

string algorithm

\* → string processing problem

Input: Two strings T and P

Problem: find if P is a substring of T

Example - (1)

Input T = gtgatcagatcact P = tcg

output: yes shift = 4, 9

Example - (2)

Input T = 189342670893, P = 1673

Output: No

Solution: (Naive Algorithm) -:

Naive algorithm (T, P)

{ suppose  $n = \text{length}(T)$ ,  $m = \text{length}(P)$ ;

for shift  $s = 0$  to  $n - m$  do

if ( $P[1..m] == T[s+1..s+m]$ ) then

print shift  $s$ ;

end algo.

Complexity:  $O((n - m + 1)m)$

①

Knuth-Morris-Pratt algorithm -:

KMP flow chart construction -:

Input: P, a string of characters; m, the length of P.

Output: fail, The array of failure links, defined for index 1, ..., m. The array is passed in and the algorithm fills it

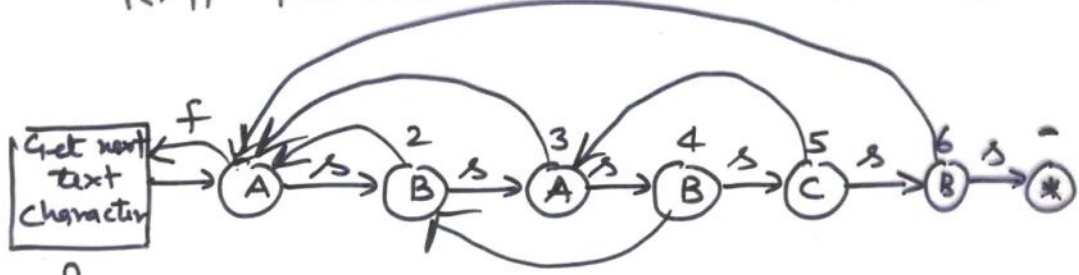
```
void KmpSetup(char[] P, int m, int[] fail)
{
    int k, s;
    fail[1] = 0;
    for (k = 2; k <= m; k++)
        s = fail[k-1];
        while (s >= 1)
            if (P[s] == P[k-1])
                break;
        s = fail[s];
    fail[k] = s+1;
}
```

---

Complexity  $O(m^2)$

Example:

KMP flowchart for  $P = \text{'A B A B C B'}$



Note:

\*  $fail[1] = 0$

\*  $fail[2]$   
 $s = fail[1] = 0$   
 $\Rightarrow fail[2] = 1$

\*  $fail[3]$   
 $s = fail[2] = 1$   
 $\because p_2 \neq p_1$   
 $s_1 = fail[1] = 0$   
 $\Rightarrow fail[3] = 1$

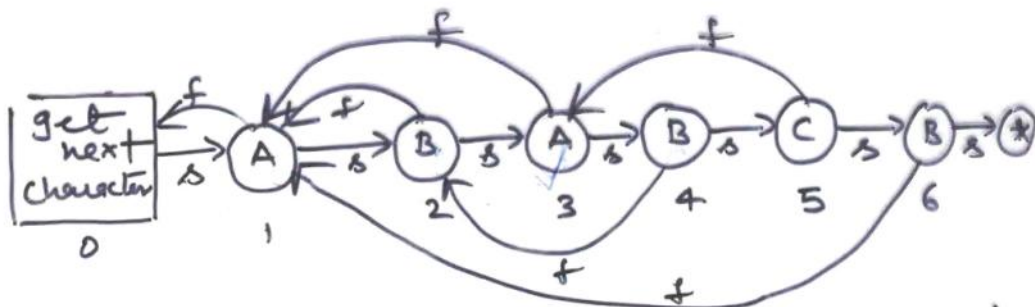
\*  $fail[4]$   
 $s = fail[3] = 1$   
 $\because p_3 = p_1$   
 $\Rightarrow fail[4] = 2$

\*  $fail[5]$   
 $s = fail[4] = 2$   
 $p_4 = p_2$   
 $\Rightarrow fail[5] = 3$

$fail[6]$

$s = fail[5] = 3$   
 $p_5 \neq p_3$   
 $\Rightarrow s_1 = fail[3] = 1$   
 $p_5 \neq p_1$   
 $\Rightarrow s_2 = fail[2] = 0$   
 $\Rightarrow fail[6] = 1$

(3)



KMP- flow chart for p = 'ABABCB'

Action of the KMP flow chart for the pattern ABABCB on the text last line

ACABAA BABA  
 1 2 3 4 5 6 7 8 9 10

} 4 // none fail ✓

KMP cell No	Text being processed		Success or failure
	Index	Character	
1	1	A	s
2	2	C	f
1	2	C	f
1	3	A	s
2	4	B	s
3	5	A	s
4	6	A	f
2	6	A	f
1	6	A	s
2	7	B	s
3	8	A	s
4	9	B	f
1	10	A	s

## Algorithm KMP scan

Input: P and T  
           $\xrightarrow{\text{P}}$        $\xrightarrow{\text{T}}$   
          pattern    text string

m = length of P,

fail: array of failure links

Output: The return value is the index in T where a copy of P begins or -1

```
int kmpscan(char[] P, char[] T, int m, int[] fail)
{
    int match;
    int j, k;
    // j indexes text and k indexes pattern & fail array
    match = -1; j = 1; k = 1
    while (endText(T, j) == fail[k])
        if (k > m)
            match = j - m // match found
            break;
        if (k == 0)
            j++;
            k = 1 // start pattern over
        else if (T[j] == P[k])
            j++;
            k++;
        else
            k = fail[k];
            // continue loop
    return match
}
```

$\Theta(m+n)$

(5)

Exercise:- Draw KMP flow chart for  
the pattern

A B A B A B C B