

## Module-II

①

### Syllabus

{ Basic Compression Technique - Run length encoding,  
RLE Text compression, statistical Methods - Prefix  
coders, Binary Huffman Coding, Non-binary  
Huffman Algorithms, Arithmetic coding.

### Run length Encoding (RLE)

Run-length Encoding (RLE) is a form of lossless data compression in which runs of data (sequences in which the same data value occurs in many consecutive data elements) are stored as a single data value and count, rather than as the original run.

If a data item 'd' occurs 'n'

Consecutive times in the input stream,

then replace the  $n$  occurrences with  
 with the - single pair  $\text{ind}$ . The  
 $n$  consecutive occurrences of a data  
 item are called a runlength of  $n$ .

### RLE Text Compression

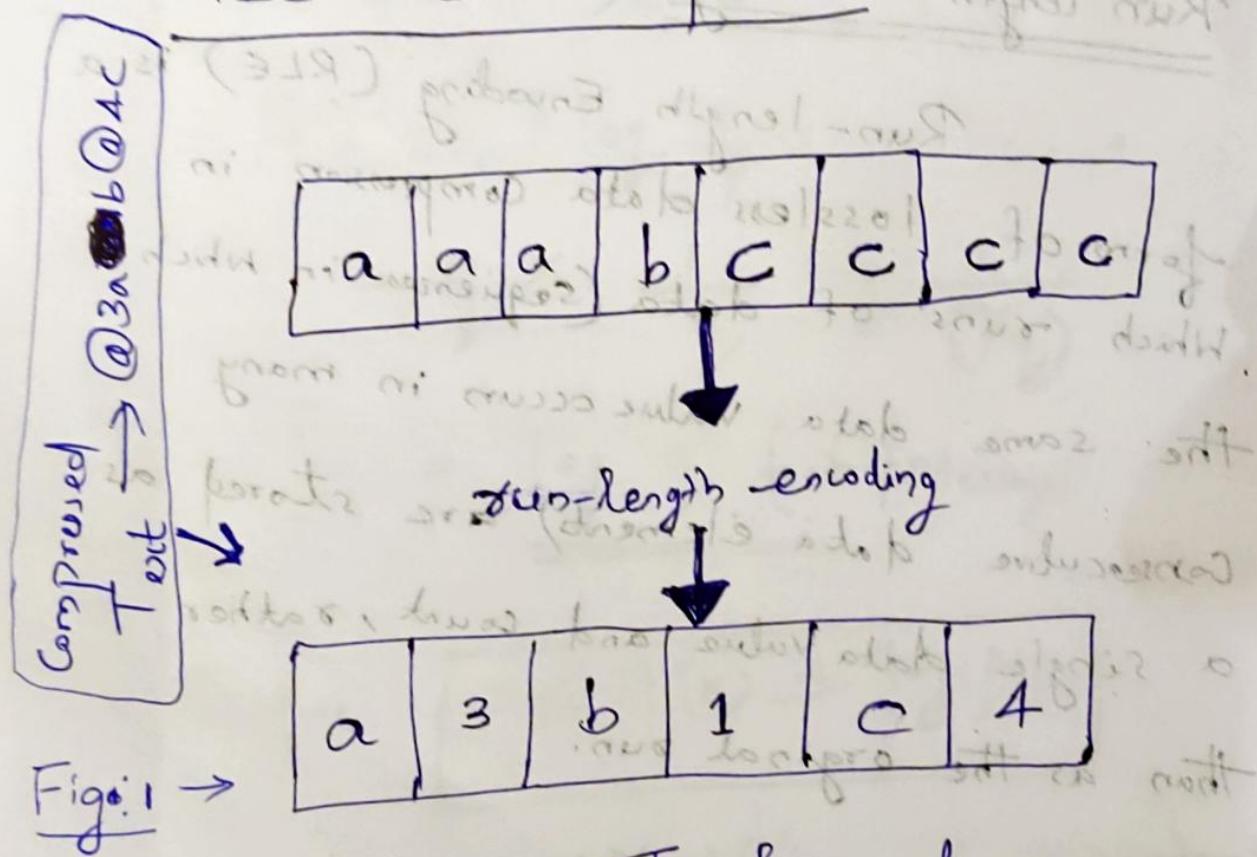


Fig:1 →

This type of compression

works best with Simple images and

(3)

animations that have a lot of

redundant pixels.

The following figure shows the flowchart for a simple run-length text compressor.

After reading the first character, the repeat-count is 1 and the character is saved. Subsequent characters are compared with the one already saved, and if they are identical to it, the repeat-count is incremented. When a different character is read, the operation depends on the value of the repeat-count. If it is small, the saved character is written on the

(3)

Compressed file, and the newly-read character is saved. Otherwise, an '@' is written followed by the repeat-count and the saved character.

Decompression is also shown in the following flow chart. When an '@' is read, the repetition count 'n' and the actual character are immediately read, and the character is written 'n' times in the output stream.

The decompressed data corresponding to fog::

$\rightarrow$  @3ab@4c.

# Flowchart of RLE Compression

(5)

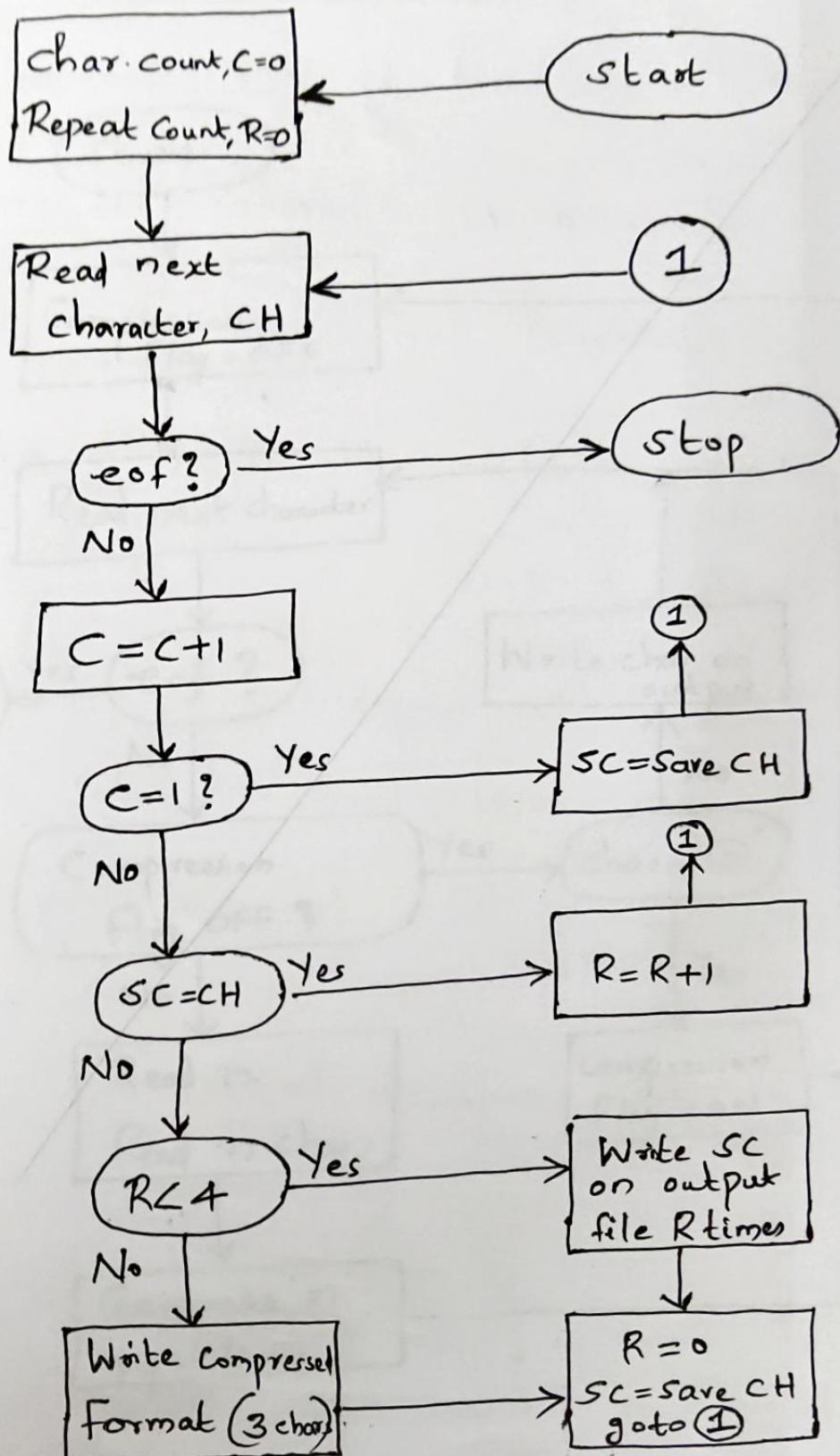


Fig:- RLE - Compression

# 6

## Flowchart of RLE Decompression

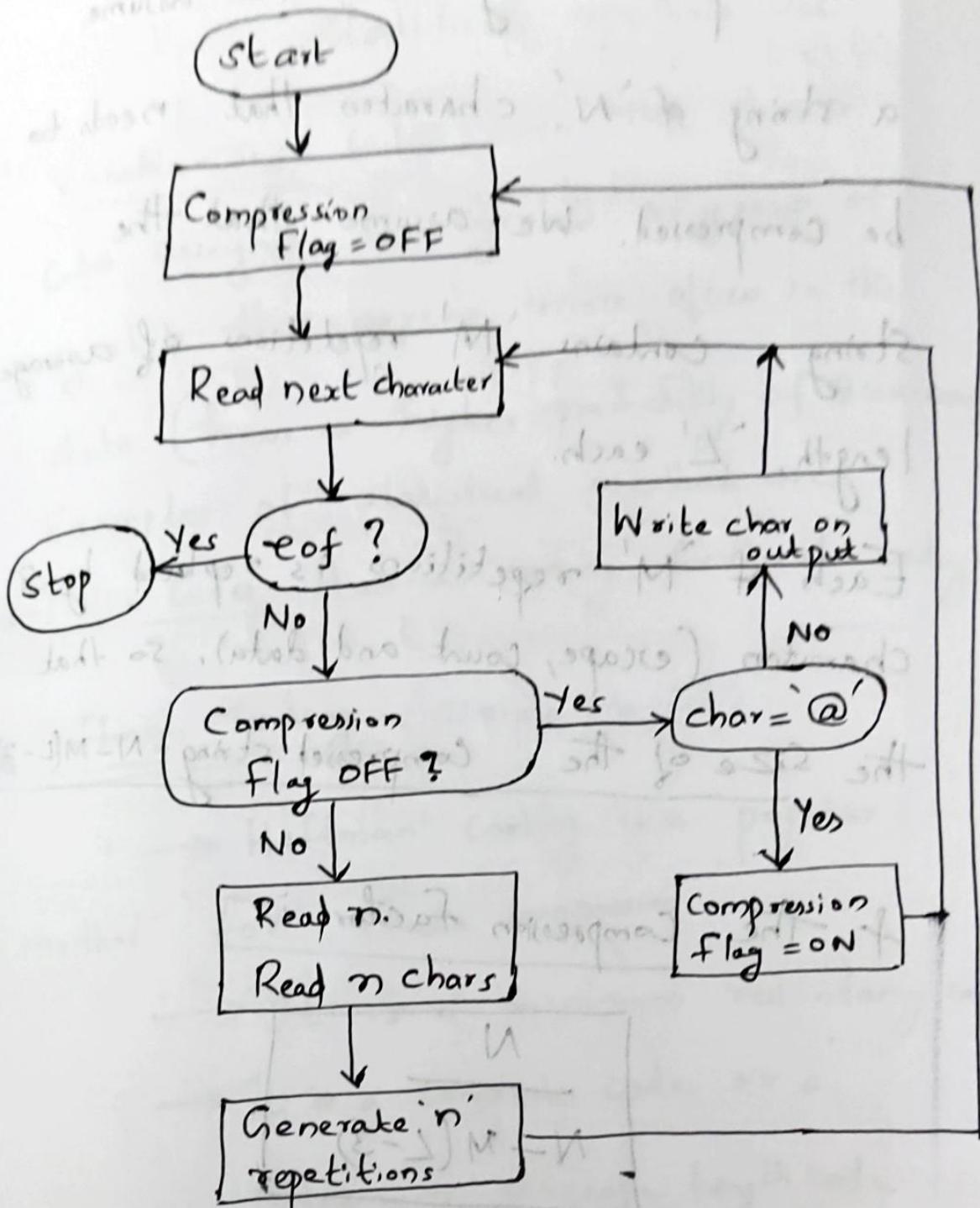


Fig- RLE Decompression.

To get an idea of compression ratios produced by RLE, we assume a string of 'N' characters that needs to be compressed. We assume that the string contains 'M' repetitions of average length, 'L' each.

Each of 'M' repetitions is replaced by 3 characters (escape, count and data), so that the size of the compressed string =  $N - M(L-3)$

- & the compression factor is

$$\frac{N}{N - M(L-3)}$$

# Statistical Methods for compression

(8)

Statistical methods use variable-size codes, with the shorter codes assigned to symbols or group of symbols that appear more often in the data (have a higher probability of occurrence). Examples of statistical methods are

Prefix Codes, Huffman Coding, Arithmetic Coding etc.

already studied  
in Module I

## Huffman Coding (Simple Method)

→ Huffman Coding is a popular method for data compression.

→ It is a minimum redundancy code.

→ It is a compact code or a minimum average length code.

Binary.

## Algorithm for Huffman coding

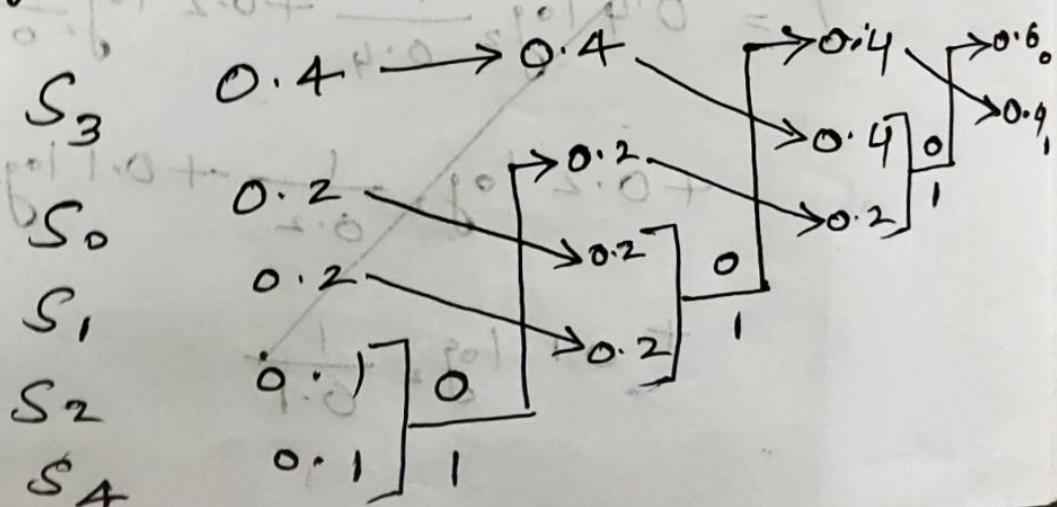
Step 1: List the source symbols in the decreasing order of probabilities. The two source symbols having the lowest probabilities are assigned 0 and 1.

Step 2: These two source symbols are regarded as being combined to a new source symbol with probability equal to the sum of two original probabilities. Probability of this new symbol is placed in the list in accordance with this value.

Step 3: This procedure is repeated until two sources will be remained for which 0 + 1 are assigned. The code for each original symbol is found by working Backwards and tracing sequences of 0's and 1's.

Q The source symbol,  $S = \{s_0, s_1, s_2, s_3, s_4\}$  has probabilities  $\{0.2, 0.2, 0.1, 0.4, 0.1\}$ . Obtain Huffman coding. Calculate the average size of this code.

Ans:- Average according to decreasing probabilities.



<u>Symbol</u>	<u>Code</u>	<u>Original code by bit reversal</u>
$S_3$	00	00
$S_0$	01	10
$S_1$	11	11
$S_2$	010	010
$S_4$	110	011

Average size of the code,

$$\bar{L} = \sum_{i=1}^q P_i l_i$$

~~$$\bar{L} = 0.4 \log_2 \frac{1}{0.4} + 0.2 \log_2 \frac{1}{0.2}$$~~

~~$$+ 0.2 \log_2 \frac{1}{0.2} + 0.1 \log_2 \frac{1}{0.1}$$~~

~~$$+ 0.1 \log_2 \frac{1}{0.1}$$~~

(13)

$$\begin{aligned} \bar{L} &= 0.4 \times 2 + 0.2 \times 2 + 0.2 \times 2 + 0.1 \times 3 \\ &\quad + 0.1 \times 3 \\ &= \underline{\underline{2.2 \text{ bits/symbol}}} \end{aligned}$$

HW The message 's' consists of 5 individual messages  $\{s_1, s_2, s_3, s_4, s_5\}$  & their probabilities are  $\{0.55, 0.15, 0.15, 0.1, 0.05\}$   
Obtain Huffman code and its efficiency.

Hint:- Efficiency  $\gamma = \frac{H(s)}{\bar{L}} \rightarrow$  Entropy  
 $\bar{L} \rightarrow$  Average Length of code

$$\text{Redundancy} = 1 - \gamma$$

Note:- Length of Huffman codes :-  $H(s) \leq \bar{L} \leq H(s) + 1$

$H(s) \rightarrow$  Entropy  
 $\bar{L} \rightarrow$  Average code word length of Huffman code.

Huffman Tree

→ The Huffman code

Can be also <sup>made</sup> created by creating

Huffman trees.



Q. Given the 8 symbols A, B, C, D, E, F, G, H (13)

With probabilities,  $\frac{1}{30}, \frac{1}{30}, \frac{1}{30}, \frac{2}{30}$ ,  
respectively.

$\frac{3}{30}, \frac{5}{30}, \frac{5}{30}, \frac{12}{30}$ . Draw

Huffman tree and calculate the average

size of code tree in bits per symbol.

Ans:-

Let,

$\frac{1}{30} \rightarrow A$

$\frac{5}{30} \rightarrow G$

$\frac{1}{30} \rightarrow B$

$\frac{12}{30} \rightarrow H$

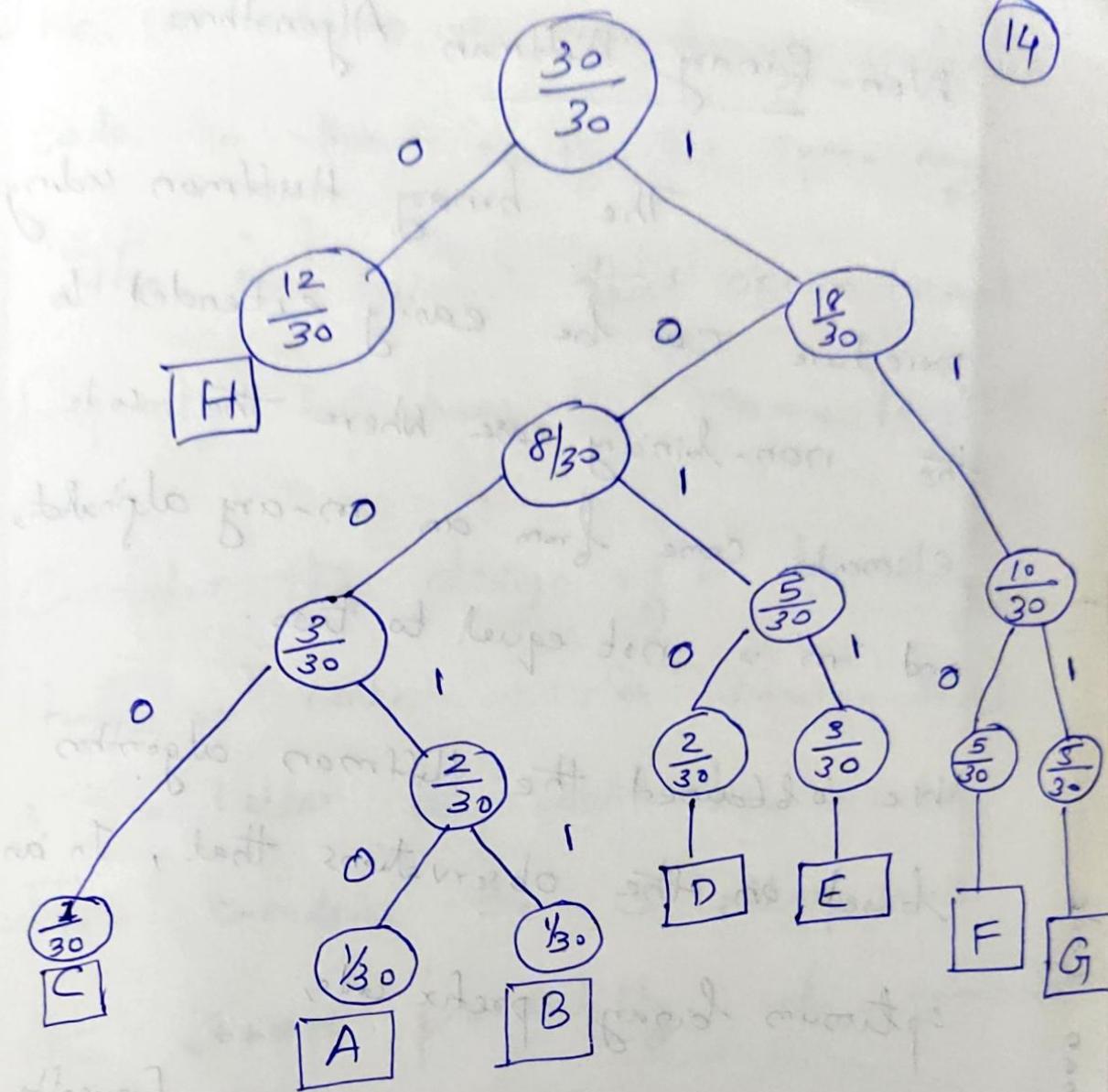
$\frac{1}{30} \rightarrow C$

$\frac{2}{30} \rightarrow D$

$\frac{3}{30} \rightarrow E$

$\frac{5}{30} \rightarrow F$

14



$$\frac{1}{30} \rightarrow$$

$$Cd\ size = \frac{1}{30} \times 5 + \frac{1}{30} \times 5 + \frac{1}{30} \times 4$$

$$+ \frac{2}{30} \times 4 + \frac{3}{30} \times 4 + \frac{5}{30} \times 3 + \frac{5}{30} \times 3$$

$$+ \frac{12}{30} \times 1 = \underline{\underline{\frac{76}{30}} \text{ bits/symbol}}$$

## Non-Binary Huffman Algorithms

The binary-Huffman coding procedure can be easily extended to the non-binary case where the code elements come from an  $m$ -ary alphabet, and ' $m$ ' is not equal to two.

We obtained the Huffman algorithm based on the observations that, In an optimum binary prefix code,

- (i) Symbols that occur more frequently will have shorter codewords than symbols that occur less frequently.
- (ii) the two symbols that occur least frequently will have the same length.

Observations from Huffman code

(16)

We can obtain a Non-binary Huffman code in almost exactly the same way.

"The 'm' symbols that occur least frequently will have the same length".

Consider the design of a Ternary Huffman code for a source with a 6-letter alphabet. We would first combine the 3 letters with the lowest probability into a composite letter. This would give us a reduced alphabet with four letters.

(17)

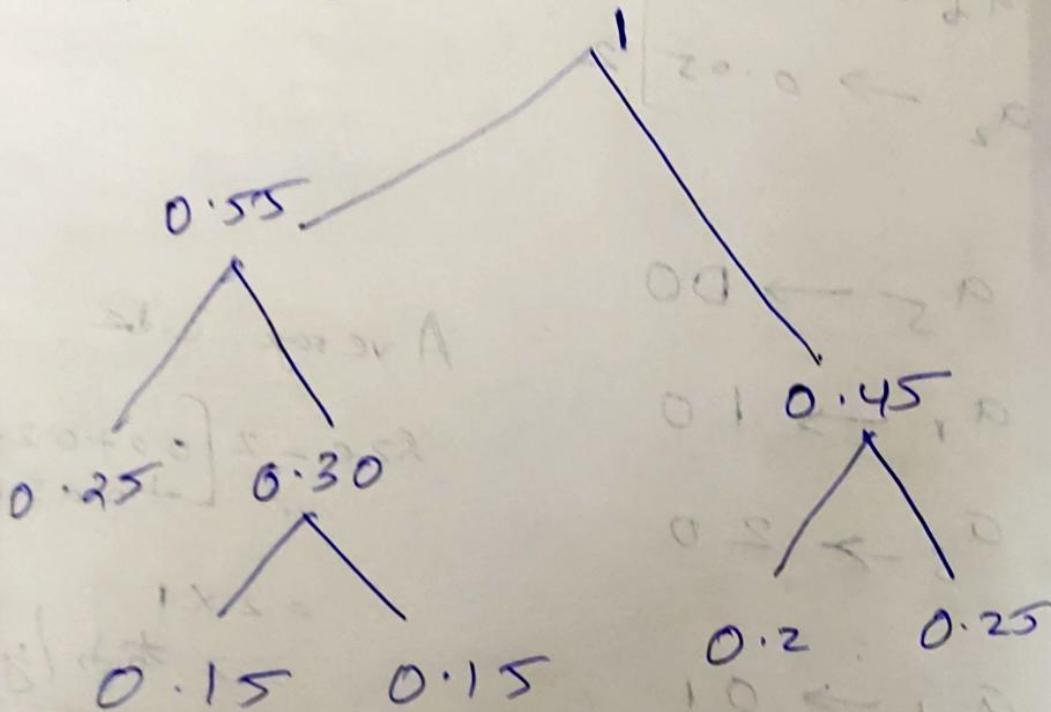
Huffman Binary Tree

Q. Generate a Huffman code tree for five symbols with probabilities

$0.15$ ,  $0.15$ ,  $0.2$ ,  $0.25$ , and  $0.25$

Find Average code size.

Ans:-



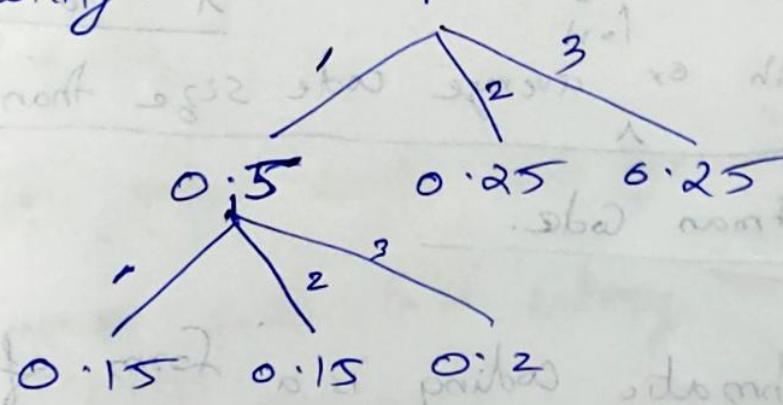
Average code size,

$$2 \times 0.25 + 3 \times 0.15 + 3 \times 0.15$$

$$+ 2 \times 0.2 + 2 \times 0.25 = 2.3 \text{ bits}$$

Ans to Question ⑯ Huffman Ternary Tree (for the same problem)

IF the tree is constructed by selecting three symbols with smallest probabilities and merging them into one parent symbol with the combined probability.



Here, Average code size

$$= 2 \times 0.15 + 2 \times 0.15 + 2 \times 0.125$$

$$+ 1 \times 0.25 + 1 \times 0.25$$

$$\underline{\underline{= 1.5 \text{ bits/symbol}}}$$

## Arithmetic Coding

The main difference between

Arithmetic coding and Huffman Coding is that

Arithmetic coding provides ~~a low average code~~

~~length or low average code size than the~~

Huffman code.

Arithmetic coding is a form of

entropy encoding used in lossless

data compression.

Arithmetic coding differs from other forms of entropy encoding such as

Huffman coding in that, ~~the~~ rather

than reoperating on the inputs into this

component symbols and replacing each

with a code, arithmetic coding encodes

the entire message into a single number,

an arbitrary precision factor  $q$ , where

$$0 < q < 1$$

→ This code works well for

Sequences with low entropy.

Eg:  $\rightarrow$  Main characteristics

- $\rightarrow$  It is a stream based encoding
- $\rightarrow$  Lossless Data compression
- $\rightarrow$  Non Block code
- $\rightarrow$  Efficient code

## Arithmetic Encoding Procedure

Q. Consider the transmission of a message "Went." comprising a string of characters with prob. of

$$e \rightarrow 0.3, n \rightarrow 0.3, t \rightarrow 0.2, o \rightarrow 0.1,$$

∴ → 0.1. Perform Arithmetic

Encoding. for lossless compression

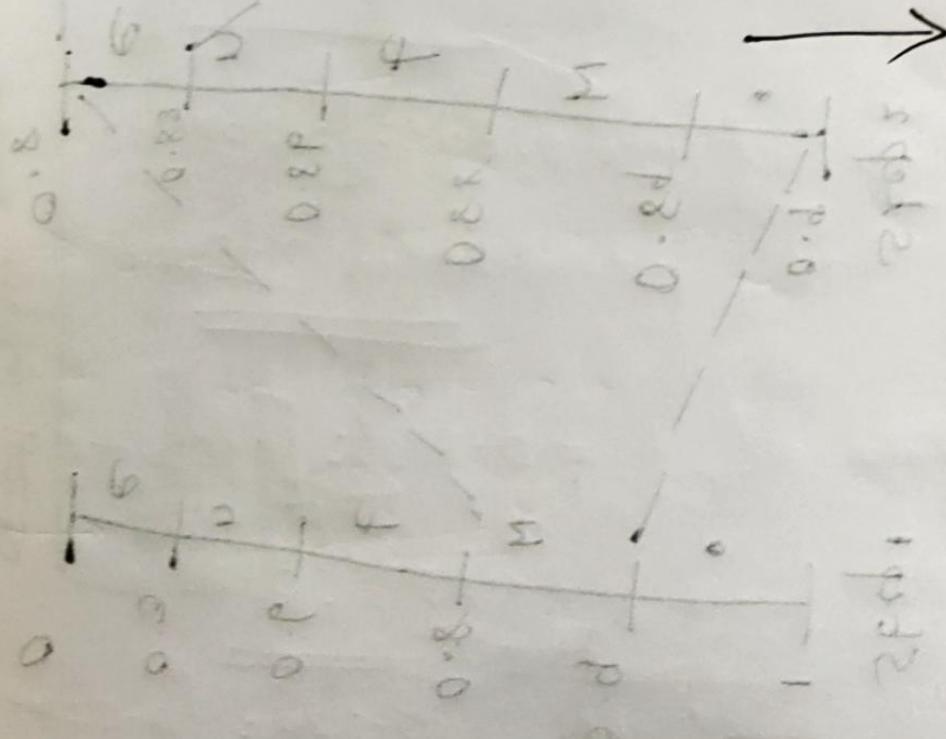
### Ans: Procedure

Step 1:- Divide the numeric range 0 to 1 into number of different symbol present in the message.

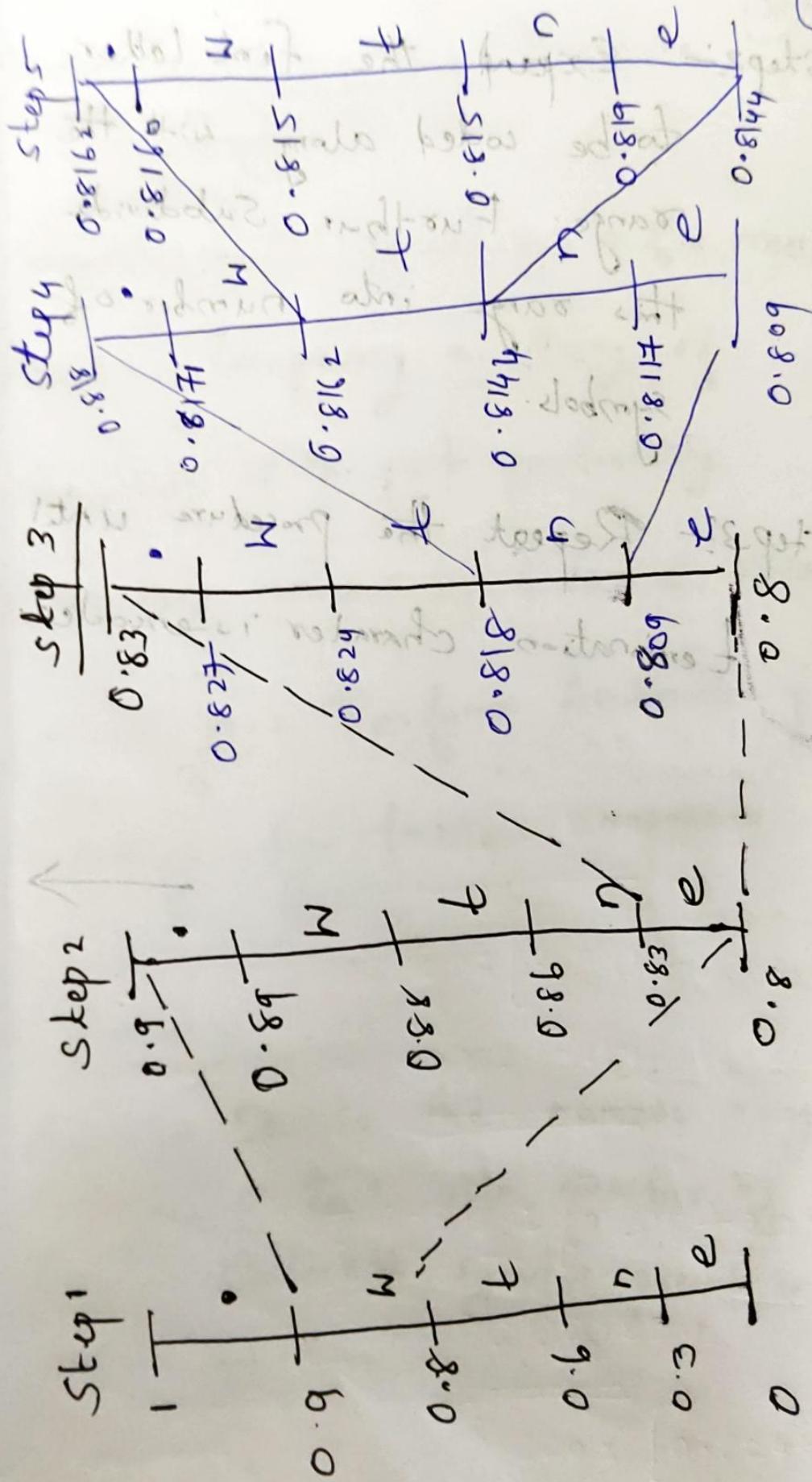
(22)

Step 2:- Expand the first letter to be coded along with the range. Further subdivide this range into number of symbols.

Step 3:- Repeat the procedure until termination character is encoded.



(23)



## Step 2 Calculations

$$d = \text{upper bound} - \text{lower bound} = 0.9 - 0.8 = 0.1$$

Range of symbol = lower limit : lower limit  
 + d (probability of symbol)

$$\text{Range of } m = 0.83 : [0.83 + 0.1(0.3)]$$

$$= 0.83 : 0.86$$

$$\text{Range of } n = 0.89 : [0.89 + 0.1(0.9)]$$

$$\text{Range of } n' = 0.88 : [0.88 + 0.1(0.1)]$$

$$= 0.88 : 0.89$$

(24)

(25)

### Step 3 Calculation

$$d = \text{Upper bound} - \text{Lower bound}$$

$$= 0.83 - 0.8 = \underline{\underline{0.03}}$$

We know that,  $\rightarrow$  Range of symbol = lower limit + d (Prob. of symbol)

$$\therefore \text{Range of } e = 0.8 : [0.8 + 0.03(0.3)] \\ = 0.8 : \underline{\underline{0.809}}$$

$$\text{Range of } n = 0.809 : [0.809 + 0.03(0.3)] \\ = 0.809 : \underline{\underline{0.818}}$$

$$\text{Range of } t = 0.818 : [0.818 + 0.03(0.2)] \\ = 0.818 : \underline{\underline{0.824}}$$

$$\text{Range of } w = 0.824 : [0.824 + 0.03(0.1)]$$

$$= 0.824 : 0.827$$

Range of ~~0.82~~ :

$$= 0.827 : [0.827 + 0.03(0.1)]$$

$$= 0.827 : 0.83$$

### Step 4 Calculations

Similarly,

We calculate,

$$\text{Range of } \bar{e} = 0.809 : 0.8117$$

$$\text{Range of } \bar{n} = 0.8117 : 0.8144$$

$$\text{Range of } \bar{t} = 0.8144 : 0.8162$$

$$\text{Range of } \bar{w} = 0.8162 : 0.8171$$

$$\text{Range of } \bar{o} = 0.8171 : 0.818$$

Step 5

(27)

Range of ' $e$ ' =  $0.8144 : 0.8194$

" " " $n$ " =  $0.81448 : 0.81548$

" " " $t$ " =  $0.81548 : 0.81584$

" " " $w$ " =  $0.81584 : 0.81602$

" " " $s$ " =  $0.81602 : 0.8162$

Finally,

the Arithmetic code

Word is the codeword

between the limits of

the termination characters.

i.e,

$0.81602 < \text{codeword} < 0.8162$

Generation of Tag:

$$\text{Tag} = \frac{(\text{Upper limit of codeword} + \text{Lower limit of codeword})}{2} \quad (28)$$

$$\begin{aligned}\text{Tag} &= \frac{(0.8162 + 0.81602)}{2} \\ &= 0.81611\end{aligned}$$

Tag = (Upper limit of code word

+ Lower limit of code word) / 2

$$\text{Tag} = \frac{(0.8162 + 6.81602)}{2}$$

$$= 0.81611$$

Comparison of Huffman code and

Arithmetic code

Arithmetic code	Huffman code
<ul style="list-style-type: none"> <li>It is not a statistical Method</li> </ul>	<ul style="list-style-type: none"> <li>It is a statistical Method</li> </ul>
<ul style="list-style-type: none"> <li>It yields an optimum result.</li> </ul>	<ul style="list-style-type: none"> <li>It doesn't yield an optimum result.</li> </ul>
<ul style="list-style-type: none"> <li>There is no one to one correspondence b/w source symbol and code word</li> </ul>	<ul style="list-style-type: none"> <li>There is one to one correspondence b/w source symbol and code word</li> </ul>

### Arihmetic code

- There is only one unique code for entire message

### Huffman code

- Separate symbols are assigned for each symbols.

Advantages of Arithmetic code over

### Huffman code

→ Arithmetic coding typically has a better compression ratio than Huffman coding, as it produces a single symbol rather than several separate codewords.

→ Arithmetic coding surpasses the Huffman techniques in its compression ability.

→ In terms of complexity, arithmetic coding is better than Huffman Coding.