

CSE 13S Spring 2021
Assignment 5: Huffman Coding
Design Document

Encoder: Find Huffman encodings if an input file, use those encodings to compress the file.

Accept command line options:

- h - print out program info message
- i - specify an infile other than stdin
- o - specify an outfile for the encoded codes only other than stdout.
- v - print compression statistics (uncompressed size, compressed size, space saving)

to stderr

$$100 \times (1 - (\text{compressed} / \text{uncompressed})) = \text{saving}$$

Encoder Flow:

- 1) Parse command line args with getopt
- 2) Construct a histogram of 256 uint64ts
- 3) Increment the count of first and last histogram elements (0 and 255)
This is to eliminate the errors associated with empty or single character infile.
- 4) Construct Huffman tree using a priority queue (build_tree)
Create a priority queue with a Node for each symbol in the histogram with a positive

frequency

While there are two or more items in the priority queue, dequeue two of them, combine them to have a parent node

Parent node frequency = left + right child frequency

The last node standing is the root.

- 5) Construct a code table, this is a simple 256 code array. (Build_codes)

Starting at root traverse Huffman tree

If current node is leaf, the code is the path from root to node

Save the code to the code table

Else the current node must be a branch, push 0 to the code and recursively apply the process to the left link.

Pop a bit from the code, as with depth first search, then recurse down the right node. (Adding 1 to the code).

Pop the most recent bit from the code.

- 6) Construct a header (ADT for tracking information that must be transferred from the input to the output.

- 7) Perform a post order traversal of the Huffman tree.

Write L followed by the byte symbol for each leaf

Write I followed by nothing for interior nodes.

- 8) For every item in the infile, use your code to write the corresponding code to the outfile. (Write_code)

Flush remaining codes *to prevent data leaks?*

Decoder: Decode the Hamming codes made by the encoder.

Print statistics (bytes processed, uncorrected and corrected errors, error rate decimal to 6 digits of precision) to stderr.

Accept command line options:

- h - print out program info message
- i - specify an infile other than stdin
- o - specify an outfile for the decoded codes other than stdout.
- v - enables printing statistics to stderr, same stats as encode

Decoder Flow:

- 1) Parse command line options with getopt
- 2) If input magic number does not match macro magic number, end program
- 3) Read the input's Huffman tree
 - Store the tree in an array of size tree_size
 - Reconstruct the tree using rebuild_tree
 - Read over the tree array (tree_dump)
 - If element of the array is L, the next element is a leaf symbol, create a node for that leaf with node_create
 - Push the node to your stack
 - If element of the array is I, pop once for RIGHT child, again for LEFT child
 - Join children to make a parent node, push this onto the stack
 - Last element is the root of the Huffman tree.
- 4) Read the infile one bit at a time (read_bit)
 - If bit = 0, walk down left child of current tree node
 - If bit = 1, walk down right child of current tree node
 - If you hit a leaf node, write the leaf's symbol to the outfile and return back to the root to repeat.
- 5) Outfile is complete when it's length matches file size passed by the infile.

The Node ADT:

These are the base units Huffman trees are made of.

Nodes contain these variables:

Pointer to left child (Node pointer)

Pointer to right child (Node pointer)

It's own symbol (uint8)

The frequency of the symbol (used by encoder) (uint64)

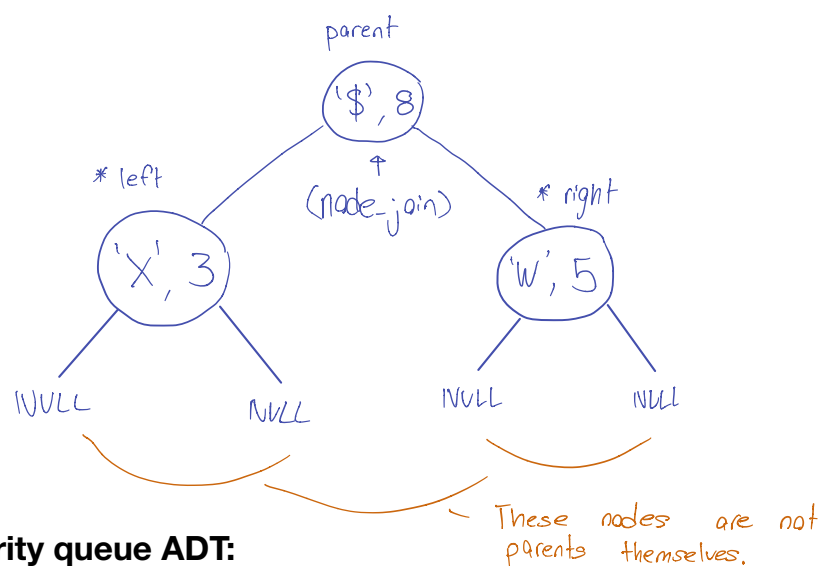
Functions to implement:

node_create - create new node

node_delete - destructor, free the pointer

node_join - join a LEFT child to a RIGHT child, return a pointer to the created parent with the children set accordingly. Parent symbol = "\$". Frequency = left frequency + right frequency

node_print - debug print function.



The priority queue ADT:

This priority queue stores nodes.

Like a queue but instead of first in first out, it is highest priority first out.

Priority is indicated by the frequency of a node (lower frequency = higher priority).

Priority queue functions:

pq_create - create a priority queue with a maximum size specified by input *capacity*

pq_delete - destructor, free the pointer

pq_empty - true if empty, false otherwise

pq_full - true if full, false otherwise

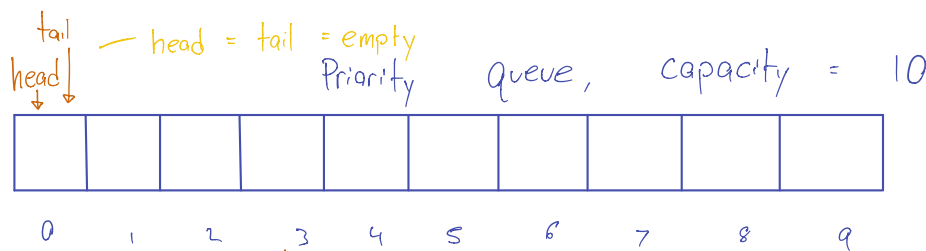
pq_size - return size

enqueue - return false if queue is full before enqueueing

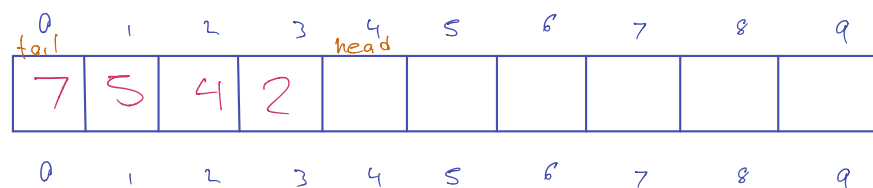
dequeue - return false if queue is full before dequeuing

pq_print - display a tree of your ordered priority queue, *works best if queue is ordered by*

insertion sort



- enqueue at tail
- dequeue at head



- degree what is at head

↳ Since tail always @ 0.

- top == bottom
AND byte(bottom) = empty

Pseudocode for code stack functions:

```
code_full()
    if code->top == code_length AND code->items(top) == 256 - aka, if all bits are full
        return true
    else return false

code_empty()
    if code->top == code_bottom AND code->items(bottom) == 0
        return true
    else return false

code_push()
    if code_full == false
        if code->items(top) == 256 - aka current byte is full
            top++
            code->items(top) << - aka shift all bits one to the left
            code-> items(top) ++ - add 1, so that all bits (with one additional one) are now set.
        else return false

code_pop()
    if code_empty == false
        if code->items(top) == 0 - aka current byte is already empty
            top- -
            code->items(top) >> - shift right (you could also just do integer division by 2)
        else return false
```

The Stack of Nodes ADT:

Same stack as from asgn3, 4 however, instead of storing an array of items, this stack stores nodes of items.

Effected stack functions:

- stack_push - push a node onto the stack, or return false if stack is full
- stack_pop - pop a node from the stack to the pointer, or retur false if stack is full

The Header ADT:

Header is simply a storage structure which contains information (metadata?) on the infile, the information on the Huffman tree the decoder will need to replicate, and the file size to aim for when decoding.

I/O functions:

- read_bytes - read bytes from an infile into a buffer, stop when the infile is parsed OR you have read the specified number of bytes - using a loop

- write_bytes - write a specified number of bytes (or all bytes) from a buffer into the outfile using a loop

- read_bit - places bytes into a static buffer of bytes - using an index to track the bit you are on, read every individual bit of the bytes in the buffer - refilling the buffer when it has been read until your infile is finished.

- write_code - set single bits from the code ADT to the buffer, when the buffer is full, send the contents to the outfile

- flush_codes - A function to clear bits remaining in the buffer after encoding is complete, *clear bits in last input byte before flushing***

Huffman.h functions:

`build_tree` - given a computed histogram construct a Huffman tree.

Histogram will have 256 indices, one for each symbol

Return the root node for your Huffman tree

build_codes - fill the code table with code for each of the symbols in the Huffman tree

constructed codes are copied to a CODE ADT table with 256 bitwise indices

rebuild_tree - using the tree_dump array, reconstruct the same Huffman tree.

Return root node for the reconstructed tree

`delete_tree` - destructor for a tree, this must free ALL nodes and set tree pointer to Null.

Example of the encoding/tree making process:

First create a histogram: Word to encode:

Doppelkupplungsgetriebe $23 = \text{length}$
14 different elements

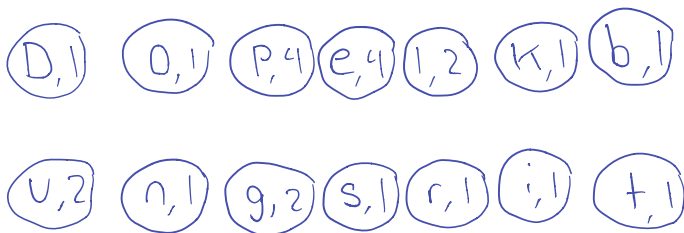
Histogram :

DOPELX	UNGSrib +	2.121111
1	1	1
1	4	2
4	4	1
2	1	1
1	b	1
	+	1

Priority queue for this histogram



Nodes for this histogram :



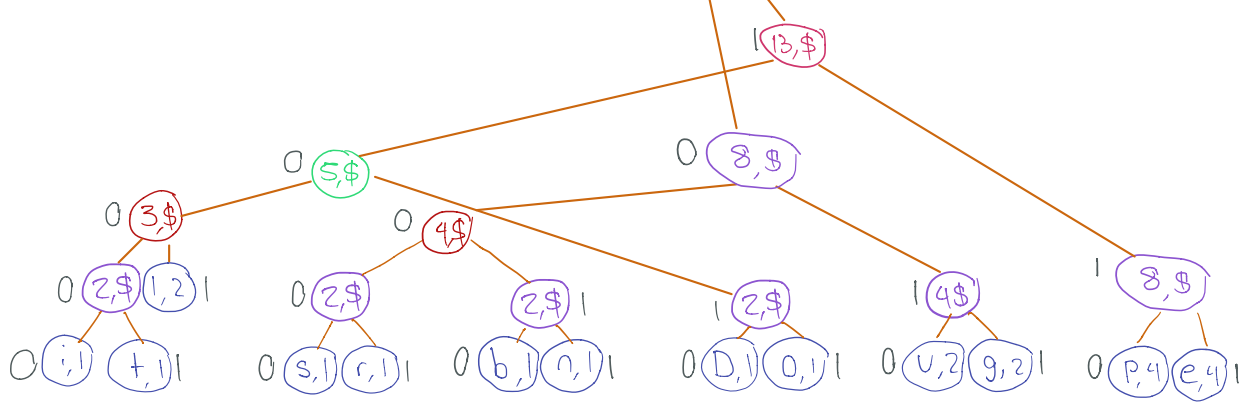
Constructing the Huffman Tree:

PQ



Dequeue 2, enqueue parent - until one node left





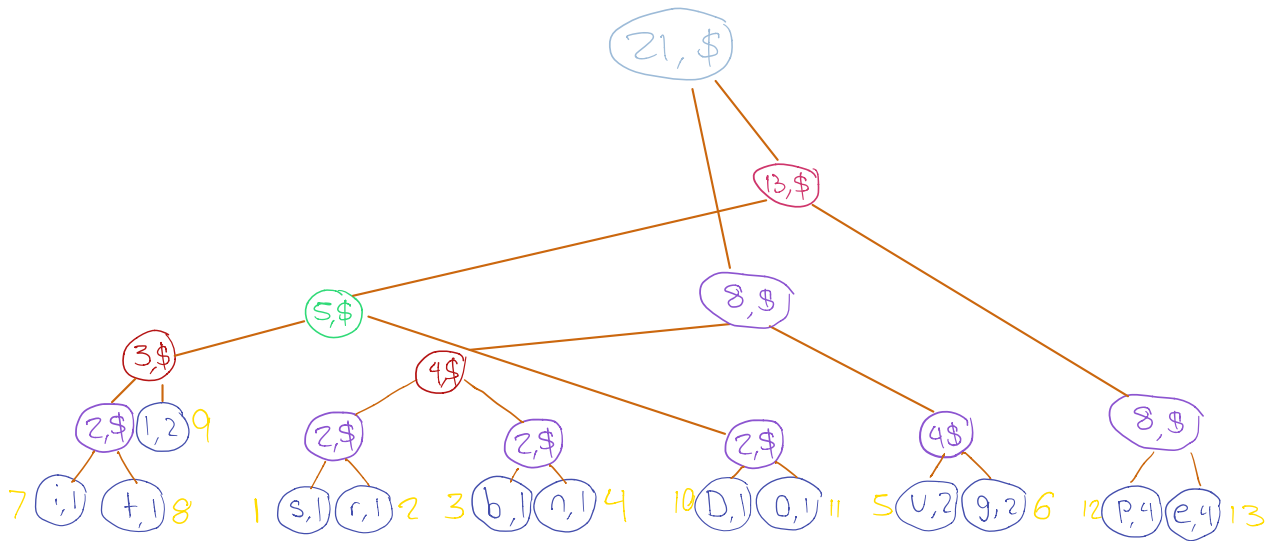
Constructing a code table

Left = 0, Right = 1

D	1010
O	1011
P	110
e	111
L	1001
+	10001

u	010
n	0011
g	011
s	0000
r	0001
i	10000
b	0010

Constructing a Tree Dump Array



Tree Dump :

[Ls, Lr, I, Lb, Ln, I, I, Lu, Lg, I, I, Li, Lt, I, Ll, I, Ld, Lo, I, I, Lp, Le, I, I, I]

Reconstructing a tree from the tree dump array (using a stack)

Stack for reconstruction:

↓

Reconstructed Tree :

(\$)

