

# Huffman Encoding Text file Compression

what is node?



contains: a character  
its frequency  
and a left and right  
pointer to other  
child node.

what is a pointer: A value which  
refers to address of an object

eg: a is a variable  
&a is address of a

then p be pointer of a  
ie,  $p = \&a$

$\therefore *p = a$  ; \* : dereference  
operator  
ie,  $*\&a = a$

what is struct min-heap?

Contains: size; no. of elements in heap  
(unique characters)

Node\*\* array; here array is a pointer  
which points to other node pointers

ie, array is a pointer which  
points to a set of pointers of  
type Node\*

eg: min-heap: size = 5  
array = 3200

3200 have 5 x nodes  
ie, \*3200 have 

a
b
c
d
e

 \* array[0]  
and a, b, c, d, e  
are pointers (Node\*) to each node  
of unique character.

Before we continue,

consider a text file, open using FILE\* and a  
pointer to beginning of file is provided  
It's run along each character to obtain frequency

count-freq  $\rightarrow$  freq[i]++ for each time char i read

ie, if 'a' is read the ASCII of a will be 31  
once then freq[31]++

$\therefore$  finally freq[0] = 0  
[1] = 0  
:  
[31] = 1  
:  
[255] = 0

no. of counts

[freq[ASCII index value]]

to form character array and frequency array of only present characters

count = 0;  
i runs from 0  $\rightarrow$  255 if freq[i] > 0 append freq[count] = freq[i];  
char[count] = (char)i;  
count++

$\therefore$  if a: 5 times, b: 3 times, d: 1 time

then freq = { a, b, d }

characters = { 5, 3, 1 }

Now min-heap to form nodes and allocate all nodes in a single  
memory location. (pointer 3136936 location - 1e)

struct minheap\* createheap (freq, characters, count)

struct minheap\* minheap = (struct minheap\*) malloc (sizeof (struct minheap));

type of value  
returned

ie pointer to a minheap struct.

new min-heap kind  
size 200 location 200000  
and pointer of location



min-heap  $\rightarrow$  size = cost;

indicate min-heap is a pointer of struct

(if min-heap was struct itself (can use min-heap-size) the same functionality is)

min-heap  $\rightarrow$  array = (struct node\*\*) malloc (min-heap  $\rightarrow$  size \* sizeof (struct node\*));

now, we read each character [i] and

return  $\rightarrow$  array,

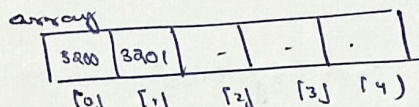
cost \* sizeof pointer of node  
 $\therefore$  5 x pointer  $\rightarrow$  means  
 small pointer

node\* new node (char [i], freq [i]) is stored at  $\rightarrow$  array [i]

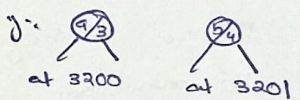
```

struct node* new node (char [i], freq [i]) {
    struct node* node = (struct node*) malloc (size of struct node);
    node  $\rightarrow$  character = char [i]; node  $\rightarrow$  l = NULL
    node  $\rightarrow$  freq = freq [i]; node  $\rightarrow$  r = NULL
    return node;
}
  
```

Now we have ; min heap: size = 5



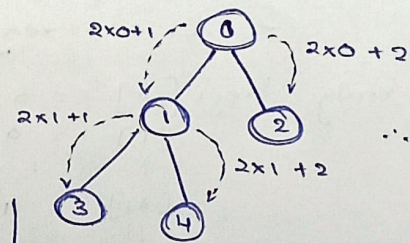
contains pointers to each node of edge character.



we need to re-arrange this in an array order in order to extract min freq node.

we need to arrange array now ready to freq order.

Note that: binary tree order is given as



Parent = i  
 $\therefore$  l:  $2i+1$   
 r:  $2i+2$

void Heapify (min-heap, i)

smallest = i  
 l =  $2i+1$   
 r =  $2i+2$

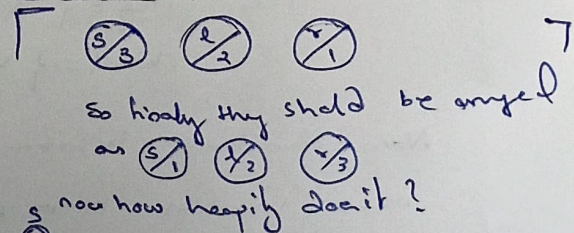
Start with  $i = \frac{\text{cost}-1}{2}$

then i-- after each heapify. the heap is going to be all nodes and heapify recursively applied.

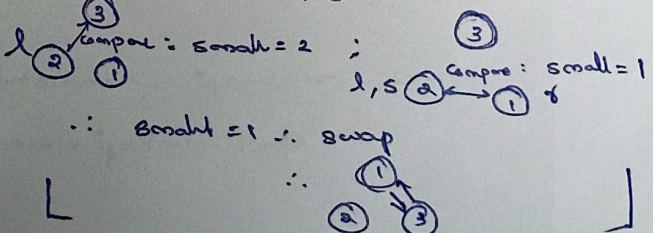
if (min-heap  $\rightarrow$  array [l]  $\rightarrow$  freq < min-heap  $\rightarrow$  array [smallest]  $\rightarrow$  freq)  
 smallest = l;

if (min-heap  $\rightarrow$  array [r]  $\rightarrow$  freq < min-heap  $\rightarrow$  array [smallest]  $\rightarrow$  freq)  
 smallest = r;

if (smallest  $\neq$  i)  $\leftarrow$  means smallest freq char was not at index i is swap array [smallest] and array [i]



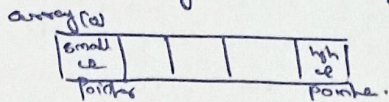
now how heapify does it?



recursively heapify (min-heap, smallest)



Now we got min heap: size = 5



Now we need to build tree

struct \*node tree (min heap) {

node \*l, \*r, \*top  
 while (!isSizeOne (min-heap)) {  
 smallest array[10] → l = extract min (min-heap);  
 next smallest array[11] → r = extract min (min-heap);

top = new node (\$, l → l → freq + r → r → freq)

min heap array character

frequency sum

top → l = l

top → r = r

insert into heap (min-heap, top)

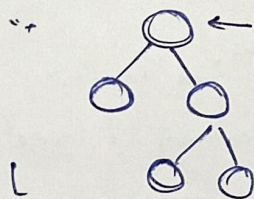
Basically we extracted 2 min node, added them up to form a new node and placed

it back in heap and do this repeatedly

until only one node left in min heap

and its pointer is returned to tree

if its pointer is tree.



[

extract\_min (min-heap)

return min heap → array[10]

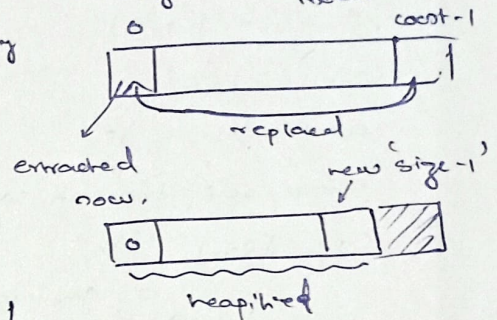
and reach min heap → array[10] with

min heap array [size - 1]

size --

and heapify it. i.e. re-align frequency

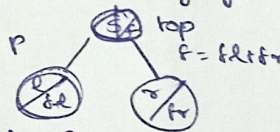
visually:



insert into heap (min-heap, top)

we need to find proper index of character to place

it according to its frequency or node top



how to find index?

increase count to fit new node.

let i = count - 1; i.e. last blank position

while (top → freq < array[i]) → freq

find swap position i.e.

top at i-1 and i

place

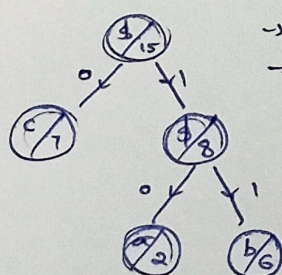
min heap → array[i] = min heap → array[i-1]

i = (i-1)/2

min heap → array[i] = top

ie, we place array with frequency greater than top to the bottom position until we have appropriate index for placing top

So now we have to generate code words for each character.

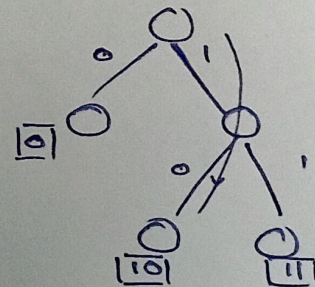


char code intermittently stores code

char codes [ ] [ ]

ASCII of character at leaf node i.e. end node

∴ codes have code word or index as ASCII number of each character.





```
void generateCode (tree, code, codes) {
```

```
    if (!l and !r then to add to code
        and code given to codes[ch][ ]
```

```
    if (tree -> l)
```

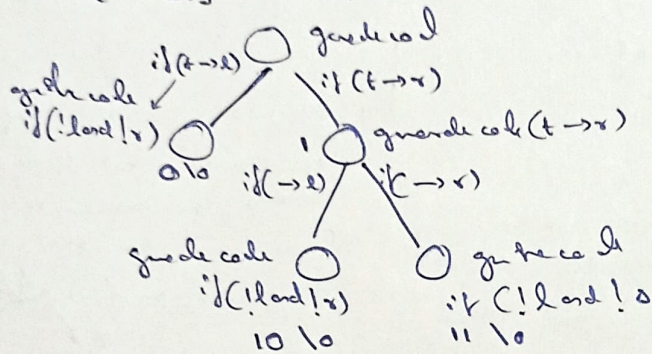
```
        len = strlen(code)
```

```
        code[len] = '0'
```

```
        code[len+1] = '\0'
```

```
        generateCode (tree -> l, code, codes)
```

```
        code[len] = '\0'
```



note: if (tree -> l) is always if (1) i.e., it will run for sure cause tree -> l is a pointer value of a new node unless it is an end node and value > 0 are all considered YES. ∴ if (1) means sure run. it will traverse along all nodes recursively.

Adding '\0' to next index doesn't really do any good in b/w as it's not considered as a real character in early replace while recursively moving along the tree.

Now just read i/p file and obtain print codes[ch] in o/p file. int ch = ASCII value of

our binary tree is already provided to decode

Decompression:

```
void decompress (tree, o/p file)
```

As pointer of file move along

```
it read bit = 0 ; curnt = curnt -> l
```

```
bit = 1 ; curnt = curnt -> r
```

```
if (!curnt -> l and !curnt -> r)
```

print (curnt -> character) to decompress file.

and reset pointer, curnt = tree

Initially curnt = tree i.e. and we replace curnt\* as we traverse

