.1 Minimum-Heap

2.2.1.1 Insert

The sketch of the function `heap_insert` is as following:

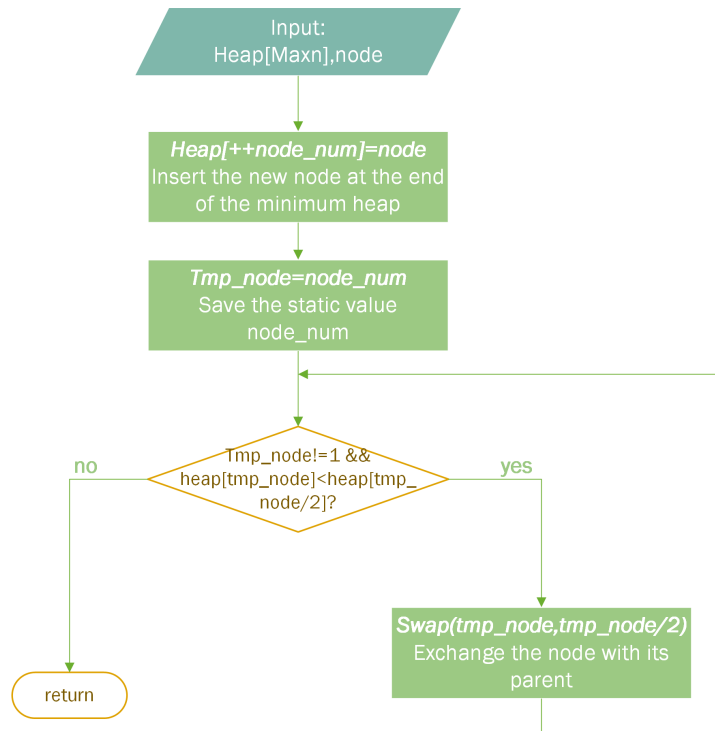


Figure2: Program Sketch of min-heap insertion

The pseudocode to insert elements in to a minimum heap is as following:

```
1 function heap_insert(ptr_huffman_node node)
2   heap[++node_num] <- node
3   this <- node_num
4   while (this != 1 and heap[this]->frequency < heap[this/2]->frequency)
5     do swap(this, this/2);
6     this /= 2;
7   end
8   return
```

For example, if we want to insert the element **0** into the minimum tree as following:

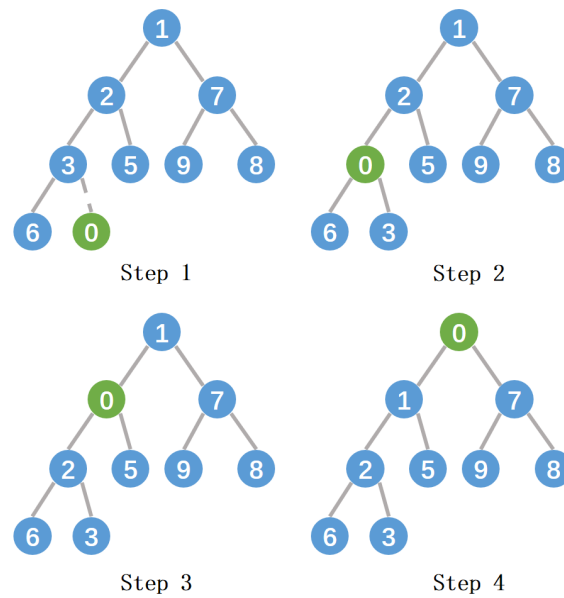


Figure3: Sample of inserting in min-heap

2.2.1.1 Delete_Min

The sketch of the function **delete_min** is as following:

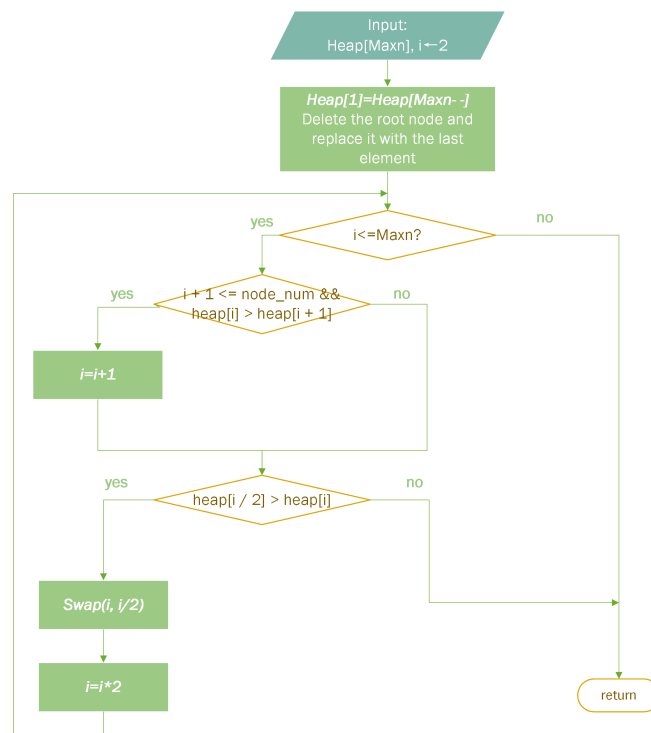


Figure4: Program Sketch of min-heap deletion

The pseudocode to delete the minimum element in a minimum heap is as following:

```

1  function heap_delete()
2      heap[1] ← heap[node_num--]
3      for i ← -2 to node_num
4          do
5              if (i + 1 ≤ node_num && heap[i] > heap[i + 1])
6                  do i ← i + 1
7              end
8              if (heap[i / 2] > heap[i])
9                  do swap(i, i/2)
10             else
11                 do break
12             end
13             i ← i * 2
14         end
15     return

```


For example, if we want to delete the minimum element **0** in the minimum tree as following:

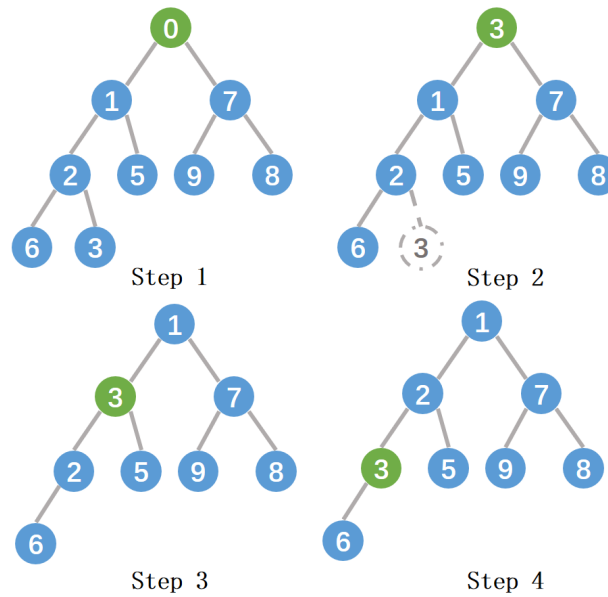


Figure5: Sample of deleting the min element

2.2.2 Construct a Huffman tree

The basic idea of Huffman coding is to construct a Huffman tree with the use frequency as the weight, and then use Huffman tree to code the characters. To construct a Huffman tree, the character to be encoded is taken as the leaf node, and the frequency of using the character in the file is taken as the weight of the leaf node. In a bottom-up way, a tree is constructed after $n - 1$ **merging** operations. The core idea is that the larger the weight, the closer the leaf is to the root.

The form of Huffman structure is like:

```

1 typedef struct huffman_node* ptr_huffman_node ;
2 struct huffman_node {
3     char chr;                // The letter of A-Z,a-z,_
4     int frequency;           // The number of this letter's appearances
5     ptr_huffman_node left;    // The pointer to the left child
6     ptr_huffman_node right;   // The pointer to the right child
7 };

```

The sketch of the function **construct_huffman** is as following:

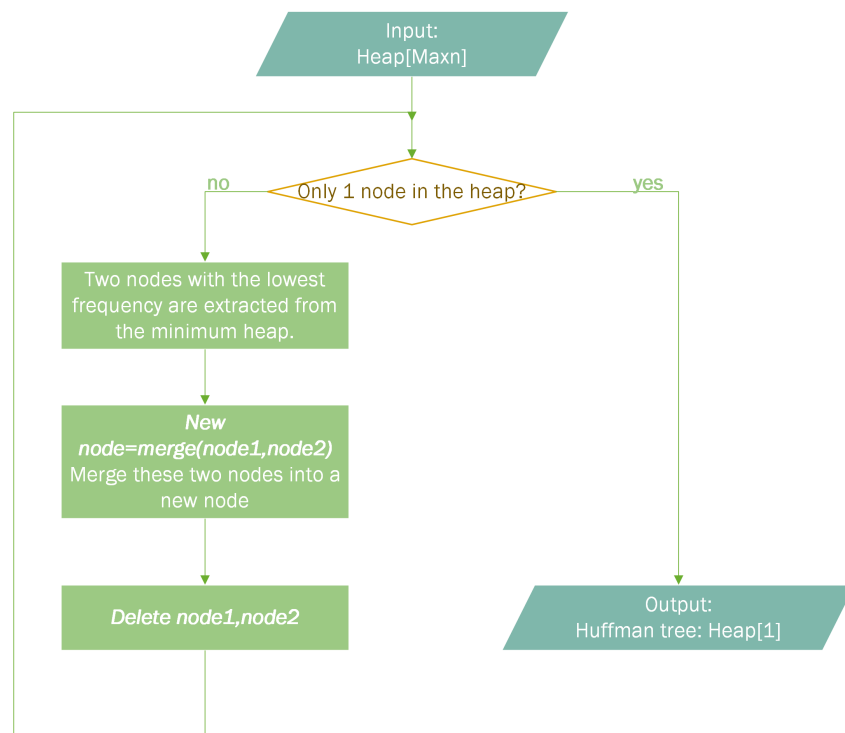


Figure6: Program Sketch of Huffman tree construction

The pseudocode to construct a Huffman tree is as following:

```

1  function construct_huffman()
2      while (there are more than 1 node in heap)
3          do
4              node1 ← Heap[1]
5              delete_min()
6              node2 ← Heap[1]
7              delete_min()
8              insert(merge(node1, node2))
9          end
10     return

```

For example, suppose that the input min-heap is as following:

Heap	1	2	3	4	5
char	A	B	C	D	E
frequency	1	2	4	3	7

Obviously, the two nodes with the least frequency are A and B. So we merge them and insert the new node into the min-heap.

Heap	1	2	3	4
char	D	A+B	C	E
frequency	3	3	4	7

Then, the two nodes with the least frequency are A+B and D. So we merge them and insert the new node into the min-heap.

Heap	1	2	3
char	C	E	D+(A+B)
frequency	4	7	6

Now, the two nodes with the least frequency are C and D+(A+B). So we merge them and insert the new node into the min-heap.

Heap	1	2
char	E	C+(D+(A+B))
frequency	7	10

Finally, merge the remain two nodes, we can get the Huffman tree as following:

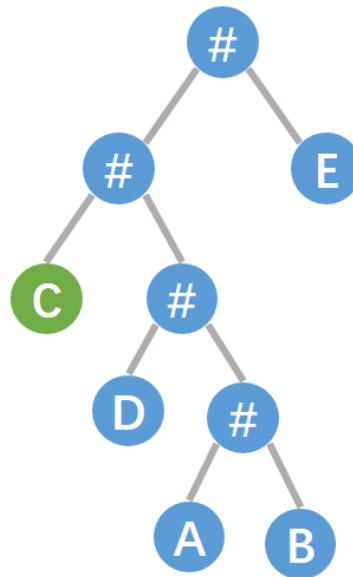


Figure6: Sample of Huffman tree

2.2.3 Compare

To judge whether the student's submission is correct, we just need to compare the weight of Huffman tree we construct before with the weight of the solution that student submits.

The pseudocode to calculate the weight of Huffman tree is as following:

```

1  function calculate_weight(ptr_huffman_node node, int level)
2      if (node == NULL) return;
3      calculate_weight(node->left, level + 1)
4      calculate_weight(node->right, level + 1)
5      if (node->chr != '#')
6          do weight += node->frequency * level
7          // Since I use '#' to mark it non-leaf
8      end
9  end

```

The weight of the sample showed in **Figure 6** is 36.

If one of the solutions submitted is as following:

char	A	B	C	D	E
frequency	1	2	4	3	7
code	0110	0111	010	00	1

The weight of it can be calculated as:

$$weight = 1 \times 4 + 2 \times 4 + 4 \times 3 + 3 \times 2 + 7 \times 1 = 36$$

So the weight of the submission is equal to the Huffman tree. In order to judge the correctness, then we also need to check whether there is an element's code is the foreword or subsequence of another element.

```
1 function check_prefix(int test_groups_num)
2   for i<-0 to test_groups_num step 1
3     do for j<-0 to test_groups_num step 1
4       if i=j continue
5       result<-strstr(test_table[i],test_table[j])
6       if result==NULL
7         do continue
8       if result=test_table[i]
9         return 0
10      end
11    end
12    return 1
13 end
```

Chapter 3 Testing Results

We partly use the c++ testing case generating program in Appendix II to generate our test case.

3.1 Test case 1

Purpose of this case

Juxtaposition, multi branches; length wrong; length right but prefix wrong; only uppercase characters

Input(sample)

```
1 7
2 A 1 B 1 C 1 D 3 E 3 F 6 G 6
3 4
4 A 00000
5 B 00001
6 C 0001
7 D 001
8 E 01
9 F 10
10 G 11
11 A 01010
12 B 01011
13 C 0100
14 D 011
15 E 10
16 F 11
17 G 00
18 A 000
19 B 001
20 C 010
21 D 011
22 E 100
23 F 101
24 G 110
25 A 00000
26 B 00001
27 C 0001
28 D 001
29 E 00
30 F 10
31 G 11
```

Output

```
1 Yes
2 Yes
3 No
4 No
```

Result

Correct

3.2 Test case 2

Purpose of this case

Lowercase letters with 0, 1 reversed and 2 points swapping; 2 points overlap

Input

```
1 | 3
2 | a 13 b 3 c 78
3 | 1
4 | a 01
5 | b 00
6 | c 1
```

Output

```
1 | Yes
```

Result

Correct

3.3 Test case 3

Purpose of this case

Unequal length of codes but all right; equal length but wrong prefix; length of codes greater than N

Input

```
1 4
2 A 4 B 4 C 6 D 6
3 1
4 A 111
5 B 1
6 C 0
7 D 001
```

Output

```
1 No
```

Result

Correct

3.4 Test case 4

Purpose of this case

Maximum $N = 63$ & $M = 1000$

Input

Refer to `max_in.txt` in the `code` folder.

Output

Refer to `max_out.txt` in the `code` folder.

Result

Correct

3.5 Test case 5

Purpose of this case

Minimum $N = 2$ & $M = 1$

Input

```
1 | 2
2 | A 0 B 1
3 | 1
4 | A 1
5 | B 0
```

Output

```
1 | Yes
```

Result

Correct

3.6 Test case 6

Purpose of this case

The number of characters to be encoded is double, while the length to be submitted is equal.

Input

```
1 | 4
2 | A 3 B 5 C 6 D 7
3 | 1
4 | A 00
5 | B 01
6 | C 10
7 | D 11
```

Output

```
1 | Yes
```

Result

Correct

3.7 Test case 7

Purpose of this case

Not Huffman coded, but correct

Input

1	3
2	A 0 B 0 C 8

Output

1	Yes
---	-----

Result

Correct

Chapter 4 Analysis and Comments

4.1 Analysis on Time Complexity

We can divide the program into three parts.

The first part is the input part. There is only one loop in this function, whose time complexity is $O(N)$.

The second part is to calculate the total weight of the huffman tree. When analyzing time complexity, Function *calculate_weight* is mainly about a two-layer loop, which is used to find the two smallest in the array of frequency, sum them up to get a new number and delete those two original numbers. Note that we use minimum heap to implement this design, so the time complexity is reduced to $O(N\log N)$.

The third part is to judge which of those M sets of codes is correct. It's also a two-layer loop. The outer loop repeats M times while the inner loop N times, so we can conclude that the time complexity of this part is $O(M*N)$ as we must traverse all M sets of codes and each of them consists of N lines .

In conclusion, when M is much greater than N, the time complexity of this program is $O(M*N)$. If not, then it's $O(N\log N)$.

4.2 Analysis on Space Complexity

Although a lot of memory is allocated dynamically in the construction of Huffman tree, most of this memory is only a linear function of N . For example, we use array to store frequencies and characters, both of which use N spaces. In conclusion, the total space complexity is $O(N)$.

4.3 Comments

The basic idea of Huffman coding is to construct a Huffman tree with the use frequency as the weight, and then use Huffman tree to code the characters. To construct a Huffman tree, the character to be encoded is taken as the leaf node, and the frequency of using the character in the file is taken as the weight of the leaf node.

In a bottom-up way, a tree is constructed after $n-1$ "merging" operations. The core idea is that the larger the weight, the closer the leaf is to the root. The greedy strategy adopted by Huffman algorithm is to take two trees without parents and with the least weight from the set of trees as left and right subtrees, which gives us a further understanding of greedy algorithm.

Appendix

Appendix I Source Code in C

```
1  #include <stdio.h>
2  #include <string.h>
3  #include <stdlib.h>
4
5  #define MAX_CODING_LENGTH 65    // Slightly larger than 63, which is the maximal length of input
6  #define ASCII 128              // Use ascii code for indexing, reducing time
7
8  // This structure describes a huffman node in the huffman tree
9  typedef struct huffman_node* ptr_huffman_node ;
10 struct huffman_node {
11     char chr;                    // The letter of A-Z,a-z,_
12     int frequency;               // The number of this letter's appearances
13     ptr_huffman_node left;      // The pointer to the left child
14     ptr_huffman_node right;     // The pointer to the right child
15 };
16
17 int weight = 0;                 // The weight of coding using huffman
18 tree
19 int node_num = 0;               // The number of nodes in the heap
20 int frequency[ASCII];           // Helps in finding the frequency
21 quickly
22 ptr_huffman_node heap[MAX_CODING_LENGTH] = {NULL}; // The heap that stores huffman nodes
23 char test_table[MAX_CODING_LENGTH][MAX_CODING_LENGTH]; // The input case in each test
24
25 /**
26  * @Author      : C01dkit
27  * @Date        : 21:10   2020/4/25
28  * @Param       : node
29  * @return      : void
30  * @Output      : none
31  * @Description : This function will insert a pointer in the heap
32  */
33 void heap_insert(ptr_huffman_node node);
34
35 /**
36  * @Author      : C01dkit
37  * @Date        : 21:35   2020/4/25
38  * @Param       : smaller_index  ----- The index of the node whose frequency is smaller in
39  the heap
40  *               larger_index  ----- The index of the node whose frequency is larger in the
41  heap
42  * @return      : void
43  * @Output      : node
44  * @Description : This function will swap two pointers in the heap using their index.
45  */
46 void swap(int smaller_index, int larger_index);
47
48 /**
49  * @Author      : C01dkit
50  * @Date        : 21:41   2020/4/25
51  * @Param       : none
52  * @return      : void
53  * @Output      : none
54  * @Description : This function will delete heap[1] and return nothing
55  */
56 void heap_delete();
57
58 /**
59  * @Author      : C01dkit
60  * @Date        : 22:11   2020/4/25
61  * @Param       : node  ----- The pointer to the root of huffman tree
62  *               level  ----- Current depth of the tree, start at 0
63  * @return      : void
64  * @Output      : none
65  * @Description : This function will do a post-order traversal and calculate total weight of the
66  huffman tree
67  */
68 void calculate_weight(ptr_huffman_node node,int level);
69
70 /**
71  * @Author      : C01dkit
72  * @Date        : 22:39   2020/4/25
73  * @Param       : test_groups_num  ----- The number of strings going to check
74  * @return      : int
75  * @Output      : none
76  * @Description : This function will check current group of test.
77  *               If it obeys prefix rule, then the function will return 1, otherwise 0.
78  */
```



```

73  */
74  int check_prefix(int test_groups_num);
75
76  int main () {
77      int N;
78      ptr_huffman_node temp;
79      scanf("%d",&N);
80      getchar();
81      // This loop will restore the initial data and construct a heap
82      for (int i = 0; i < N; i++) {
83          temp = (ptr_huffman_node)malloc(sizeof(struct huffman_node));
84          if (i == 0) scanf("%c %d",&temp->chr,&temp->frequency);
85          else scanf(" %c %d",&temp->chr,&temp->frequency);
86          temp->left = NULL;
87          temp->right = NULL;
88          frequency[(int)temp->chr] = temp->frequency;
89          heap_insert(temp);
90      }
91      // This loop will construct a huffman tree using the heap
92      while (node_num > 1) { // If the heap only contains one node, then it is the root of the
huffman tree
93          temp = (ptr_huffman_node)malloc(sizeof(struct huffman_node));
94          temp->chr = '#'; // Marking this node a non-leaf node
95          temp->frequency = heap[1]->frequency;
96          temp->left = heap[1];
97          heap_delete(); // Pop one min-node
98          temp->frequency += heap[1]->frequency;
99          temp->right = heap[1];
100         heap_delete(); // Pop another min-node
101         // Since the minimal 2 nodes have been popped and connected to the temp node
102         // Then push the node back into the heap for the next loop
103         heap_insert(temp);
104     }
105     calculate_weight(heap[1],0); // Calculate the weight using huffman algorithm
106     int M, test_weight;
107     char test_char,test_string[MAX_CODING_LENGTH];
108     scanf("%d",&M);
109     // This loop will read M test cases, the answer will be given right away
110     for (int i = 0; i < M; i++) {
111         test_weight = 0;
112         for (int j = 0; j < N; j++) {
113             getchar();
114             scanf("%c %s",&test_char,test_string);
115             strcpy(test_table[j], test_string); // This makes checking prefix more quickly
116             test_weight += strlen(test_string) * frequency[(int)test_char]; // Calculate the
weight of this case
117         }
118         // After input, there are two test.
119
120         // If test result cost more, then it fails immediately.
121         if (test_weight != weight) printf("No\n");
122
123         // Otherwise, check if one string starts at the beginning of any other one
124         else if (check_prefix(N) == 0) printf("No\n");
125
126         // If this case is able to pass the previous two tests, then it is OK.
127         else printf("Yes\n");
128     }
129     return 0;
130 }
131
132 int check_prefix(int test_groups_num) {
133     char* str_result;
134     for (int i = 0; i < test_groups_num; i++) {
135         for (int j = 0; j < test_groups_num; j++) {
136             if (j == i) continue; // No need to test itself
137             str_result = strstr(test_table[i], test_table[j]); // Check1 : substring
138             if (str_result == NULL) continue; // If it is not its substring, then continue
139             if (test_table[i] == str_result) { // Check2 : prefix
140                 return 0; // If it starts from the same address as the other one, then it
doesn't obey prefix rule
141             }
142         }
143     }
144     return 1;
145 }
146
147 void calculate_weight(ptr_huffman_node node, int level) {
148     if (node == NULL) return;
149     calculate_weight(node->left,level + 1);
150     calculate_weight(node->right,level + 1);
151     if (node->chr != '#') weight += node->frequency * level; // Since I use '#' to mark it
non-leaf
152 }
153
154 void heap_delete() {
155     heap[1] = heap[node_num--]; // Move the last one to the top, node_num counts down

```

```

156 // This loop will check the property of min-heap from top
157 for (int i = 2; i <= node_num; i *= 2) {
158     if (i + 1 <= node_num && heap[i]->frequency > heap[i + 1]->frequency) i++;
159     if (heap[i / 2]->frequency > heap[i]->frequency) swap(i, i / 2);
160     else break;
161 }
162 }
163
164 void heap_insert(ptr_huffman_node node) {
165     heap[++node_num] = node; // Insert one node to the last, node_num counts up
166     int this = node_num;
167     // This loop will check the property of min-heap from bottom
168     while (this != 1 && heap[this]->frequency < heap[this/2]->frequency) {
169         swap(this, this/2);
170         this /= 2;
171     }
172 }
173
174 void swap(int smaller_index, int larger_index) {
175     // Simply swap two nodes
176     ptr_huffman_node temp;
177     temp = heap[smaller_index];
178     heap[smaller_index] = heap[larger_index];
179     heap[larger_index] = temp;
180 }

```

Appendix II Testing Set Generating Script in C++

```
1  /*
2  Test cases generator
3  Input: N & M
4  Output:out.txt
5  */
6
7  #include <bits/stdc++.h>
8  using namespace std;
9
10 #define MAXN 63
11 #define MAXM 1000
12
13 std::map<int, char> Code_dict;
14 char sym[MAXN];
15 int freq[MAXN];
16
17 int Binary(int i)
18 {
19     if (i < 2) return i;
20     else
21     {
22         return (Binary(i / 2) * 10 + i % 2);
23     }
24 }
25
26 class Node {
27 public:
28     char c;
29     int frequency;
30     Node* left;
31     Node* right;
32     Node(char _c, int f, Node* l = NULL, Node* r = NULL)
33         :c(_c), frequency(f), left(l), right(r) { }
34     bool operator<(const Node& node) const { // reload
35         return frequency > node.frequency;
36     }
37 };
38
39 void initNode(priority_queue<Node>& q, int N)
40 {
41     for (int i = 0; i < N; i++)
42     {
43         Node node(sym[i], freq[i]);
44         q.push(node);
45     }
46 }
47
48 //construct huffmanTree
49 void huffmanTree(priority_queue<Node>& q) {
50     while (q.size() != 1) {
51         Node* left = new Node(q.top()); q.pop();
52         Node* right = new Node(q.top()); q.pop();
53         Node node('R', left->frequency + right->frequency, left, right);
54         q.push(node);
55     }
56 }
57
58 // print huffmanCode
59 void huffmanCode(Node* root, string& prefix, map<char, string>& result) {
60     string m_prefix = prefix;
61     if (root->left == NULL) return;
62     //left subtree
63     prefix += "0";
64     if (root->left->left == NULL) result[root->left->c] = prefix;
65     else huffmanCode(root->left, prefix, result);
66     prefix = m_prefix;
67     //right subtree
68     prefix += "1";
69     if (root->right->right == NULL) result[root->right->c] = prefix;
70     else huffmanCode(root->right, prefix, result);
71 }
72
73 void testResult(map<char, string> result) {
74     //recur
75     map<char, string>::const_iterator it = result.begin();
76     while (it != result.end())
77     {
78         cout << it->first << ' ' << it->second << endl;
```

```

79         ++it;
80     }
81 }
82
83 int main()
84 {
85     priority_queue<Node> q;
86     int N, M;
87     //initialize
88     cout << "Please input the number of N and M:";
89     cin >> N >> M;
90
91     stringstream coutBuf = cout.rdbuf();
92     ofstream of("out.txt");
93     stringstream fileBuf = of.rdbuf();
94     cout.rdbuf(fileBuf);
95
96     cout << N << endl;
97     //i==key, c==Code_dict[i]
98     int c = 48;
99     int* Flag = (int*)malloc(sizeof(int) * 63);
100     for (int i = 0; i < 10 + 26 * 2 + 1; i++, c++)
101     {
102         if (i == 10) c = 65;
103         if (i == 36) c = 95;
104         if (i == 37) c = 97;
105         Code_dict[i] = c;
106         //0 48 A 65 _ 95 a 97
107         Flag[i] = 0;
108     }
109     //random char and frequency
110     srand((unsigned)(time(NULL)));
111     for (int i = 0; i < N; i++)
112     {
113         int a = rand() % 63;
114         if (!Flag[a])
115         {
116             sym[i] = Code_dict[a];
117             int b = rand() % 1000 + 1;
118             freq[i] = b;
119             cout << sym[i] << ' ' << freq[i] << ' ';
120             Flag[a] = 1;
121         }
122         else i--;
123     }
124     free(Flag);
125     cout << endl;
126     cout << M << endl;
127
128     //calculate the correct codes
129     initNode(q, N);
130     huffmanTree(q);
131
132     Node root = q.top();
133     string prefix = "";
134     map<char, string> result;
135     huffmanCode(&root, prefix, result);
136     //test result
137     for (int j = 1; j <= M / 2; j++) testResult(result);
138
139     //calculate bits
140     int num = 1, k = N;
141     while ((k = k / 2) != 0)
142     {
143         num++;
144     }
145     //print codes
146     for (int j = M / 2 + 1; j <= M; j++)
147     {
148         for (int i = 1; i <= N; i++)
149         {
150             int code = 0;
151             code = Binary(i);
152             cout << sym[i - 1] << " ";
153             if (j == M && i == N)
154             {
155                 cout.width(num);
156                 cout.fill('0');
157                 cout << code;
158             }
159             else
160             {
161                 cout.width(num);
162                 cout.fill('0');
163                 cout << code << endl;
164             }
165         }
166     }

```

```
166     }
167 }
168
169 of.flush();
170 of.close();
171 cout.rdbuf(coutBuf);
172 return 0;
173 }
174
```

Appendix III Declaration and Signatures

Declaration

We hereby declare that all the work done in this project titled "*Huffman Codes*" is of our independent effort as a group.

Signatures