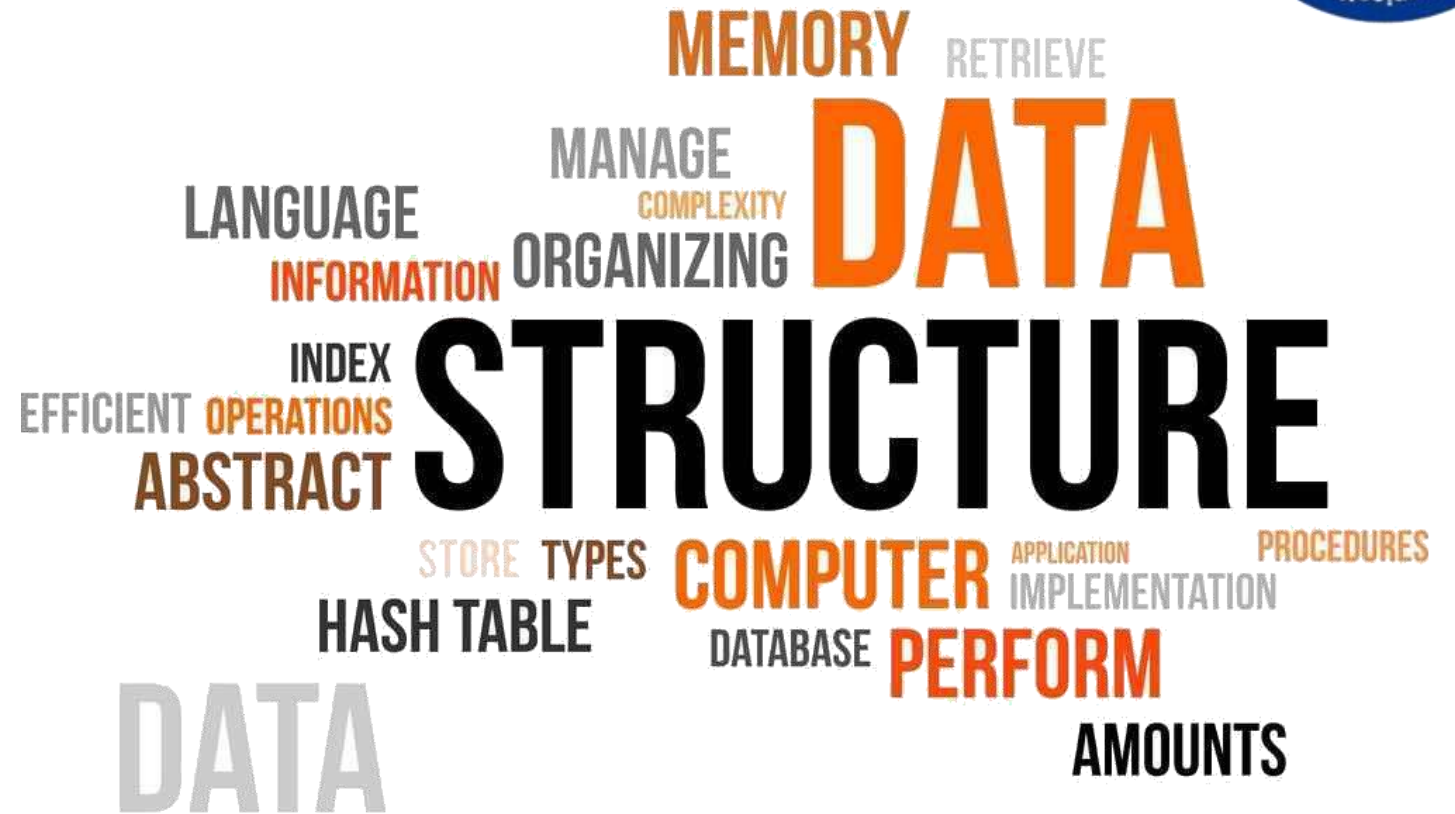




# Data Structures

Course code: IT623



**Dr. Rahul Mishra**  
**Assistant Professor**  
**DA-IICT, Gandhinagar**

# Lectures 13



## *String Processing...*

## Second Pattern Matching Algorithm:-

(21)

\* The second pattern matching algorithm uses a table which is derived from a particular pattern  $P$  but is independent of the text  $T$ . For definiteness, suppose

$\Rightarrow$  "  $P = aabab$  " pattern

\* First we give the reason for the table entries and how they are used.

\* Let  $T = T_1 T_2 T_3 \dots$ , where  $T_i$  denotes the  $i^{\text{th}}$  character of  $T$ ; and suppose the first two characters of  $T$  match those of  $P$ ; i.e. suppose  $T = aa \dots$ . Then  $T$  has one of the following three forms

(i)  $T = aab \dots$ ,

(ii)  $T = aaa \dots$ ,

(iii)  $T = aa\alpha$

Where  $\alpha$  is any character different from  $a$  or  $b$ .



\* Let we read  $T_3$  and find that  $T_3 = b$ . Then we next read  $T_4$  to see if  $T_4 = a$ ,  
I  $\rightarrow$  which will give a match of  $P$  with  $W_1$ .

\* Let  $T_3 = a$  then we know that  $P \neq W_1$ ; but we also know that  $W_2 = aa \dots$ , i.e., that  
II  $\rightarrow$  the first two characters of the substring  $W_2$  match those of  $P$ . We next read  $T_4$  to  
see if  $T_4 = b$ .

\*

III  $\rightarrow$  Let  $T_3 = \alpha$  then  $P \neq W_1$  but we also know that  $P \neq W_2$  and  $P \neq W_3$ . Since  $\alpha$   
does not appear in  $P$ . We next read  $T_4$  to see if  $T_4 = a$  i.e., to see if the first  
character of  $W_4$  matches the first character of  $P$ .

Important:  $\rightarrow$  ... Conclusion from above three cases:

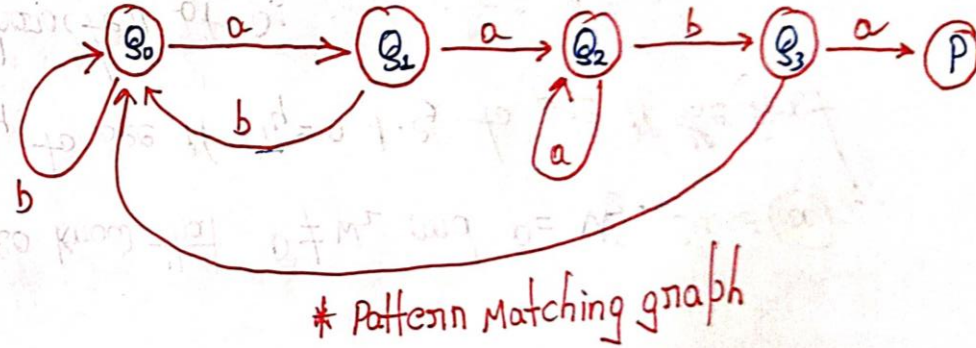
(a) When we read  $T_3$  we need only compare  $T_3$  with those character  
which appear in  $P$ . If none of these match, we are in the last  
case of a character  $\alpha$  which does not appear in  $P$

⑥ After reading and checking  $T_3$ , we next read  $T_4$ ; we do not have to go back again in the text  $T$ . (24)

Consider on example ÷

	a	b	x
$Q_0$	$Q_1$	$Q_0$	$Q_0$
$Q_1$	$Q_2$	$Q_0$	$Q_0$
$Q_2$	$Q_2$	$Q_3$	$Q_0$
$Q_3$	P	$Q_0$	$Q_0$

\* Pattern matching table



The Figure contains the table that is used in our second pattern matching algorithm for the pattern  $P = aaba$ .

We always use italic letters to represent pattern in a string.



\* The table is obtained as follows.

\* Let  $Q_i$  denote the initial substring of  $P$  of length  $i$ ; hence

$Q_0 = \Lambda$	, $Q_1 = a$	, $Q_2 = a^2$	, $Q_3 = a^2b$	, $Q_4 = a^2ba = P$
-----------------	-------------	---------------	----------------	---------------------

$Q_0 = \Lambda$  is the empty string.

> The rows of the table are labeled by these initial substring of  $P$  excluding  $P$  itself.

> The columns of the table are labeled  $a$ ,  $b$ , and  $x$ , where  $x$  represents any character that doesn't appear in the pattern  $P$ .

\* Let  $f$  be the function determined by the table; i.e., let denote the entry in row  $Q_i$  and column  $x$ .

(26)

The entry  $f(Q_i, t)$  is defined to be the longest  $Q$  that appears as a terminal substring in the string  $Q_i t$ .

- $a^2$  is the longest  $Q$  that is a terminal substring of  $Q_2 a = a^3$ , so  $f(Q_2, a) = Q_2$
- $\Lambda$  is the longest  $Q$  that is a terminal substring of  $Q_1 b = ab$ , so  $f(Q_1, b) = Q_0$
- $a$  is the longest  $Q$  that is a terminal substring of  $Q_0 a = a$ , so  $f(Q_0, a) = Q_1$
- $\Lambda$  is the longest  $Q$  that is a terminal substring of  $Q_3 a = a^3 b x$ , so  $f(Q_3, x) = Q_0$

\*  $Q_1 = a$  is a terminal substring of  $Q_2 a = a^3$  we have  $f(Q_2, a) = Q_2$  because  $Q_2$  is also a terminal substring  $Q_2 a = a^3$ .



(27)

→ Pattern  $P = aaba$ . Let  $T = T_1 T_2 T_3 \dots T_N$  denote that  $n$  character-string text which is searched for the pattern  $P$ .

→ Beginning with the initial state  $Q_0$  and using the text  $T$ , will be obtain a sequence of states  $s_1, s_2, s_3, \dots$  as follows.

→ We let  $s_1 = Q_0$  and we read the first character  $T_1$ . From either the table or the graph in Figure,  $-- (s_1, T_1)$  yield new state, and so on.

Notably: \* (a) Some state  $s_k = P$ , the desired pattern. In this case, (a) does not appear in  $T$  and its index is  $k - \text{LENGTH}(P)$ . Example: aabcaba

(b) No state  $s_1, s_2, \dots, s_{N+1}$  is equal to  $P$ . In this case,  $P$  does not appear in  $T$ . Example: abcaabaca



(Pattern Matching). The pattern matching table  $F(Q_1, T)$  of a pattern  $P$  is in memory, and the input is an  $N$ -character string  $T = T_1T_2 \dots T_N$ . This algorithm finds the INDEX of  $P$  in  $T$ .

1. [Initialize.] Set  $K := 1$  and  $S_1 = Q_0$
2. Repeat Steps 3 to 5 while  $S_K \neq P$  and  $K \leq N$ .
3.     Read  $T_K$ .
4.     Set  $S_{K+1} := F(S_K, T_K)$ . [Finds next state.]
5.     Set  $K := K + 1$ . [Updates counter.]
- [End of Step 2 loop.]
6. [Successful?]  
    If  $S_K = P$ , then:  
        INDEX =  $K - \text{LENGTH}(P)$ .  
    Else:  
        INDEX = 0.  
    [End of If structure.]
7. Exit

## Pattern Matching Second Algorithm

\* The running time of the above algorithm is proportional to the number of times the step 2 loop is executed.

\* The worst case occurs when all of the text **T** is read, i.e., when the loop is executed  $n = \text{LENGTH}(T)$  times.

\* Thus, we can conclude that the complexity of this pattern matching algorithm is  $O(n)$ .

Two numerical examples -

①

②



Consider the pattern  $P = aaabb$ . Construct the table and the corresponding labeled directed graph used in the "fast," or second pattern matching, algorithm.

First list the initial segments of  $P$ :

$$Q_0 = \Lambda, \quad Q_1 = a, \quad Q_2 = a^2, \quad Q_3 = a^3, \quad Q_4 = a^3b, \quad Q_5 = a^3b^2$$

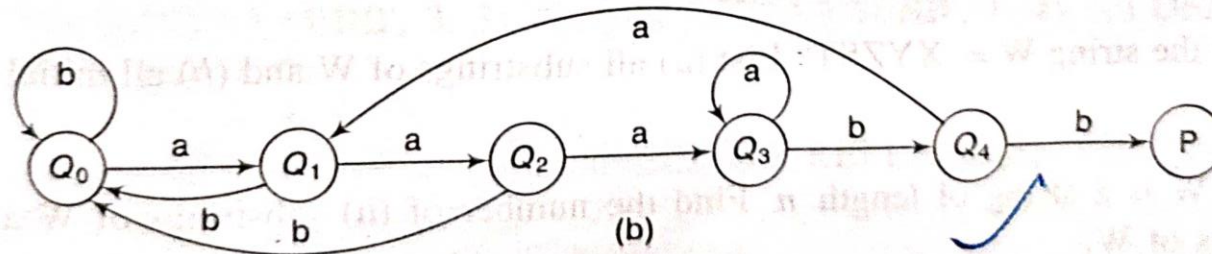
For each character  $t$ , the entry  $f(Q_i, t)$  in the table is the largest  $Q$  which appears as a terminal substring in the string  $Q_i t$ . We compute:

$$\begin{array}{lllll} f(\Lambda, a) = a, & f(a, a) = a^2, & f(a^2, a) = a^3, & f(a^3, a) = a^3, & f(a^3b, a) = a \\ f(\Lambda, b) = \Lambda, & f(a, b) = \Lambda, & f(a^2, b) = \Lambda, & f(a^3, b) = a^3b, & f(a^3b, b) = P \end{array}$$

Hence the required table appears in Fig. 3.10(a). The corresponding graph appears in Fig. 3.10(b), where there is a node corresponding to each  $Q$  and an arrow from  $Q_i$  to  $Q_j$  labeled by the character  $t$  for each entry  $f(Q_i, t) = Q_j$  in the table.

	a	b
$Q_0$	$Q_1$	$Q_0$
$Q_1$	$Q_2$	$Q_0$
$Q_2$	$Q_3$	$Q_0$
$Q_3$	$Q_3$	$Q_4$
$Q_4$	$Q_1$	$P$

(a)



(b)

Find the table and corresponding graph for the second pattern matching algorithm where the pattern is  $P = ababab$ .

The initial substrings of  $P$  are:

$$Q_0 = \Lambda, \quad Q_1 = a, \quad Q_2 = ab, \quad Q_3 = aba, \quad Q_4 = abab, \quad Q_5 = ababa, \quad Q_6 = ababab = P$$

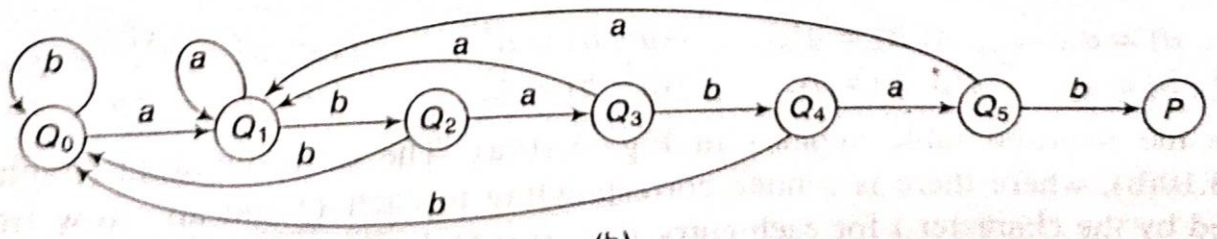
The function  $f$  giving the entries in the table follows:

$f(\Lambda, a) = a$	$f(\Lambda, b) = \Lambda$
$f(a, a) = a$	$f(a, b) = ab$
$f(ab, a) = aba$	$f(ab, b) = \Lambda$
$f(aba, a) = a$	$f(aba, b) = abab$
$f(abab, a) = ababa$	$f(abab, b) = \Lambda$
$f(ababa, a) = a$	$f(ababa, b) = P$

The table appears in Fig. 3.11(a) and the corresponding graph appears in Fig. 3.11(b).

	a	b
$Q_0$	$Q_1$	$Q_0$
$Q_1$	$Q_1$	$Q_2$
$Q_2$	$Q_3$	$Q_0$
$Q_3$	$Q_1$	$Q_4$
$Q_4$	$Q_5$	$Q_0$
$Q_5$	$Q_1$	$P$

(a)



(b)