

## Unit IV

11L

Text Processing: Sting Operations, Brute-Force Pattern Matching, The Boyer- Moore Algorithm, The Knuth-Morris-Pratt Algorithm, Standard Tries, Compressed Tries, Suffix Tries, The Huffman Coding Algorithm, The Longest Common Subsequence Problem (LCS), Applying Dynamic Programming to the LCS Problem.

## Unit IV

Text Processing: Sting Operations, Brute-Force Pattern Matching, The Boyer- Moore Algorithm, The Knuth-Morris-Pratt Algorithm, Standard Tries, Compressed Tries, Suffix Tries, The Huffman Coding Algorithm, The Longest Common Subsequence Problem (LCS), Applying Dynamic Programming to the LCS Problem.

ADVISED TO WATCH ALL VIDEOS IN ATLEAST 1.5x SPEED

### String Intro

<https://www.youtube.com/watch?v=JvKMLSa0Rq0>

### String Operations

Operation	Description
Concatenation	The <b>Addition</b> operator, "+", can be used to concatenate strings together.
Formatting Data	The <b>STRING</b> function is used to format data into a string. The <b>READS</b> procedure can be used to read values from a string into IDL variables.
Case Folding	The <b>ToLower()</b> method returns a copy of the string converted to lowercase. Similarly, the <b>ToUpper()</b> method returns a copy converted to uppercase. The <b>CapWords()</b> method returns a copy where the first letter of each word is capitalized. See also the <b>STRLOWCASE</b> and <b>STRUPCASE</b> functions.
White Space Removal	The <b>Compress()</b> and <b>Trim()</b> methods can be used to eliminate unwanted white space. See also the <b>STRCOMPRESS</b> and <b>STRTRIM</b> functions.
Length	The <b>Strlen()</b> method (or <b>STRLEN</b> function) returns the length of the string.
Substrings	See the <b>CharAt()</b> , <b>Extract()</b> , <b>IndexOf()</b> , <b>Insert()</b> , <b>LastIndexOf()</b> , <b>Remove()</b> , <b>Replace()</b> , <b>Reverse()</b> , and <b>Substring()</b> methods. See also the <b>STRPOS</b> , <b>STRPUT</b> , and <b>STRMID</b> routines.
Splitting and Joining	The <b>Split()</b> method is used to break strings apart, and the <b>Join()</b> method can be used to join them. See also the <b>STRSPLIT</b> and <b>STRJOIN</b> functions.
Comparing Strings	The <b>"EQ" operator</b> can be used to directly compare two strings. The <b>Contains()</b> , <b>EndsWith()</b> , <b>Matches()</b> , and <b>StartsWith()</b> methods can be used to compare portions of strings. See also the <b>STRCMP</b> , <b>STRMATCH</b> , and <b>STREGEX</b> functions.

<https://www.youtube.com/watch?v=-pSyzCWsBA8>

### String/Text Processing

Insert- <https://www.youtube.com/watch?v=Rj1Tcb5RFzU>

Delete- <https://www.youtube.com/watch?v=SYeMm5HiJQ0>

Replace- <https://www.youtube.com/watch?v=uLZUfggh-DQ>

**Boyer Moore** : check the pdf that I sent u/ attached below.

[Use only last function is there not bad match table (according to sir) so, don't waste your time watching YouTube videos on this, coz max are of bad match table.]

[https://www.youtube.com/watch?v=4Oj\\_ESzSNck](https://www.youtube.com/watch?v=4Oj_ESzSNck)

**Brute-Force Pattern Matching**: <https://www.youtube.com/watch?v=yMJLpdKV0BQ>

**Knuth-Morris-Pratt Algorithm**:

Easy Engineering classes- ADA

Eg on KMP Algorithm:

T: a b a b c a b c a b a b a b d

i: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

P: a b a b d

j: 0 1 2 3 4 5

a	b	a	b	d
0	0	1	2	0

$T[s] = P[s]$  mismatch

IMP

- i) Take two Variables i and j
- $i = \text{String}(T(1)), j = P[0]$
- ii) Compare  $T(i)$  with  $P(j+1)$ 
  - i) if Match is found (Move both i and j to right)
  - ii) if Mismatch (move j to the loc<sup>n</sup> as per  $\pi$  table index)
  - iii) if  $j=0$  } move i to the right

EEC Classes

<https://www.youtube.com/watch?v=V5-7GzOfADQ> [better]

<https://www.youtube.com/watch?v=H3HFstXmnP8> [if less time, then see this]

**Failure function visualization**: <https://www.educative.io/edpresso/what-is-the-knuth-morris-pratt-algorithm>

**Standard Tries, Compressed Tries, Suffix Tries**: TRIES PPT {UPLOADED IN DRIVE} [prefer this]

<https://drive.google.com/file/d/1saZrFjk5mWUdw9q0PS9apCina4DX6Yct/view?usp=sharing>

[https://www.geeksforgeeks.org/types-of-](https://www.geeksforgeeks.org/types-of-tries/#:~:text=A%20Compressed%20Trie%20is%20an,of%20redundant%20nodes%20is%20performed.)

[tries/#:~:text=A%20Compressed%20Trie%20is%20an,of%20redundant%20nodes%20is%20performed.](https://www.geeksforgeeks.org/types-of-tries/#:~:text=A%20Compressed%20Trie%20is%20an,of%20redundant%20nodes%20is%20performed.) [2<sup>nd</sup> resource]

Huffman coding:

## Huffman Algorithm

- **Step 1: Get Frequencies**
  - Scan the file to be compressed and count the occurrence of each character
  - Sort the characters based on their frequency
- **Step 2: Build Tree & Assign Codes**
  - Build a Huffman-code tree (binary tree)
  - Traverse the tree to assign codes
- **Step 3: Encode (Compress)**
  - Scan the file again and replace each character by its code
- **Step 4: Decode (Decompress)**
  - Huffman tree is the key to decompress the file

Slide 9

[https://www.youtube.com/watch?v=co4\\_ahEDCho&t=158s](https://www.youtube.com/watch?v=co4_ahEDCho&t=158s) [For understanding]

<https://www.programiz.com/dsa/huffman-coding> [Algorithm]

**The Longest Common Subsequence Problem (LCS):**

LONGEST-COMMON-SUBSEQUENCE( $x, m, y, n$ )

```
1  for  $i \leftarrow -1$  to  $m - 1$ 
2      do  $T[i, -1] \leftarrow 0$ 
3  for  $j \leftarrow -1$  to  $n - 1$ 
4      do  $T[-1, j] \leftarrow 0$ 
5  for  $i \leftarrow 0$  to  $m - 1$ 
6      do for  $j \leftarrow 0$  to  $n - 1$ 
7          do if  $x_i = y_j$ 
8              then  $T[i, j] \leftarrow T[i - 1, j - 1] + 1$ 
9              else  $T[i, j] \leftarrow \max(T[i, j - 1],$ 
                                    $T[i - 1, j])$ 
10 return  $T$ 
```

<https://www.youtube.com/watch?v=sSno9rV8Rhg>

Applying Dynamic Programming to the LCS Problem.

APV

11/2/21:  
Thursday      The Boyer Moore Algorithm:

Text: a b a c a a b a d c a b a c a b a a b b

pattern: a b a c a b

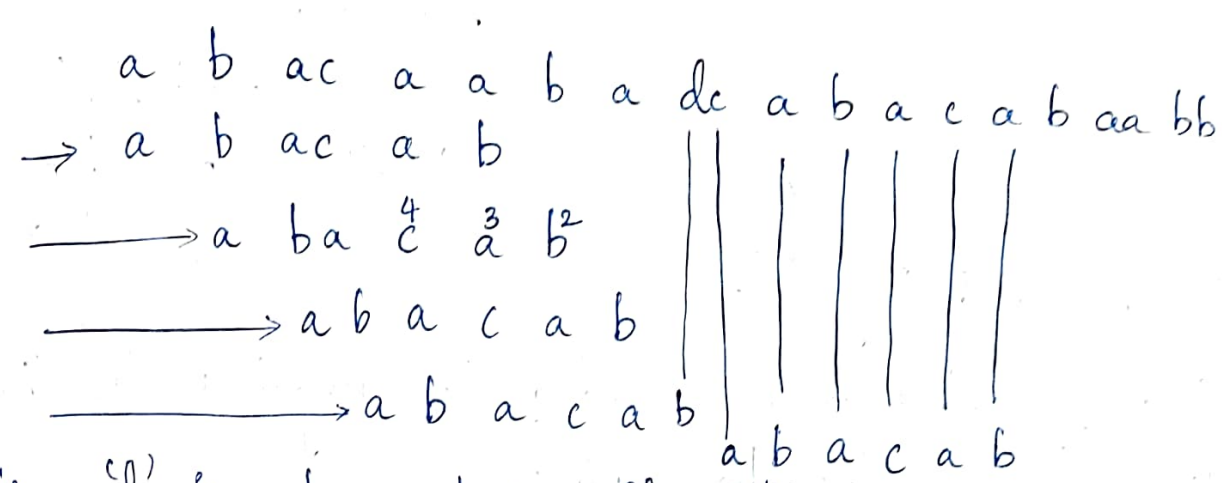
Maximum index of 'a' in pattern: 4

" " " 'b' " " : 5

" " " 'c' " " : 3

" " " 'd' " " : 0

If 'd' is encountered in text, the entire pattern until 'd' is skipped.



Since 'd' is not present in pattern, the pattern is not compared with 'd'. And hence we skip '6' positions.



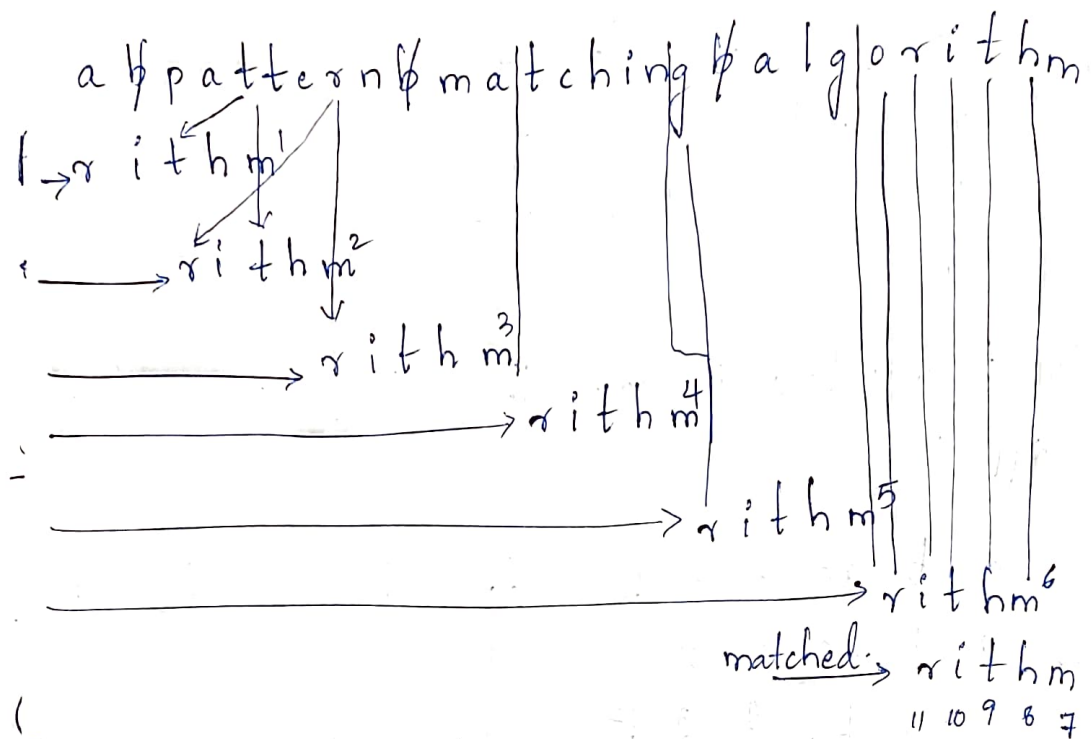
↯ → blank space

Text: a ↯ pattern ↯ matching ↯ algorithm

pattern: rithm

a ↯ pattern ↯ matching ↯ algorithm

char | a p t e r n m i h g l o  
last(char)



Drawback with Boyer Moore

If no. of symbols are more; and pattern length is less then less non-zero values. less skips.

Test = 0 1 2 3 4 5 6  
a a a a a a a

P = 0 1 2 3 4  
b a a a a  
~~b a a a a~~  
~~a a a a a~~

Prefix matches and prefix mismatches.

And so the shift occurs only by one position.

Compute the last function.

$i \leftarrow m-1$

$j \leftarrow m-1$

repeat

if  $P(j) = T(i)$  then

if ( $j == 0$ ) return  $i$

else

$i \leftarrow i - 1$

$j \leftarrow j - 1$

else

$i \leftarrow i + m - \min(j, 1 + \text{last}(\text{TC}(i)))$

$j \leftarrow m - 1$

until  $i > n - 1$

return 'pattern not in text'

Boyer Moore for  $\Sigma = \{-, a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z\}$

a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v
30	9	7	33	24	14	35	27	6	17	8	29	19	13	34	20	4	10	21	26	18	23

w	x	y	z	-
12	16	32	31	28

String: The quick brown fox jumps over the lazy dog  
 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32  
 33 34 35