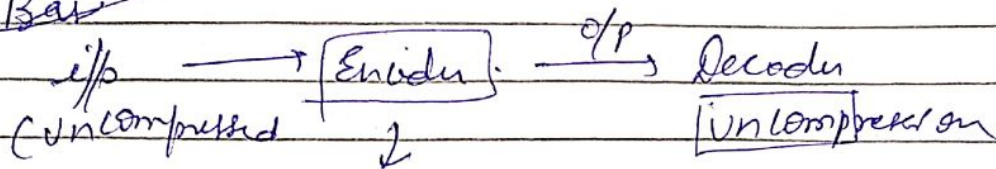# Data Compression

- process of reducing the amount of data required to represent a given quantity of info.

→ 2 GB — 750 mb — Compressed.
→ time not → 2 hrs — into 1 hr.
→ amount of disk space required.

## Huffman code Data Compression -

Bar

i/p ———→ [Encode] —o/p→ Decode
(uncompressed)        ↓            [uncompression]

Data Compression

                                                    Code Word
                                                    length

Lossless                    Lossy.        entropy.
    ↓                          ↓
Text Encryption          Some info is lost
    ↓                          ↓
exact replica           less imp. info from
                        the media is removed.
                                ↓
Compression ratio = Compressed file / Original file

= 750/2000                      Image, Video, Audio
                                        ↓
Compress factor = 1 / Compression ratio      not exact replica
(integer)

Compression time      Time to compress (ms)

Decompression time    Time to Decompress (ms)

Fixed length Encoding

a - 97 → 7 bits (1100001)

b
c
:
g

7X7 = 49 bits

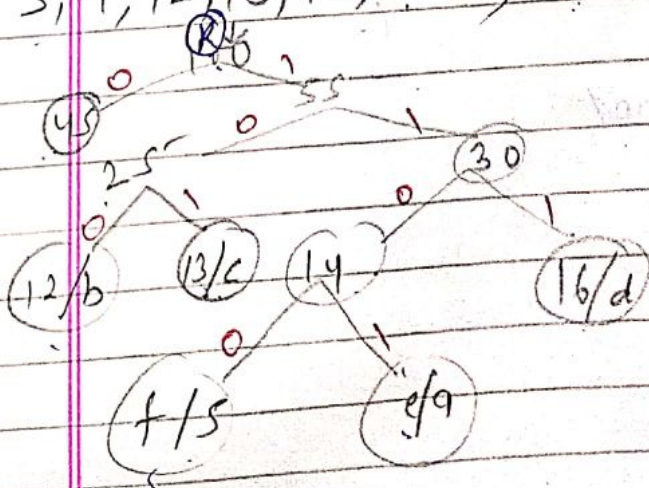Lossless Compression Technique -
(no loss of info)

Based on binary tree freq-sorting method that allow to encode
any msg sequence into shorter encoded msg.
Huffman code is a particular type of optimal
prefex code. o/p can be viewed as
variable - length code.

1. Create Sorted nodes on frequency/probability
2. Start
3. find & remove two smallest probability node
4. Create new node, weight[node] = W(A) + W(B)
5. Insert new node, back to sorted list.
6. Repeat the loop until the list consist of the
   only last node

| Char | frequency | fixed length code | Var. length code | Code length | |
|---|---|---|---|---|---|
| a | 45 | 000 | 0 | 1 | highest Compression |
| b | 13 | 001 | 101 | 3 | E Var logn |
| c | 12 | 010 | 100 | 3 | |
| d | 16 | 011 | 1111 | 4 | |
| e | 9 | 100 | 1101 | 4 | |
| f | 5 | 101 | 1100 | 4 | |

5, 9, 12, 13, 16, 45 , 12, 13, 14, 16, 45 , 14, 16, 25, 45
                                                          25, 30, 45

*lossless*

# Run length Encoding — Simplest.

In this runs of data are stored as a single data value, and count rather than original run.
↓ Sequence of same symbol / data value

AAAAAAAA ———→ A ⑧ count
                        data value

e.g.

$\underset{9}{\underline{BBBBBBBBB}}$ $\underset{5}{\underline{AAAAA}}$ $\underset{1}{N}$ $\underset{2}{44}$ $\underset{3}{\underline{mmm}}$ ———→ B09A05
                                                                  N01602M03

                        ㉀                                      ⑮

but for 0 & 1

$\underset{14}{\underline{00000000000000}}$ | $\underset{4}{\underline{0000}}$ | $\underset{0}{}$ | $\underset{12}{\underline{000000000000}}$

        ↓                                      ↓   ↓              ↘
       111 0                                 0100  0000          1100

**Drawback** — $\underset{3}{\underline{xyz}}$ ——→ $\underset{6}{\underline{x01 y01 z01}}$

When repeating value not .

# Arithmatic coding
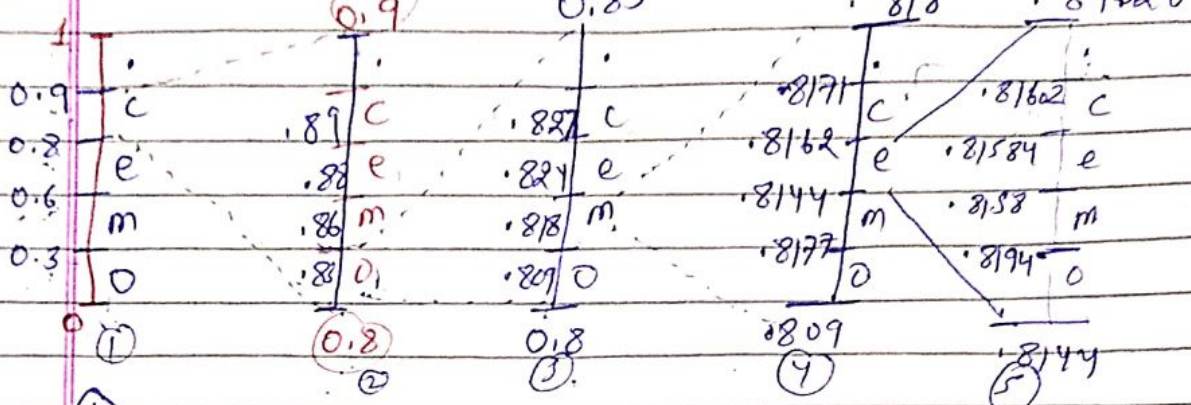
- form of entropy encoding used in lossless data.

- ✓ Efficient code — not depend on size #
- ✓ Non block code — block of data on the probability of data
  <u>Stream based Encoding</u>

- ✓ Work very well for sequence = low entropy.

$0 = 0.3 \qquad m = 0.3 \qquad e = 0.2 \qquad c = 0.1 \qquad (\therefore = \cdot 1)$

CODE



① 

② low limit + diff (Prob.)

$$= 0.8 + 0.1 (0.3)$$

$0 = = 0.8 + 0.3 = \cdot 83 = 0$

$m = \cdot 83 + 0.1 (0.3) - \cdot 86$

$e = \cdot 86 + 0.1 (0.2) = \cdot 88$

$c = \cdot 88 + \cdot 1 (0.1) = \cdot 89$

③ $0 = \cdot 8 + \cdot 03 (0.3) - \cancel{\cdot 89} \; \cdot 8 + \cdot 009 = \cdot 809$

$m = \cdot 809 + \cdot 03 (0.3) = \cdot 818$

$e = \cdot 818 + \cdot 03 (0.2) = \cdot 824$

$c = \qquad \qquad = \cdot 827$

⑤

$\cdot 81602 < \text{Codeword Range} < \cdot 81620$

$$\text{Generation tag} = \frac{L \cdot L + U \cdot L}{2} = \frac{\cdot 81602 + \cdot 81620}{2}$$

$$= \frac{\cdot 81611}{}$$

## Decoding

Codeword = 0.572

Symbol    |     E    C    E
prob      |    .1   .4   .5

### Step 1



based on probability

diff. d = upper bound − lower bound = 0.9 − 0.4 = 0.5

Range of symbol = lower limit; lower limit + d (prob. of symbol)

Range of C = 0.4; 0.4 + 0.5(0.4) = 0.6

E = ~.85

.572. in between. 0.57 & 0.58

$\underbrace{\qquad}$ range of termination

### ECE!

finite precision Arithmetic