

Date: 02/05/21

# Black Board

Design and Analysis of Algorithms

Topics:

- Greedy Algorithms II
  - Huffman Encoding ✓

# Fixed Length Encoding

ASCII

A — 65 → 01000011

B — 66 → ...

8-bit

All codes have the same length.

# File Compression using **fixed length encoding**

- Example: a file with 6 unique characters

	a ✓	b ✓	c ✓	, d ✓	e ✓ <sup>2<sup>8</sup></sup>	f ✓
Frequency (in thousands)	45	13	12	16	9	5
Fixed-length codeword	000	001	010	011	100	101

100,000 characters  $\Rightarrow$  800,000 bits.

aacdeb...

000000010...

$$\begin{array}{c} \overline{\quad} \quad \overline{\quad} \quad \overline{\quad} \\ 2 \times 2 \times 2 \\ = 8 \end{array}$$

# Compression Ratio

	a	b	c	d	e	f
Frequency (in thousands)	45	13	12	16	9	5
Fixed-length codeword	<u>000</u>	<u>001</u>	010	011	100	101

Compression Ratio:  $r = 1 - \frac{\text{Size of compressed file}}{\text{Size of original file}}$

AGCT

$$= 1 -$$

$$\frac{300,000 \text{ bits}}{800,000 \text{ bits}}$$

$$= 1 - \frac{3}{8} = 0.625$$

62.5%

# Variable Length Encoding

	a	b	c	d	e	f
Frequency (in thousands)	45	13	12	16	9	5
Fixed-length codeword	000	001	010	011	100	101
Variable-length codeword	<u>0</u>	<u>101</u>	<u>100</u>	<u>111</u>	<u>1101</u>	<u>1100</u>

Compression Ratio:

$\frac{a}{0}, \frac{b}{001}$

abb...

000/001...

$$r_2 = 1 - \frac{1000 \times (45.1 + 13.3 + 12.3 + 16.3 + 9.4 + 5.4)}{800,000 \text{ bits}}$$

$$= 1 - \frac{224,000}{800,000} = 72\%$$

$a=0, \frac{C=100}{b=001}$

consider

00100

aac  
brr

Encoding must be **Prefix-free**  
for unambiguous decoding

✓  
→ Variable length encoding  
(based on "shorter codes for more frequent  
characters.")

+  
Prefix free

# The Huffman Encoding Algorithm

$$|F| = n$$

- ✓ Finds a prefix-free, variable length encoding for an input file.  $O(n)$
- The encoding is optimal, i.e. achieves the highest compression ratio over all prefix-free, variable length encodings
  - (Proof of optimality uses induction on the file size)



# Iterations of the Huffman Algorithm

Step 1: Compute a frequency table and create “PQT entries” for each unique original character.

$H = n$

$O(n)$

Each PQT or “Priority Queue” is a **struct type tree and priority queue node** containing:

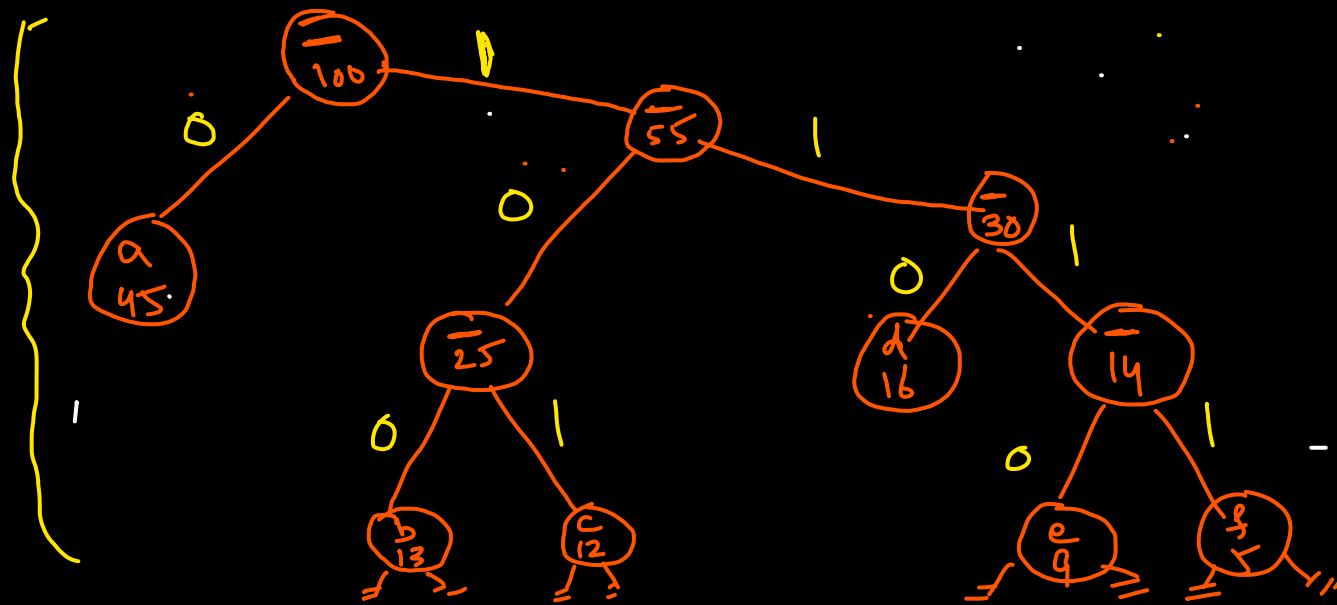
- key (used to maintain the min-heap)
- Sym: (symbol of the character, either original or -)
- lchild, rchild pointers (tree)
- Necessary comparison operators

convert  
Min Heap | Binary Tree



- Step 2: Build Priority Queue of the PQTs.

	a	b	c	d	e	f
Frequency (in thousands)	45	13	12	16	9	5
Fixed-length codeword	000	001	010	011	100	101
Variable-length codeword	0	101	100	111	1101	1100



- Step 3: Build Compression Tree from the Priority Queue in  $n$  iterations, where  $n$  is the number of unique symbols in the file.

→ getMin(), getMin()  
extractMin(), extractMin()

→ Make a parent node with  
accumulative frequency.

→ Insert parent node into the heap.

---







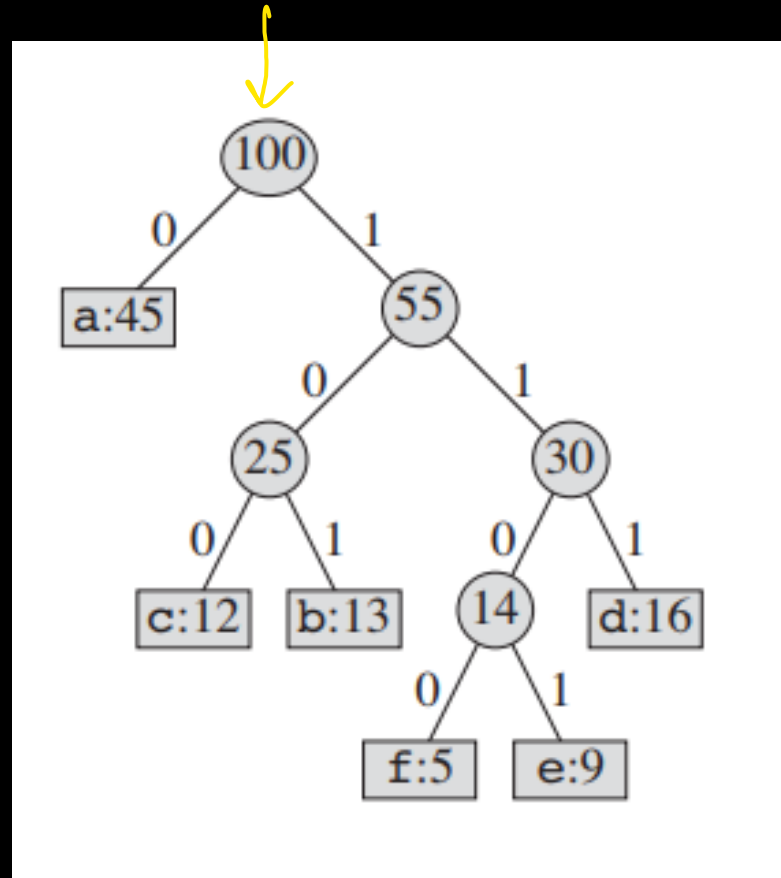
- Step 4: Compress the file using the compression tree: compression table, bit level packing, etc.



# Decoding

001001100...

aacf.....



$|F|=n$   
 $n \rightarrow$  size of file  
 $N \rightarrow$  # of unique  
 signs.

$T \rightarrow$  tree  
 $C(i) =$  begin  
 of the code  
 of charact  $i$ .

What is the running time  
 & Decompression program.

of  
 Compression