



## طراحی و تحلیل الگوریتم‌ها

مبحث دوازدهم

مطالعه‌ی موضوعی الگوریتم‌ها

تطابق رشته‌ها

String Matching

کاظم فولادی قلعه

دانشکده مهندسی، دانشکدگان فارابی

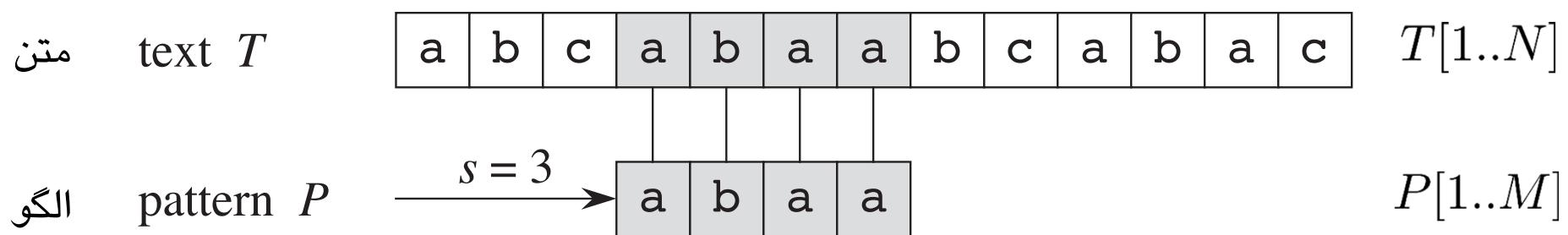
دانشگاه تهران

<http://courses.fouladi.ir/algorithm>

## تطابق رشته‌ها

STRING MATCHING

تطابق رشته‌ها: یافتن همهٔ موارد وقوع یک الگو در یک متن



An example of the string-matching problem, where we want to find all occurrences of the pattern  $P = \text{abaa}$  in the text  $T = \text{abcabaabcabac}$ . The pattern occurs only once in the text, at shift  $s = 3$ , which we call a valid shift. A vertical line connects each character of the pattern to its matching character in the text, and all matched characters are shaded.

## تطابق رشته‌ها

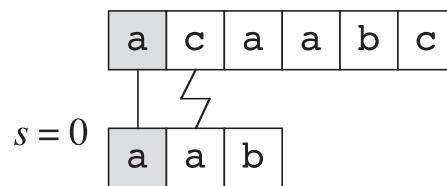
مهم‌ترین الگوریتم‌های موجود

STRING MATCHING

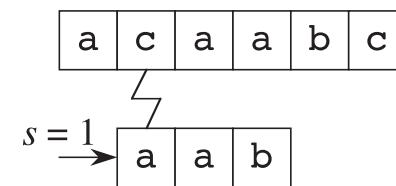
Algorithm	Preprocessing time	Matching time
Naive	0	$O((n - m + 1)m)$
Rabin-Karp	$\Theta(m)$	$O((n - m + 1)m)$
Finite automaton	$O(m  \Sigma )$	$\Theta(n)$
Knuth-Morris-Pratt	$\Theta(m)$	$\Theta(n)$

## تطابق رشته‌ها

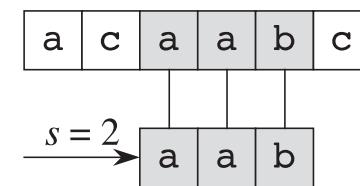
الگوریتم ساده: مثال

NAIVE ALGORITHM

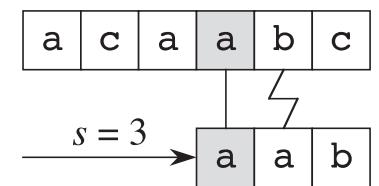
(a)



(b)



(c)



(d)

The operation of the naive string matcher for the pattern  $P = \text{aab}$  and the text  $T = \text{acaabc}$ . We can imagine the pattern  $P$  as a template that we slide next to the text. (a)–(d) The four successive alignments tried by the naive string matcher. In each part, vertical lines connect corresponding regions found to match (shown shaded), and a jagged line connects the first mismatched character found, if any. The algorithm finds one occurrence of the pattern, at shift  $s = 2$ , shown in part (c).

## تطابق رشته‌ها

الگوریتم جستجوی کورکورانه

### BRUTE-FORCE ALGORITHM

الگوریتم کورکورانه، الگو را با متن مقایسه می‌کند  
(هر بار یک کاراکتر)  
تا یک کاراکتر نامطابق پیدا شود.

## تطابق رشته‌ها

الگوریتم جستجوی کورکورانه: مثال

TWO ROADS DIVERGED IN A YELLOW WOOD  
ROADS

## تطابق رشته‌ها

الگوریتم جستجوی کورکورانه: شبکه

**do**

**if** (text letter == pattern letter)

    compare next letter of pattern to next  
    letter of text

**else**

    move pattern down text by one letter

**while** (entire pattern found or end of text)

## تطابق رشته‌ها

الگوریتم جستجوی کورکورانه: مثال

tetththeheehttehttheththehehtht  
**the**

tetththeheehttehttheththehehtht  
**the**

te**t**ththeheehttehttheththehehtht  
**the**

tet **th**ttheheehttehttheththehehtht  
**the**

tetth**t**theheehttehttheththehehtht  
**the**

tetth**the**heehttehttheththehehtht  
**the**

## تطابق رشته‌ها

الگوریتم جستجوی کورکورانه: زمان اجرا

$$O(N - M + 1)$$

تعداد مقایسه‌ها:

$$O(NM)$$

زمان اجرای بدترین حالت:

## تطابق رشته‌ها

الگوریتم جستجوی کورکورانه: زمان اجرا: مثال

- 1) *AAAAA*  
*AAAAH*      **5 comparisons made**
  - 2) *AAAAA*  
*AAAAH*      **5 comparisons made**
  - 3) *AAAAA*  
*AAAAH*      **5 comparisons made**
  - 4) *AAAAA*  
*AAAAH*      **5 comparisons made**
  - 5) *AAAAA*  
*AAAAH*      **5 comparisons made**
- ....
- N) *AAAAAAAAAAAAAA*  
**5 comparisons made**      *AAAAH*

## تطابق رشته‌ها

الگوریتم کورکورانه

**NAIVE-STRING-MATCHER( $T, P$ )**

```

1    $n = T.length$ 
2    $m = P.length$ 
3   for  $s = 0$  to  $n - m$ 
4       if  $P[1..m] == T[s + 1..s + m]$ 
5           print “Pattern occurs with shift”  $s$ 

```

الگوریتم کورکورانه، الگو را با متن مقایسه می‌کند  
 (هر بار یک کاراکتر)  
 تا یک کاراکتر نامطابق پیدا شود.

## تطابق رشته‌ها: الگوریتم رابین - کارپ

### RABIN-KARP ALGORITHM

الگوریتم جستجوی رشته‌ی رابین - کارپ،  
برای الگو

و

برای هر زیردنباله‌ی  $M$  کاراکتری (که باید مقایسه شود)  
یک مقدار هش را محاسبه می‌کند.

اگر مقادیر هش مساوی نبودند،

الگوریتم مقدار هش را برای زیردنباله‌ی  $M$  کاراکتری بعدی محاسبه می‌کند.

اگر مقادیر هش مساوی بودند،

الگوریتم بین الگو و زیردنباله‌ی  $M$  کاراکتری مقایسه‌ی کورکورانه انجام می‌دهد.

## تطابق رشته‌ها: الگوریتم رابین - کارپ

مثال

Hash value of “AAAAAA” is 37

Hash value of “AAAAAH” is 100

1) **AAAAAA**AAAAAAAAAAAAAAA  
AAAAH

$37 \neq 100$     **1 comparison made**

2) A**AAAAAA**AAAAAAAAAAAAAAA  
AAAAH

$37 \neq 100$     **1 comparison made**

3) AA**AAAAA**AAAAAAA  
AAAAH

$37 \neq 100$     **1 comparison made**

4) AAA**AAAAA**AAAAAAA  
AAAAH

$37 \neq 100$     **1 comparison made**

...

N) AAAAAAAA  
AAAAH

**6 comparisons made**

$100 = 100$

## تطابق رشته‌ها: الگوریتم رابین - کارپ

شبه کد

*pattern is M characters long***hash\_p**=hash value of pattern**hash\_t**=hash value of first M letters in  
body of text**do****if (hash\_p == hash\_t)**brute force comparison of pattern  
and selected section of text**hash\_t** = hash value of next section of  
text, one character over**while** (end of text **or**

brute force comparison == true)

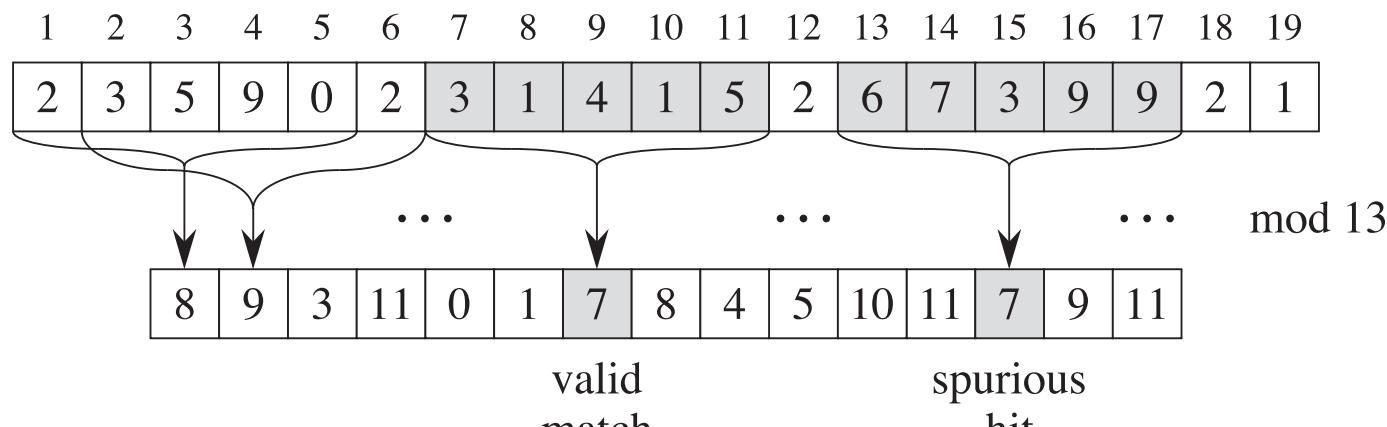
## تطابق رشته‌ها: الگوریتم رابین - کارپ

مثال: محاسبه و مقایسه هش‌ها

 $\text{mod } 13$ 

7

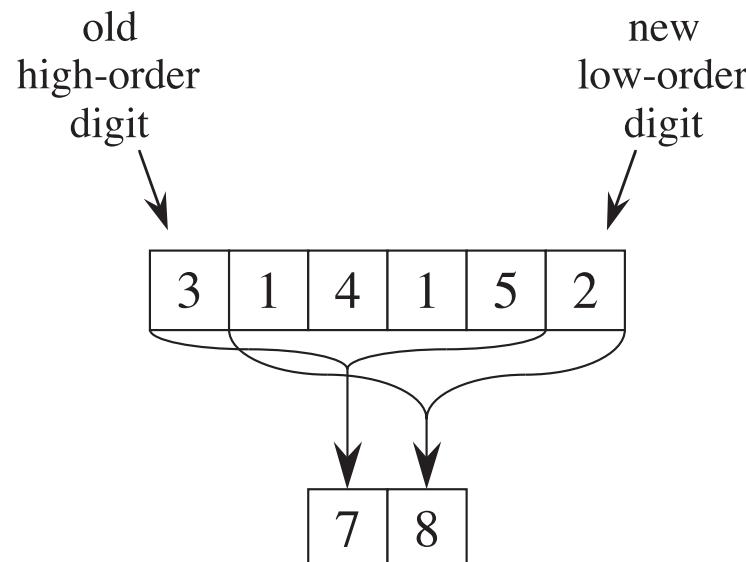
(a)



(b)

## تطابق رشته‌ها: الگوریتم رابین - کارپ

محاسبه‌ی مقدار هش



(c)

$$\begin{aligned}
 14152 &\equiv (31415 - 3 \cdot 10000) \cdot 10 + 2 \pmod{13} \\
 &\equiv (7 - 3 \cdot 3) \cdot 10 + 2 \pmod{13} \\
 &\equiv 8 \pmod{13}
 \end{aligned}$$

## تطابق رشته‌ها: الگوریتم رابین - کارپ

شبه‌کد

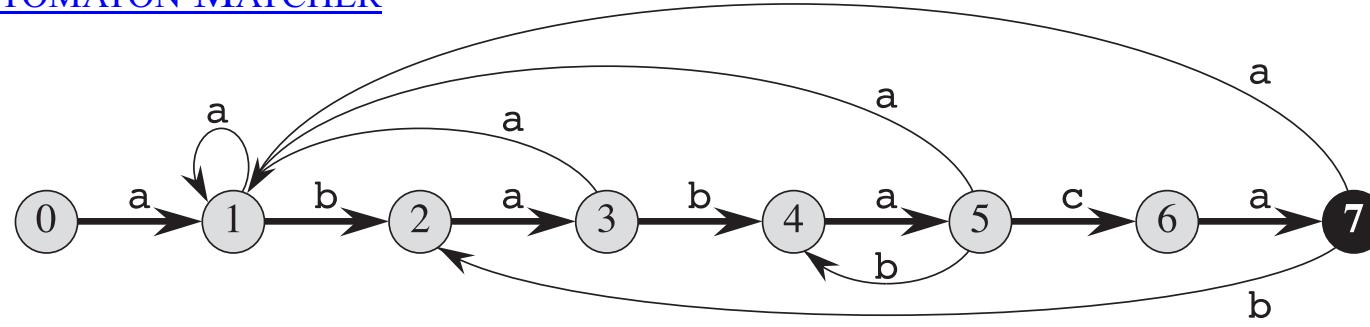
RABIN-KARP-MATCHER( $T, P, d, q$ )

```

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

```

## تطابق رشته‌ها: الگوریتم تطابق با آutomaton متناهی

FINITE AUTOMATON MATCHER

(a)

state	input			$P$
	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

(b)

$i$	—	1	2	3	4	5	6	7	8	9	10	11
$T[i]$	—	a	b	a	b	a	c	a	b	a	b	
state $\phi(T_i)$	0	1	2	3	4	5	4	5	6	7	2	3

(c)

## تطابق رشته‌ها: الگوریتم تطابق با آutomaton متناهی

شبہ کد

FINITE AUTOMATON MATCHERFINITE-AUTOMATON-MATCHER( $T, \delta, m$ )

```

1    $n = T.length$ 
2    $q = 0$ 
3   for  $i = 1$  to  $n$ 
4        $q = \delta(q, T[i])$ 
5       if  $q == m$ 
6           print “Pattern occurs with shift”  $i - m$ 

```

## تطابق رشته‌ها: الگوریتم کنوث - موریس - پرت

### KNUTH-MORRIS-PRATT ALGORITHM (KMP)

الگوریتم KMP از ردگیری اطلاعات حاصل از مقایسه‌های قبلی استفاده می‌کند:  
یک تابع پیشوند  $\pi$  محاسبه می‌شود  
تا نشان دهد در صورت شکست، چه میزان از مقایسه‌ی قبلی می‌تواند استفاده‌ی مجدد شود.

## تطابق رشته‌ها: الگوریتم کنوث - موریس - پرت

تابع پیشوند

### KNUTH-MORRIS-PRATT ALGORITHM (KMP)

#### تابع پیشوند $\pi$

$\pi(i)$  = طول بلندترین پیشوند الگوی  $P[1..i]$  که پسوند  $P[2..i]$  باشد.

نشان می‌دهد که

چه میزان از ابتدای رشته

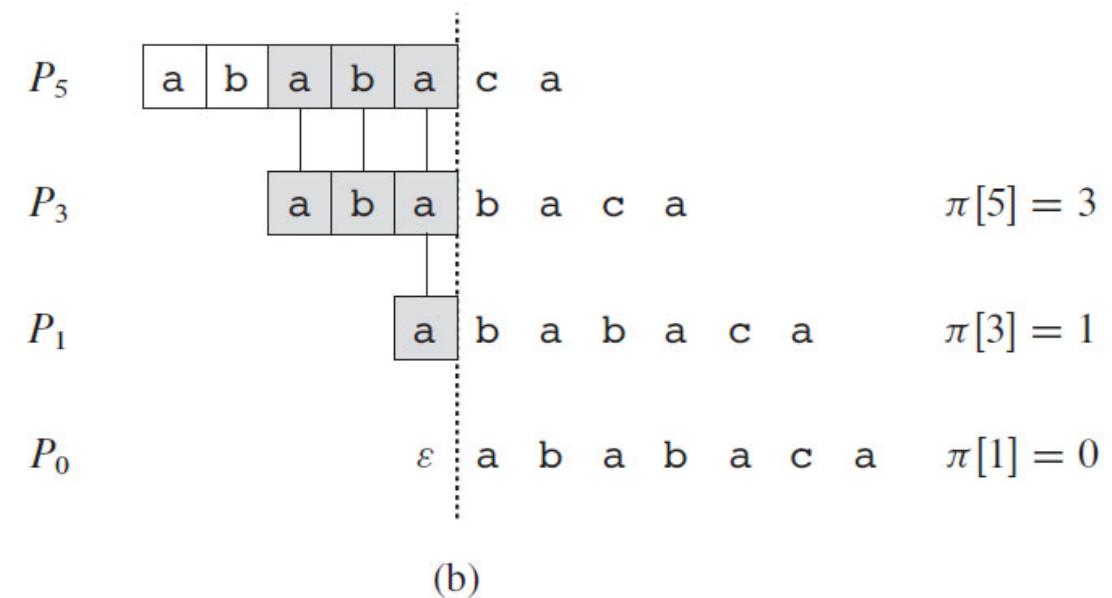
با بخشی که بلا فاصله قبل از مقایسه شکست‌خورده قرار دارد،  
تطابق می‌یابد.

## تطابق رشته‌ها: الگوریتم کنوث - موریس - پرت

KNUTH-MORRIS-PRATT ALGORITHM (KMP)

$i$	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	c	a
$\pi[i]$	0	0	1	2	3	0	1

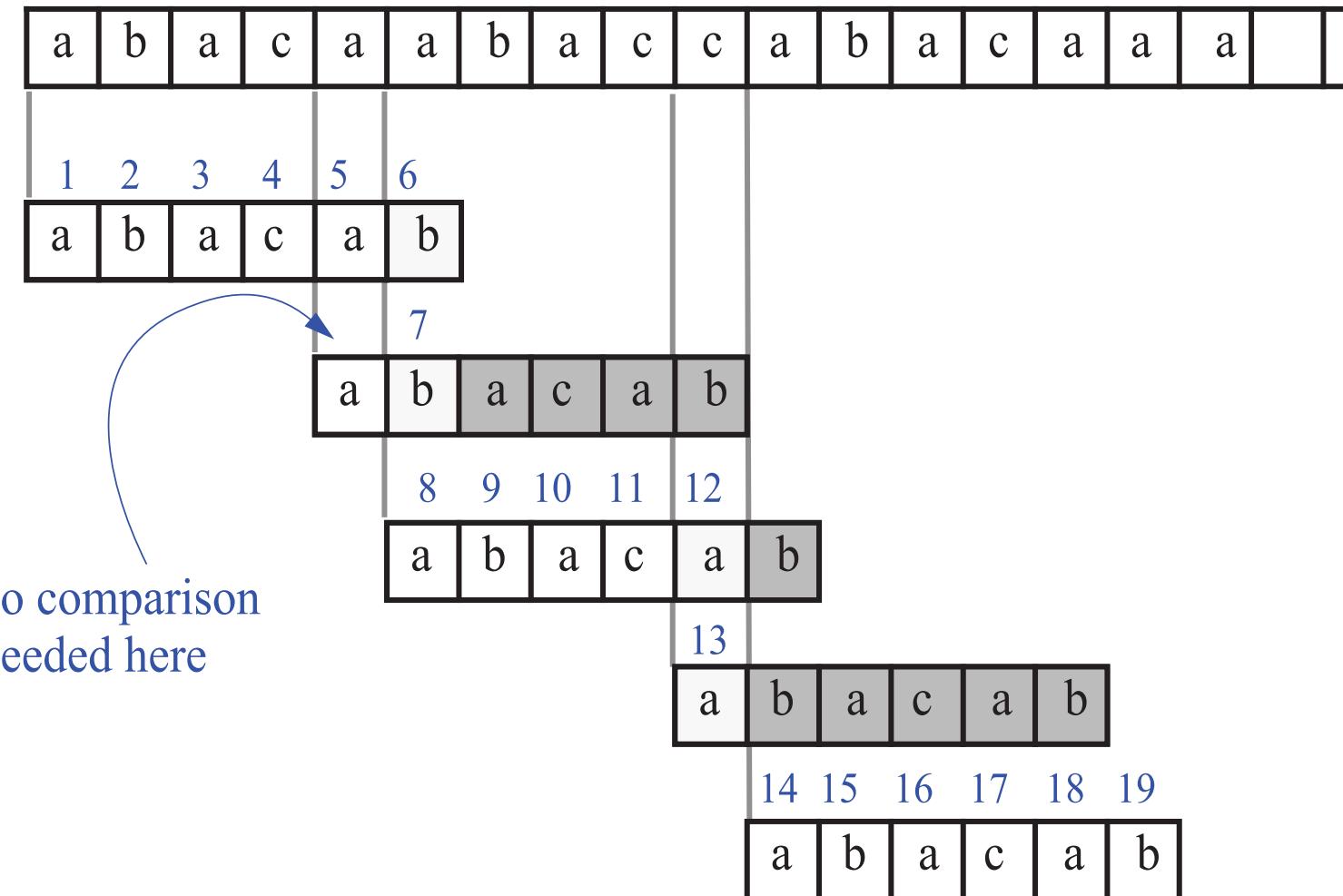
(a)



(b)

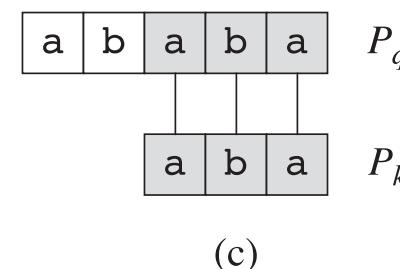
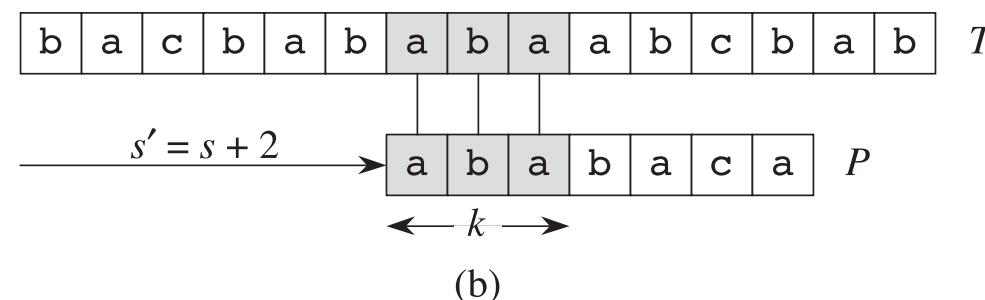
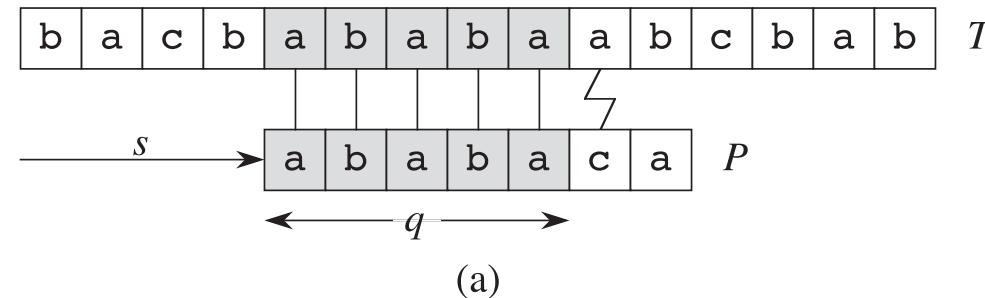
## تطابق رشته‌ها: الگوریتم کنوث - موریس - پرت

مثال

KNUTH-MORRIS-PRATT ALGORITHM (KMP)

## تطابق رشته‌ها: الگوریتم کنوث - موریس - پرت

مثال

KNUTH-MORRIS-PRATT ALGORITHM (KMP)

## تطابق رشته‌ها: الگوریتم کنوث - موریس - پرت

شبہ کد

KNUTH-MORRIS-PRATT ALGORITHM (KMP)KMP-MATCHER( $T, P$ )

```

1   $n = T.length$ 
2   $m = P.length$ 
3   $\pi = \text{COMPUTE-PREFIX-FUNCTION}(P)$ 
4   $q = 0$                                 // number of characters matched
5  for  $i = 1$  to  $n$                   // scan the text from left to right
6    while  $q > 0$  and  $P[q + 1] \neq T[i]$ 
7       $q = \pi[q]$                       // next character does not match
8      if  $P[q + 1] == T[i]$ 
9         $q = q + 1$                     // next character matches
10     if  $q == m$                       // is all of  $P$  matched?
11       print "Pattern occurs with shift"  $i - m$ 
12        $q = \pi[q]$                   // look for the next match

```

## تطابق رشته‌ها: الگوریتم کنوث - موریس - پرت

شبه‌کد:تابع محاسبه‌ی پیشوند

### KNUTH-MORRIS-PRATT ALGORITHM (KMP)

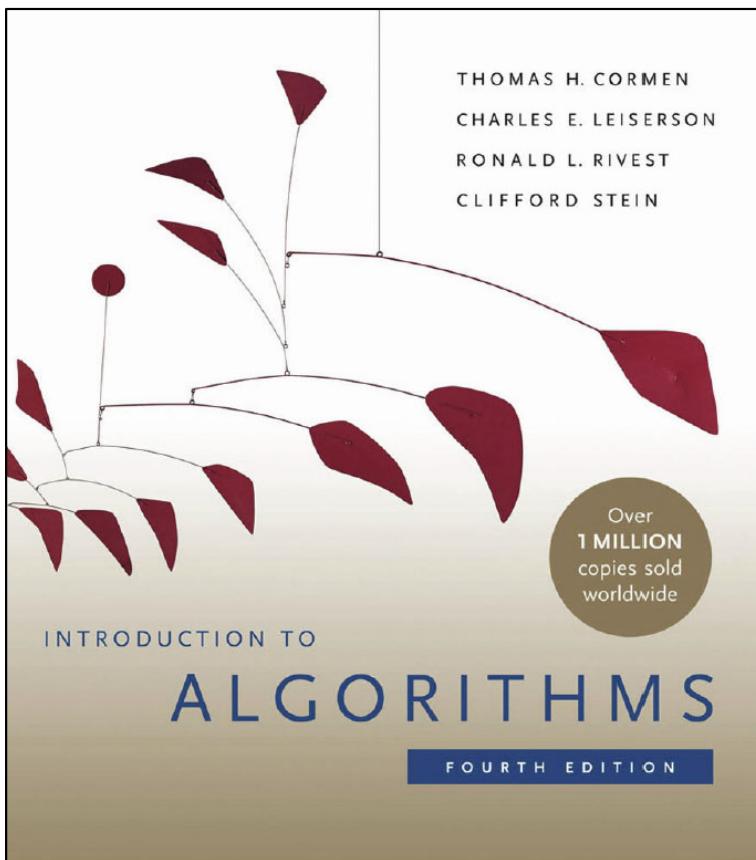
#### COMPUTE-PREFIX-FUNCTION( $P$ )

```

1    $m = P.length$ 
2   let  $\pi[1..m]$  be a new array
3    $\pi[1] = 0$ 
4    $k = 0$ 
5   for  $q = 2$  to  $m$ 
6       while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
7            $k = \pi[k]$ 
8           if  $P[k + 1] == P[q]$ 
9                $k = k + 1$ 
10           $\pi[q] = k$ 
11      return  $\pi$ 

```

## مرجع



T.H. Cormen, C.E. Leiserson, R.L. Rivest, C. Stein,  
**Introduction to Algorithms**,  
 4<sup>th</sup> Edition, MIT Press, 2022.

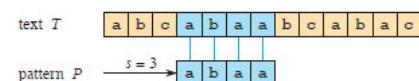
## Chapter 32

### 32 String Matching

Text-editing programs frequently need to find all occurrences of a pattern in the text. Typically, the text is a document being edited, and the pattern searched for is a particular word supplied by the user. Efficient algorithms for this problem —called “string matching”—can greatly aid the responsiveness of the text-editing program. Among their many other applications, string-matching algorithms search for particular patterns in DNA sequences. Internet search engines also use them to find web pages relevant to queries.

The string-matching problem can be stated formally as follows. The text is given as an array  $T[1 : n]$  of length  $n$ , and the pattern is an array  $P[1 : m]$  of length  $m \leq n$ . The elements of  $P$  and  $T$  are characters drawn from an alphabet  $\Sigma$ , which is a finite set of characters. For example,  $\Sigma$  could be the set  $\{0, 1\}$ , or it could be the set  $\{a, b, \dots, z\}$ . The character arrays  $P$  and  $T$  are often called *strings* of characters.

As Figure 32.1 shows, pattern  $P$  occurs with shift  $s$  in text  $T$  (or, equivalently, that pattern  $P$  occurs beginning at position  $s + 1$  in text  $T$ ) if  $0 \leq s \leq n - m$  and  $T[s + 1 : s + m] = P[1 : m]$ , that is, if  $T[s + j] = P[j]$ , for  $1 \leq j \leq m$ . If  $P$  occurs with shift  $s$  in  $T$ , then  $s$  is a *valid shift*, and otherwise,  $s$  is an *invalid shift*. The *string-matching problem* is the problem of finding all valid shifts with which a given pattern  $P$  occurs in a given text  $T$ .



**Figure 32.1** An example of the string-matching problem to find all occurrences of the pattern  $P = abaa$  in the text  $T = abcabaabcaabac$ . The pattern occurs only once in the text, at shift  $s = 3$ , which is a valid shift. A vertical line connects each character of the pattern to its matching character in the text, and all matched characters are shaded blue.