

Huffman coding is a lossless data compression algorithm. In this algorithm, a variable-length code is assigned to input different characters. The code length is related to how frequently characters are used. Most frequent characters have the smallest codes and longer codes for least frequent characters.

There are mainly two parts. First one to create a Huffman tree, and another one to traverse the tree to find codes.

- Huffman Coding is a famous Greedy Algorithm.
- It is used for the lossless compression of data.
- It uses variable length encoding.
- It assigns variable length code to all the characters.
- The code length of a character depends on how frequently it occurs in the given text.
- The character which occurs most frequently gets the smallest code.
- The character which occurs least frequently gets the largest code.
- It is also known as **Huffman Encoding**.

Huffman Coding

BCCABBDDEAECCBBAEDDCC

length=20

ASCII - 8-bit

$8 \times 20 = 160 \text{ bits}$

A	65	01000001
B	66	01000010
C	67	:
D	68	:
E	69	:

## Huffman Coding

Message  $\rightarrow$  BCCABBBDDAECCBBBAEDDCC

character	count/frequency	Code	0/1 bit
A	3 $3/20$	000	0 0 0 bit
B	5 $5/20$	001	0 0 1 bit
C	6 $6/20$	010	0 1 0 bit
D	4 $4/20$	011	0 1 1 bit
E	2 $2/20$	100	1 0 0 bit

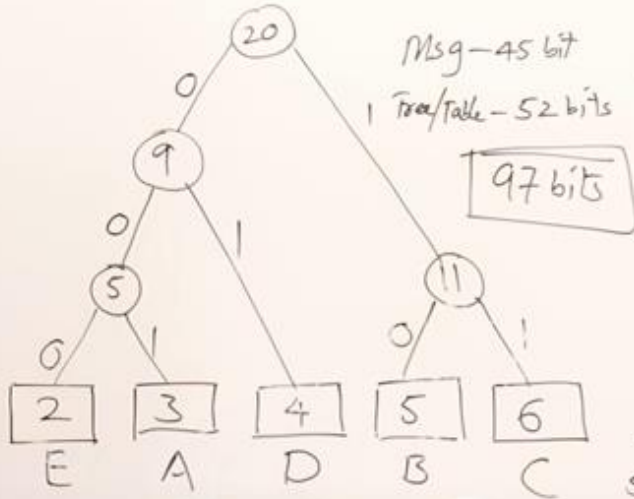
## Huffman Coding

Message  $\rightarrow$  BCCABBBDDAECCBBBAEDDCC  
001010 . . .

character	count/frequency	Code	20x3=60 bits
A	3 $3/20$	000	$5 \times 8 \text{ bit}$ $5 \times 3$
B	5 $5/20$	001	$\uparrow$ characters $\uparrow$ codes
C	6 $6/20$	010	$40 + 15 = 55$
D	4 $4/20$	011	Msg - 60 bits
E	2 $2/20$	100	Table - 55 bits
20			115 bits

# Huffman Coding

Message  $\rightarrow$  BCCABBDDECCBBAEDDCC  
 10 11 11 00 110 10 01 01 - - -

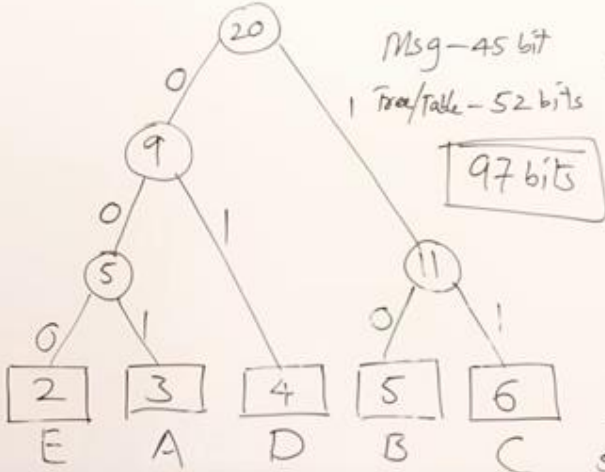


char	count	Code	
A	3	001	$3 \times 3 = 9$
B	5	10	$5 \times 2 = 10$
C	6	11	$6 \times 2 = 12$
D	4	01	$4 \times 2 = 8$
E	2	000	$2 \times 3 = 6$
		20	
5x 8bit		40 bits	
		12 bit	
		45 bit	

Msg-45 bit  
 Tree/Table-52 bits  
**97 bits**

# Huffman Coding

Message  $\rightarrow$  BCCABBDDECCBBAEDDCC  
 10 11 11 00 110 10 01 01 - - -



char	count	Code	
A	3	001	$3 \times 3 = 9$
B	5	10	$5 \times 2 = 10$
C	6	11	$6 \times 2 = 12$
D	4	01	$4 \times 2 = 8$
E	2	000	$2 \times 3 = 6$
		20	
5x 8bit		40 bits	
		12 bit	
		45 bit	

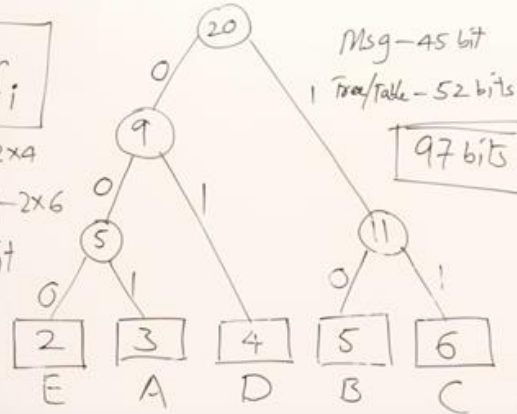
Msg-45 bit  
 Tree/Table-52 bits  
**97 bits**

# Huffman Coding

Message  $\rightarrow$  BCCABBDDECCBBAEDDCC  
 10 11 11 00 110 10 01 01 - - -

$$\sum d_i * f_i$$

3x2 + 3x3 + 2x4  
 1x5 + 2x6  
 5 bit

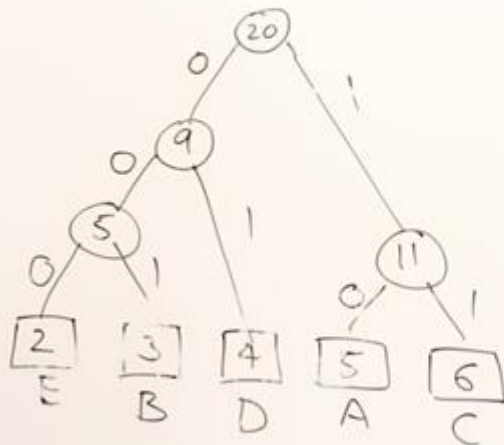


Char	count	Code
A	3	001
B	5	10
C	6	11
D	4	01
E	2	000

5x 8 bit 40 bits + 5 bit 45 bits

# Huffman Coding

Message  $\rightarrow$  BCCDACCBDABCCDEAAEDA  
 001 11 11 01 10 00 11 11 01 - - -



BCCD

Example 2

# Huffman Coding

Huffman Coding is a technique of compressing data to reduce its size without losing any of the details. It was first developed by David Huffman.

Huffman Coding is generally useful to compress the data in which there are frequently occurring characters.

---

## How Huffman Coding works?

Suppose the string below is to be sent over a network.

B	C	A	A	D	D	D	C	C	A	C	A	C	A	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Initial string

Each character occupies 8 bits. There are a total of 15 characters in the above string. Thus, a total of  $8 * 15 = 120$  bits are required to send this string.

Using the Huffman Coding technique, we can compress the string to a smaller size.

Huffman coding first creates a tree using the frequencies of the character and then generates code for each character.

Once the data is encoded, it has to be decoded. Decoding is done using the same tree.

Huffman Coding prevents any ambiguity in the decoding process using the concept of **prefix code** ie. a code associated with a character should not be present in the prefix of any other code. The tree created above helps in maintaining the property.

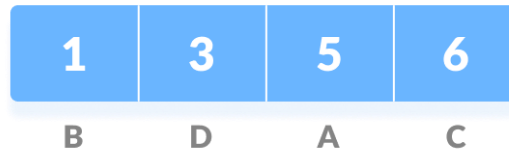
Huffman coding is done with the help of the following steps.

1. Calculate the frequency of each character in the string.



Frequency of string

2. Sort the characters in increasing order of the frequency. These are stored in a priority

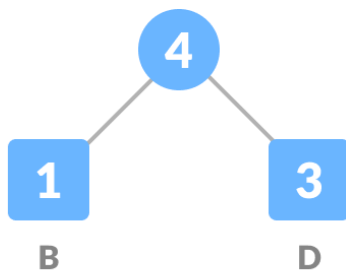
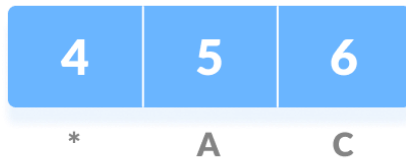


queue  $Q$ .

to the frequency

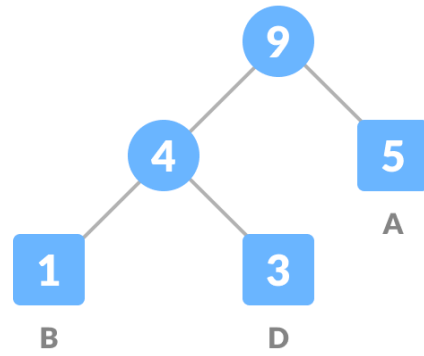
Characters sorted according

3. Make each unique character as a leaf node.
4. Create an empty node  $z$ . Assign the minimum frequency to the left child of  $z$  and assign the second minimum frequency to the right child of  $z$ . Set the value of the  $z$  as the sum of the above two minimum frequencies.

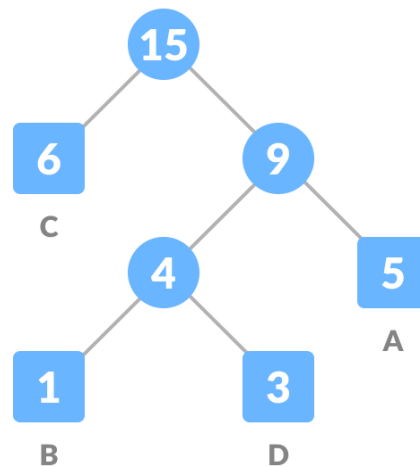


Getting the sum of the least numbers

5. Remove these two minimum frequencies from  $Q$  and add the sum into the list of frequencies (\* denote the internal nodes in the figure above).
6. Insert node  $z$  into the tree.

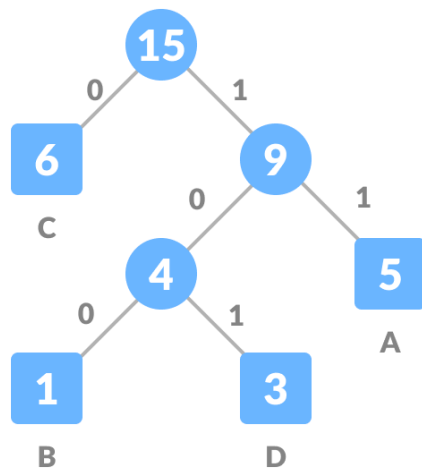


7. Repeat steps 3 to 5 for all the characters.



Repeat steps 3 to 5 for all the characters.  
Repeat steps 3 to 5 for all the characters.

8. For each non-leaf node, assign 0 to the left edge and 1 to the right edge.



Assign 0 to the left edge and 1 to the right edge

For sending the above string over a network, we have to send the tree as well as the above compressed-code. The total size is given by the table below.

Character	Frequency	Code	Size
A	5	11	$5 * 2 = 10$
B	1	100	$1 * 3 = 3$
C	6	0	$6 * 1 = 6$
D	3	101	$3 * 3 = 9$
$4 * 8 = 32$ bits			15 bits
			28 bits

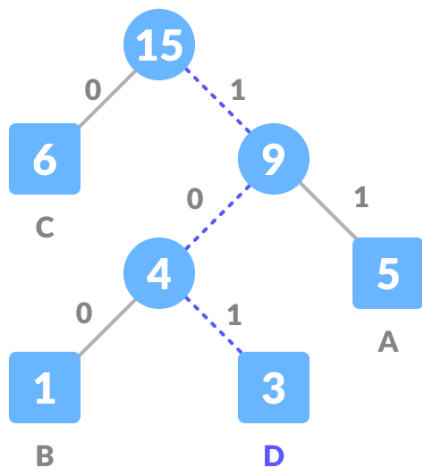
Without encoding, the total size of the string was 120 bits. After encoding the size is reduced to  $32 + 15 + 28 = 75$ .

## Decoding the code

For decoding the code, we can take the code and traverse through the tree to find the character.

Let 101 is to be decoded, we can traverse from the root as in the figure below.





Decoding