

# Assignment 6

## Huffman Coding

By Jiancheng (Jason) Xiong  
CSE 13S Spring 2023

Due June 8th 2023, 11:59 PM

### 1 Purpose

The purpose of this assignment is to create a program that takes an input data file and compresses it using the Huffman algorithm. We will use a program huff that computes the huffman code of an input file and encodes an input data file. We will also have the capability to decode the file using a huffman decoder, which has been provided.

### 2 Using the Program

In order to use the program, compile and run the file huff.c. In order to run the program, run ./huff and specify the commands in which you want to run. The available options are;

- i (inputfile), specifying the input image on which the program will run on.
- o (outputfile), specifying the output file in which the image will be stored
- h; displays a help message detailing program usage.

### 3 Program Design

The main program will be in the file huff.c, while each of the separate data structures will be in;

- io.a for the reading and writing of various uint sizes, includes functions that open and close files unbuffered, along with header file io.h
- node.c for the creation and manipulation of the node ADT; includes functions that create, delete, and display nodes.
- pq.c for the creation and manipulation of the priority queue ADT; includes functions that create, delete, and modify entries of a priority queue.
- Makefile for the compilation and joining of the separate files

## 4 Data Structures

There are various implemented data structures, being;

The bitwriter struct is used to facilitate the writing of files. It stores a buffer it uses for storing text to be written, a uint8\_t byte, as well as a uint8\_t that saves the bit position.

The Node struct is used to save the properties of a node, such as the symbol of the node, the weight of the node, as well as the current code that represents the node, as well as the current length of the code. It also stores a pointer to its left and right children nodes.

The ListElement struct used to create a linked list ADT, storing a pointer to its node as well as a pointer to the next node in the list.

The PriorityQueue ADT that represents the actual queue. We use a second struct to represent the queue so we can always have a pointer to the queue.

The Code ADT that saves the code in the huffman tree that represents the current node, as well as an uint8\_t that represents the length of the code.

## 5 Algorithms/Pseudocode

### **For Bitwriter Functions:**

Bit\_write\_open(filename)

- Calloc some memory for the bitwrite structure

- Create a buffer using write open

- Set the underlying stream in the bitwrite structure to the buffer made using write open

- Return a pointer to the new bitwrite structure

Bit\_write\_close(filename)

- If there is data left in the buffer: flush bytes using write\_uint8 until it is empty

- write\_close() the underlying stream buffer

- Free the bitwriter memory

Set its pointer to null

Bit\_write\_bit(buffer, x)

If there is an entire byte in the buffer, write the byte to underlying stream using write\_uint8

Set the bit\_position pointer to 0

Leftshift bit by byte if necessary

Increment bit position

Bit\_write\_uint8(buffer, x)

Call bit\_write\_bit 8 times, starting from the right bit

Bit\_write\_uint16(buffer, x)

Call bit\_write\_bit 16 times, starting from the right bit

Bit\_write\_uint32(buffer, x)

Call bit\_write\_bit 32 times, starting from the right bit

### **For Node Functions:**

Node\_create(symbol, weight)

Calloc a node

Set node->symbol to symbol

Set node->weight to weight

Set node->left and node->right to NULL

Return a pointer to the node

Node\_free(node)

Free the node

Set the pointer to the node to NULL

### **For Priority Queue Functions:**

pq\_create(void)

Calloc some memory for the priority queue structure

Return a pointer to the priority queue object

pq\_free(priority\_queue)

free(priority\_queue)

Set the pointer to the priority\_queue to NULL

pq\_is\_empty(priority\_queue)

Check if the queue's list field is NULL, if so, return true otherwise return false

pq\_size\_is\_1(priority\_queue)

Check the first element in the linked list field

If the first element-> is not NULL (ie there is another element in the linked list), return false, otherwise return true

enqueue(priority\_queue, tree)

Allocate memory for a new Listelement e, and set e->tree to tree

Check the current queue's state;

If the queue is empty, set q->list to e

If the weight of the new tree is less than the weight of the head, add it to the beginning of the queue by setting e->next to q->list and setting q->list to e

If the weight of the new tree is greater than the weight of the head, start at the head of the tree, and traverse the linked list until you find an element with a weight greater than the weight of the new tree, then add it to the priority queue before that element

dequeue(priority\_queue, tree)

Save the first element from priority\_queue into tree, update the pointers, and free the first node

pq\_less\_than(node1, node2)

If node1->weight is less than node2->weight return true

If node1->weight is greater than node2->weight return false

Otherwise return if node1->symbol < node2->symbol

## **For Huffman Encoding**

fill\_histogram(inbuf, histogram)

Initialize filesize to 0

- Iterate through the inbuf and add the read bytes to the histogram, incrementing filesize while doing so
- Add 0x00 and 0xff to histogram as a hack
- Return filesize

Create\_tree(histogram, num\_leaves)

- Create a priority queue pq
- Initialize num\_leaves to 0
- Create a node for every symbol that appears (has a frequency > 0), setting its symbol and adding it to the priority queue

- While the size of the pq is greater than 1:
  - Dequeue a node into left
  - Dequeue a node into right
  - Make a parent node with the weights of both children added together
  - Set parent node->left to left
  - Set parent node->right to right
  - Enqueue the parent node

- Return the remaining node

Fill\_code\_table(code\_table, node, code, code\_length)

- For each child of a node (if it exists)
  - Recursively call this function on that child with a code\_length + 1
- Else:
  - Set the code and code\_length of node in the code\_table

huff\_compress\_file(outbuf, inbuf, filesize, num\_leaves, code\_tree, code\_table)

- Write uint8\_t 'H'
- Write uint8\_t 'C'
- Write uint32\_t filesize
- Write uint16\_t num\_leaves
- Huff\_write\_tree(outbuf, code\_tree)
- For every byte b from inbuff
  - Look up the code and code\_length
  - For i from 0 to code\_length - 1

Write the rightmost bit  
Code  $\gg= 1$

## 6 Function Descriptions

### Bit\_write\_open

Takes in const char \*filename

Returns a BitWriter struct

Opens a file, creates and returns a pointer to a BitWriters structure

### Bit\_write\_close

Takes in a Bitwriter \*\*pbuf

Returns nothing

Flushes data in the byte buffer, frees the buffer, and sets pointer to NULL

### Bit\_write\_bit

Takes in a Bitwriter \*buf and uint8\_t x

Returns nothing

Write a single bit using values in the BitWriter struct pointed to by buf, if there are 8 bits in the buffer

### Bit\_write\_uint8, 16, and 32

takes in a Bitwriter \*buf and uint8/16/32\_t x

Returns nothing

Writes 8/16/32 bits using bit\_write\_bit

### Node\_create

Takes in a uint8\_t symbol and double weight

Returns a Node

Creates a node and returns a pointer to it

### Node\_free

Takes in a Node \*\*node

Returns nothing

Frees the \*node and sets it to NULL

### Pq\_create

Takes in nothing  
Returns a PriorityQueue  
Allocates a PriorityQueue object and returns a pointer to it

Pq\_free

Takes in a PriorityQueue \*\*q  
Returns nothing  
Frees q and sets pointer to null

Pq\_is\_empty

Takes in a PriorityQueue \*q  
Returns a bool  
Checks if pq is empty

Pq\_size\_is\_1

Takes in a PriorityQueue \*q  
Returns a bool  
Checks if pq has only one element

Enqueue

Takes in a PriorityQueue \*q and a Node \*tree  
Returns nothing  
Adds the tree to the priority queue in the correct position

Dequeue

Takes in a PriorityQueue \*\*q and a Node \*\*tree  
Returns bool  
Removes the first element, rearranges the pointers, and returns true

Pq\_less\_than

Takes in two nodes Node \*n1 and Node \*n2  
Returns bool  
Returns if n1 is < n2, if they're the same weight the ascii of their symbol is used as a tiebreaker

Fill\_histogram

Takes in a Buffer \*inbuf and a double \*histogram  
Returns a uint64\_t

Create\_tree

Takes in a double \*histogram and a uint16\_t \*num\_leaves

Returns a Code

Creates and returns a pointer to a new huffman tree

Fill\_code\_table

Takes in Code \*code\_table, Node \*node, uint64\_t code, uint8\_t code\_length

Returns nothing

Fills a code table for the huffman code

Huff\_compress\_file

Takes in BitWriter \*outbuf, Buffer \*inbuf, uint32\_t filesize, uint16\_t num\_leaves, Node \*code\_tree, Code \*code\_table

Returns nothing

Compresses a file using huffman encoding

## 7 Results

Pictures of the program running:

## 8 References

“Getopt() Function in C to Parse Command Line Arguments.”

*GeeksforGeeks*, 10 Sept. 2018,

[www.geeksforgeeks.org/getopt-function-in-c-to-parse-command-line-arguments/](http://www.geeksforgeeks.org/getopt-function-in-c-to-parse-command-line-arguments/).

Open(2) - Linux Manual Page, 27 Aug. 2021,

[man7.org/linux/man-pages/man2/open.2.html](http://man7.org/linux/man-pages/man2/open.2.html).

“BMP File Format.” Wikipedia, 16 May 2023,

[en.wikipedia.org/wiki/BMP\\_file\\_format](https://en.wikipedia.org/wiki/BMP_file_format).