

# Problems in String Matching and Pattern Recognition

Dr. Gur Saran Adhar

Reference clrs, Chapter 32, Page 906-  
Robert Sedgewick, Chapter 19  
Udi Manber Page 148-

+

+

## Problem Statement: String Matching

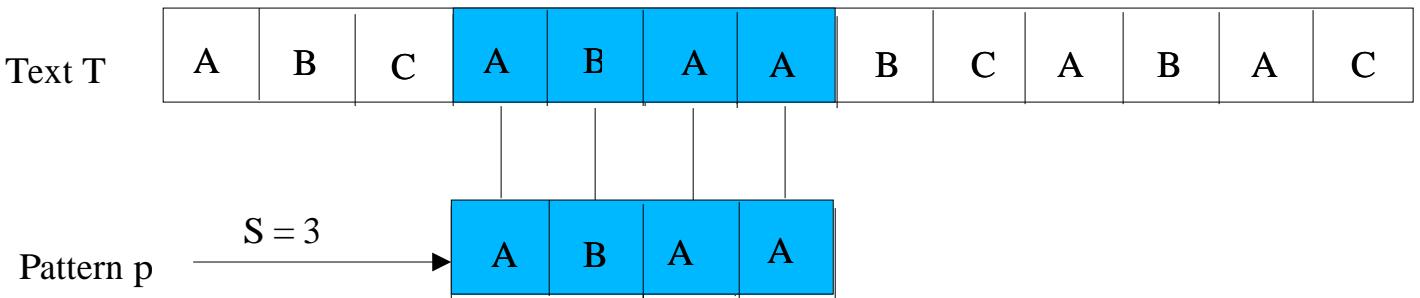
Given a **text**  $T[1 \dots n]$  of length  $n$  and a **Pattern**  $P[1 \dots m]$  of length  $m \leq n$ , where the elements of  $T$  and  $P$  are drawn from an alphabet  $\Sigma$ , the **string matching problem** is to find an occurrence of the pattern within the text.

+

1

+

+

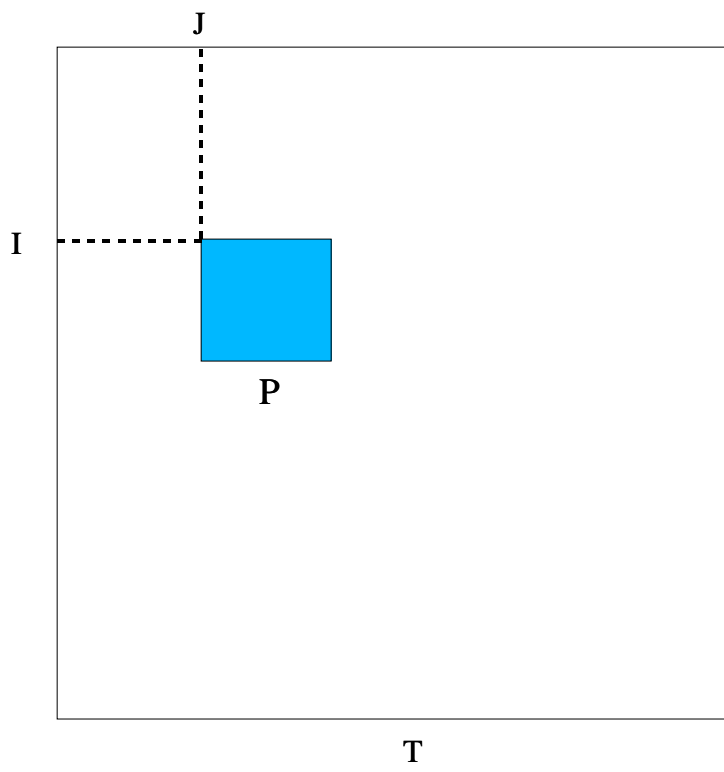


+

2

+

+



+

3

+

+

## List Of Common String Matching Algorithms

1.	Naive String-Matching Algorithm	clrs 909-
2.	Rabin Karp Algorithm	clrs 911-
3.	String Matching with Finite Automata	clrs 916-
4.	Knuth-Morris-Pratt (KMP) Algorithm	clrs 923-
5.	Boyer-Moore (BM) Algorithm	Sedgewick 28

+

4

+

+

## Notation: String Matching

**Prefix:**  $w$  is a **prefix** of string  $x$ , denoted  $w \sqsubseteq x$ , if  $x = wy$  for some string  $y \in \Sigma^*$ .

**Suffix:**  $w$  is a **suffix** of string  $x$ , denoted  $w \sqsupseteq x$ , if  $x = yw$  for some string  $y \in \Sigma^*$ .

+

5

+

+

### **Naive String Matcher**( $T, P$ )

1  $n \leftarrow \text{length}[T]$

2  $m \leftarrow \text{length}[P]$

3 **for**  $s \leftarrow 0$  to  $n - m$

4     **do if**  $P[1 \dots m] = T[s + 1 \dots s + m]$

5         **then** print "Pattern occurs with shift"  $s$

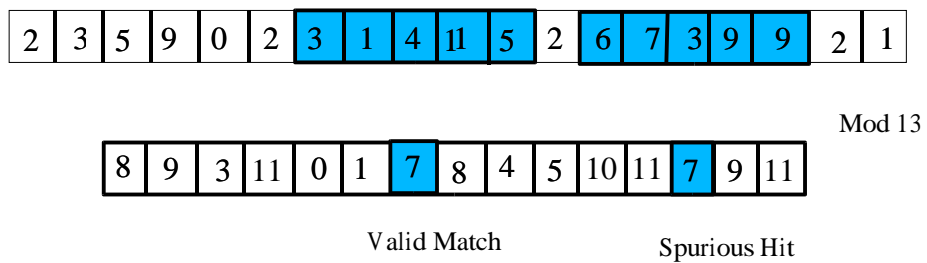
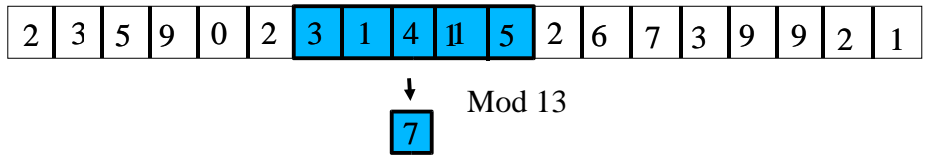
Reference clrs 909

+

6

+

+



3	1	4	1	5	2
---	---	---	---	---	---

$14152 = (31415 \sim 3.10000).10+2 \pmod{13}$

7	8
---	---

+

7



+

+

//  $p$  denotes the decimal value of the pattern  $P[1 \dots m]$ .  
 //  $t_s$  denotes the decimal value of the  $m$ -length substring  
 //  $T[1 + s \dots s + m]$ .

### **Rabin-Karp-Matcher**( $T, P, d, q$ )

```

1   $n \leftarrow \text{length}[T]$ 
2   $m \leftarrow \text{length}[P]$ 
3   $h = d^{m-1} \text{mod} q$ 
4   $p \leftarrow 0$ 
5   $t_0 \leftarrow 0$ 
6  for  $i \leftarrow 1$  to  $m$ 
7      do  $p \leftarrow (dp + P[i]) \text{mod} q$ 
8           $t_0 \leftarrow (dt_0 + T[i]) \text{mod} q$ 
9  for  $s \leftarrow 0$  to  $n - m$ 
10     do if  $p = t_s$ 
11         then if  $P[1 \dots m] = T[s + 1 \dots s + m]$ 
12             then print "Pattern occurs with shift"  $s$ 
13     if  $s < n - m$ 
14         then  $t_{s+1} \leftarrow (d(t_s - T[s + 1])h + T[s + m + 1]) \text{mod} q$ 

```

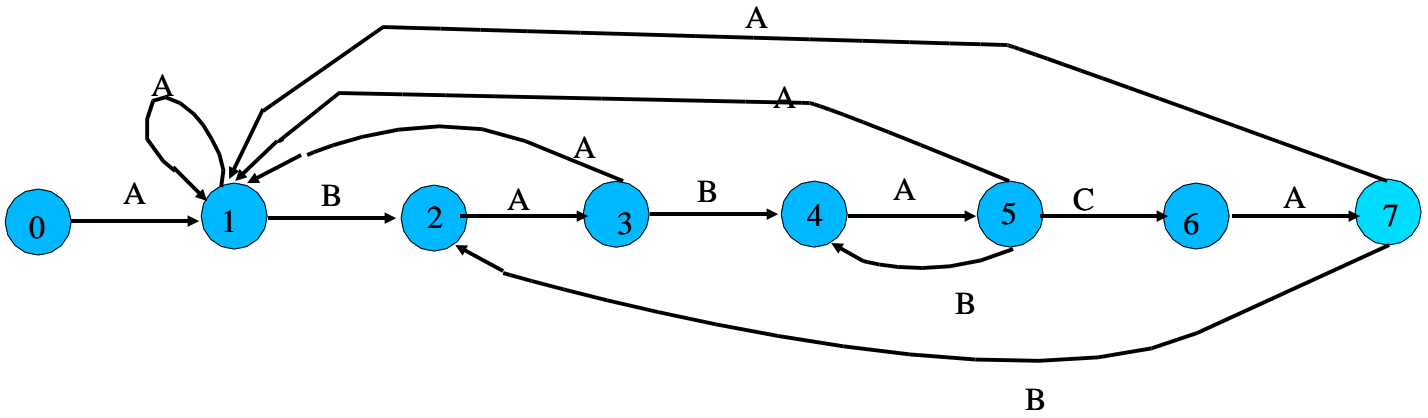
Reference clrs 909

+

8

+

+



	INPUT		
	A	B	C
0	1	0	0
1	1	2	0
2	3	0	0
3	1	4	0
4	5	0	0
5	1	4	6
6	7	0	0
7	1	2	0

P -the pattern

A

B

A

B

A

C

A

i—	1	2	3	4	5	6	7	8	9	10	11	
I[i]—	A	B	<u>A</u>	<u>B</u>	<u>A</u>	<u>B</u>	<u>A</u>	C	A	B	A	
State—	0	1	2	3	4	5	4	5	6	<b>7</b>	2	3

+

+

+

**Suffix Function:**  $\sigma$  for a pattern  $P[1\dots m]$  is the *length* of the longest prefix of  $P$  that is a suffix of  $x$ :

$$\sigma(x) = \max\{k : P_k \sqsupseteq x\}$$

The function  $\sigma$  is a mapping from  $\Sigma^*$  to  $\{1, 2, \dots, m\}$  such that  $\sigma(x)$  is the length of the prefix of  $P$  which is a suffix of  $x$ .

Observe: In the example Finite Automaton  $\delta(5, b) = 4$ , this transition is made because if the automaton reads a  $b$  in state  $q = 5$ , then  $P_q b = ababab$ , and the longest **prefix of  $P = ababab$  which is also a suffix of  $ababab$  is  $P_4 = abab$** .

The transition function  $\delta$  is defined by the following equation, for any state  $q$  and any character  $a$ .

$$\delta(q, a) = \sigma(P_q a)$$

+

+

+

### **Finite-Automaton-Matcher**( $T, \delta, m$ )

1  $n \leftarrow \text{length}[T]$

2  $q \leftarrow 0$

3 **for**  $i \leftarrow 1$  to  $n$

4     **do**  $q \leftarrow \delta(q, T[i])$

5         **then if**  $q = m$

6             **then** print "Pattern occurs with shift"  $i - m$

Reference clrs 919

+

+

+

### COMPUTE\_TRANSITION\_FUNCTION( $P, \Sigma$ )

```
1  $m \leftarrow \text{length}[P]$ 
2 for  $q \leftarrow 0$  to  $m$ 
3   do for each character  $a \in \Sigma$ 
4     do  $k \leftarrow \min(m + 1, q + 2)$ 
5       repeat  $k \leftarrow k - 1$ 
6         until  $P_k \sqsupseteq P_q a$ 
7          $\delta(q, a) \leftarrow k$ 
8 return  $\delta$ 
```

Reference clrs 922

+

+

+

## Observation: Overlapping suffixes

Suppose  $x$ ,  $y$ , and  $z$  are strings such that,  $x \sqsupseteq z$  and  $y \sqsupseteq z$ . That is, both  $x$  and  $y$  are suffix of  $z$ .

If  $|x| \leq |y|$  then  $x \sqsupseteq y$ .

That is,  $x$  is suffix of  $y$

If  $|x| \geq |y|$  then  $y \sqsupseteq x$ .

That is,  $y$  is suffix of  $x$

If  $|x| = |y|$  then  $x = y$ .

That is,  $x$  is equal to  $y$

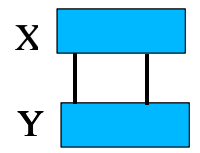
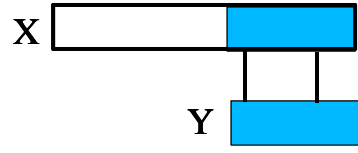
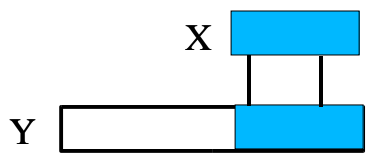
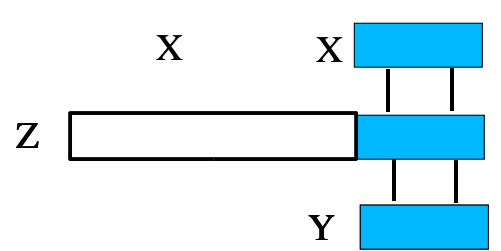
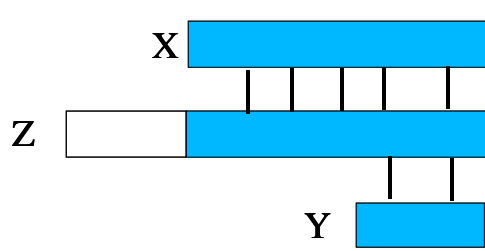
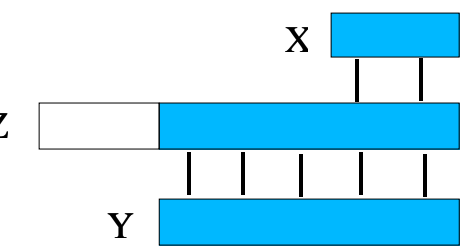
Also note, for any strings  $x$  and  $y$  and any character  $a$  if  $x \sqsupseteq y$  then  $xa \sqsupseteq ya$ .

That is if  $x$  is suffix of  $y$  then so is  $xa$  of  $ya$ .

+

+

+



+

+

+

## Knuth Morris Pratt (KMP) Algorithm Main Idea

### Prefix Function:

Prefix function  $\pi$  encapsulates the knowledge about how a pattern matches against shifts of itself. This avoids useless shifts in the naive algorithm.

$\pi[5]$  for example gives the *next potentially valid shift* given that five characters have matched (and the sixth does not match).

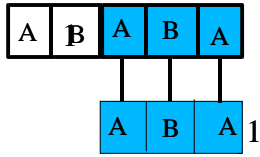
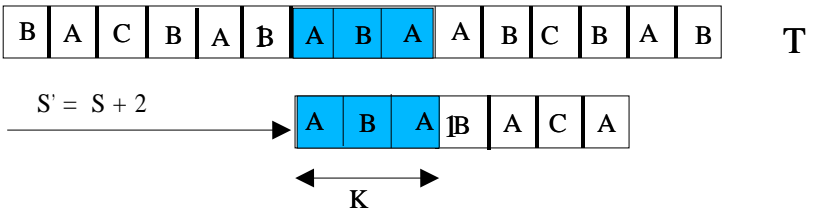
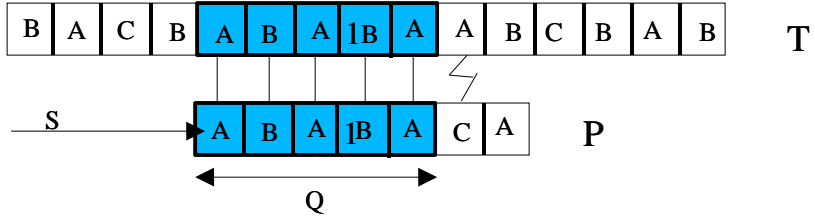
Reference clrs 923-

+



+

+



+

+

+

**Prefix Function:**  $\pi$  for a pattern  $P[1\dots m]$  is the function  $\pi : \{1, 2, \dots, m\} \rightarrow \{0, 1, 2, \dots, m-1\}$  such that:

$$\pi(q) = \max\{k : k < q \text{ and } P_k \sqsupseteq P_q\}$$

That is  $\pi[q]$  is the length of the longest prefix of  $P$  that is a proper suffix of  $P_q$ .

+

+

+

### **KMP-Matcher**( $T, P$ )

```
1  $n \leftarrow \text{length}[T]$ 
2  $m \leftarrow \text{length}[P]$ 
3  $\pi \leftarrow \text{COMPUTE\_PREFIX\_FUNCTION}(P)$ 
4  $q \leftarrow 0$ 
5 for  $i \leftarrow 1$  to  $n$ 
6     do while  $q > 0$  and  $P[q + 1] \neq T[i]$ 
7         do  $q \leftarrow \pi[q]$ 
8     if  $P[q + 1] = T[i]$ 
9         then  $q \leftarrow q + 1$ 
10    if  $q = m$ 
11        then print "Pattern occurs with shift"  $i - m$ 
12         $q \leftarrow \pi[q]$ 
```

Reference clrs 926

+

+

+

### COMPUTE\_PREFIX\_FUNCTION( $P$ )

```
1  $m \leftarrow \text{length}[P]$ 
2  $\pi[1] \leftarrow 0$ 
3  $k \leftarrow 0$ 
4 for  $q \leftarrow 2$  to  $m$ 
5     do while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6         do  $k \leftarrow \pi[k]$ 
7     if  $P[k + 1] = P[q]$ 
8         then  $k \leftarrow k + 1$ 
9      $\pi[q] = k$ 
10 return  $\pi$ 
```

Reference clrs 926

+