# Chapter 32

# String Matching

songyou@buaa.edu.cn

# VII  Selected Topics

# 32　String Matching

**char  \*strstr(char \*text, char \*pattern);**

string-match

Previous　　Next

### 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.

We formalize the string-matching problem as follows. We assume that the text is an array $T[1..n]$ of length $n$ and that the pattern is an array $P[1..m]$ of length $m \leq n$. We further assume that the elements of $P$ and $T$ are characters drawn from a finite alphabet $\Sigma$. For example, we may have $\Sigma = \{0,1\}$ or $\Sigma = \{a,b,\dots,z\}$. The character arrays $P$ and $T$ are often called **strings** of characters.

Referring to Figure 32.1, we say that 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 we call $s$ a **valid shift**; otherwise, we call $s$ 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 $T$.

**Finding all occurrences of a pattern in a text is a problem that arises frequently in text-editing programs.**

algorithm

上一个　　下一个

Contents

# 32  String Matching

**char  *strstr(char *text, char *pattern);**

- **Efficient algorithms for this problem can greatly aid the responsiveness(响应性) of the text-editing program.**

- **Applications: String-matching algorithms are also used, for example, to search for particular patterns in DNA sequences.**

**char *strstr(char *text, char *pattern);**

Find (1/53) ×

string-match ⚙

Previous    Next

**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.

We formalize the string-matching problem as follows. We assume that the text is an array $T[1 . . n]$ of length $n$ and that the pattern is an array $P[1 . . m]$ of length $m \le n$. We further assume that the elements of $P$ and $T$ are characters drawn from a finite alphabet $\Sigma$. For example, we may have $\Sigma = \{0, 1\}$ or $\Sigma = \{a, b, \ldots, z\}$. The character arrays $P$ and $T$ are often called *strings* of characters.

Referring to Figure 32.1, we say that pattern $P$ *occurs with shift s* in text $T$ (or, equivalently, that pattern $P$ *occurs beginning at position s + 1* in text $T$) if $0 \le s \le n - m$ and $T[s+1 . . s+m] = P[1 . . m]$ (that is, if $T[s+j] = P[j]$, for $1 \le j \le m$). If $P$ occurs with shift $s$ in $T$, then we call $s$ a *valid shift*; otherwise, we call $s$ 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 $T$.

The pattern occurs only once in the text, at shift $s = 3$. The shift $s = 3$ is said to be a valid shift.

## String-matching problem:

- **Text: $T[1 .. n]$, Pattern: $P[1 .. m]$, $m \leq n$.**

- **Finite alphabet: $\Sigma$, for example, $\Sigma = \{0, 1\}$ or $\Sigma = \{a, b, \ldots, z\}$.**

- **$P_i \in \Sigma$, $T_i \in \Sigma$.**

- **$P$ occurs with shift $s$ in $T$ if $0 \leq s \leq n\text{-}m$ and $T[s+1 .. s+m] = P[1 .. m]$ (that is, if $T[s+j] = P[j]$, for $1 \leq j \leq m$). (or, equivalently, that $P$ occurs beginning at position $s+1$ in $T$).**

- **Valid shift $s$: if $P$ occurs with shift $s$ in $T$; otherwise, $s$ is an invalid shift.**

- **Finding all valid shifts with which a given $P$ occurs in a given $T$.**

- $\Sigma^*$ : **the set of all finite-length strings formed using characters from the alphabet $\Sigma$. Example:**

    $\Sigma = \{a, b, c\}$

    $\Sigma^* = \{\varepsilon, a, b, c, ab, bc, ac, abc, acb, aabbc, \ldots\ldots\}$

- $\varepsilon$ : **The zero-length empty string, also belongs to $\Sigma^*$.**

- $|x|$ : **The length of $x$.**

- **The concatenation of two strings $x$ and $y$, denoted $xy$, has length $|x| + |y|$ and consists of the characters from $x$ followed by the characters from $y$.**

# Notation and terminology

- $\omega \sqsubseteq x$ : string $\omega$ is a **prefix** of $x$, if $x = \omega y$ for some $y \in \Sigma^*$.



prefix             suffix

- $\omega \sqsupseteq x$ : $\omega$ is a **suffix** of $x$, if $x = y\omega$ for some $y \in \Sigma^*$.

  - If $\omega \sqsubseteq x$ or $\omega \sqsupseteq x$, then $|\omega| \leq |x|$.

  - The empty string $\varepsilon$ is both a suffix and a prefix of every string.

  - For example, we have **ab** $\sqsubseteq$ **abcca** and **cca** $\sqsupseteq$ **abcca**.

  - For any strings $x$ and $y$ and any character $a$, we have $x \sqsupseteq y$ if and only if $xa \sqsupseteq ya$.

  - $\sqsubseteq$ and $\sqsupseteq$ are transitive relations.

**Lemma 32.1: (Overlapping-suffix lemma)**

**Suppose that $x$, $y$, and $z$ are strings such that $x \sqsupset z$ and $y \sqsupset z$. If $|x| \leq |y|$, then $x \sqsupset y$. If $|x| \geq |y|$, then $y \sqsupset x$. If $|x| = |y|$, then $x = y$.**

*Proof* **See Fig for a graphical proof.**



(a)          (b)          (c)

- **For brevity of notation, we denote the $k$-character prefix $P[1 .. k]$ of the pattern $P[1 .. m]$ by $P_k$**

  - **Thus, $P_0 = \varepsilon$ and $P_m = P = P[1 .. m]$**

- **Similarly, we denote the $k$-character prefix of the text $T$ as $T_k$**

- **string-matching problem:**
  **finding all shifts $s$ in the range $0 \leq s \leq n\text{-}m$ such that $P \sqsupset T_{s+m}$**

$T_{s+m}$



$s = 4$

$P$

字符串匹配过程：从头到尾依序扫描文本 $T$，扫描到的字符串都是 $T$ 的前缀 $T_x$，若有 $P$ 匹配，则此时 $P$ 为 $T_x$ 的一个后缀。也就是，求文本 $T$ 的前缀的 $P$ 后缀。

- **Primitive operation: comparing characters**

**The naive algorithm finds all valid shifts using a loop that checks the condition $P[1 .. m] = T[s+1 .. s+m]$ for each of the $n-m+1$ possible values of $s$.**

NAIVE-STRING-MATCHER($T$, $P$)
1 $n \leftarrow$ length[$T$]
2 $m \leftarrow$ length[$P$]
3 for $s \leftarrow 0$ to $n-m$
4     if $P[1 .. m] = T[s+1 .. s+m]$
5         print "Pattern occurs with shift" $s$



- **The procedure can be interpreted graphically as sliding a "template" containing the pattern over the text.**

- **Line 3 considers each possible shift explicitly.**

- **The test on line 4 determines whether the current shift is valid or not; this test involves an implicit loop** （第4行包括一个隐式的循环）.

# 32.1 The naive string-matching algorithm

NAIVE-STRING-MATCHER(*T*, *P*)
1 *n* ← length[*T*]
2 *m* ← length[*P*]
3 for *s* ← 0 to *n-m*
4     if *P*[1 .. *m*] = *T*[*s*+1 .. *s*+*m*]
5         print "Pattern occurs with shift" *s*



(a)                  (b)

(c)                  (d)

// 返回首次匹配位置，用库函数实现
strstr(T, P);

思考题：
- 所有匹配（按伪代码规则）都找出来，怎么实现？
- 怎么实现strrstr？（最后一次匹配的位置）

```
// 返回首次匹配位置，自定义函数实现
char * _strstr(const char *T, const char *P)
{
    if(T == NULL)
        return NULL;
    int n = strlen(T), m = strlen(P), s, i;
    for(s=0; s<=n-m; s++)
    {
        for(i=0; i<m; i++)
            if(P[i] != T[s+i]) break;
        if(i == m) return T+s;
    }
    return NULL;
}
```

(a)  (b)  (c)  (d)

**Running time ?**

NAIVE-STRING-MATCHER(*T*, *P*)
1 $n \leftarrow$ length[*T*]
2 $m \leftarrow$ length[*P*]
3 for $s \leftarrow 0$ to *n-m*
4     if $P[1 .. m] = T[s+1 .. s+m]$
5         print "Pattern occurs with shift" *s*

# Exercise  32.1-2, 32.1-4

$T$

| a | b | c | a | d | d | b | c |
|---|---|---|---|---|---|---|---|

| b | c | a | d |
|---|---|---|---|

$P$

R-K algorithm performs well in practice and that also generalizes to other algorithms for related problems, such as two-dimensional pattern matching.
用了Hash和简单数论的方法。对 $T_s$ 计算 $p$ 时，还有更有效的方法（后一个 $m$ 个 $T_s$ 的 $p$ 值跟前面计算的结果有关系，可充分利用前面的计算信息，加快计算速度）。
计算 $p(P)$，可以用Horner's rule.

**coding rule: {a, b, c, d} ⇨ {0, 1, 2, 3},**

**then $p(P) = 1*4^3 + 2*4^2 + 0*4^1 + 3*4^0 = 99$**

$( (12345)_{10} = 1*10^4 + 2*10^3 + 3*10^2 + 4*10^1 + 5*10^0 )$

**if $p(P) == p(T_{s+m})$ ( $T_{s+m} = T[s+1 .. s+m]$, s = 0, 1, ..), then match.**

**or, if ( $p(P) \bmod q$ ) == ( $p(T_s) \bmod q$ ), check if $P == T_s$**

**preprocessing time: $\Theta(m)$**
**worst-case running time: $\Theta((n-m+1)m)$**

扩展阅读chapter31
Number-Theoretic Algorithms

$T$

| a | b | c | a | d | d | b | c |
|---|---|---|---|---|---|---|---|

| b | c | a | d |
|---|---|---|---|

$P$

**coding rule: {a, b, c, d} ⇨ {0, 1, 2, 3},**

**then $p(P) = 1*4^3 + 2*4^2 + 0*4^1 + 3*4^0 = 99$**

$( (12345)_{10} = 1*10^4 + 2*10^3 + 3*10^2 + 4*10^1 + 5*10^0 )$

## Efficient randomized pattern-matching algorithms

RM Karp, MO Rabin - IBM journal of research and development, 1987 - ieeexplore.ieee.org
We present randomized algorithms to solve the following string-matching problem and some
of its generalizations: Given a string X of length n (the pattern) and a string Y (the text), find
the first occurrence of X as a consecutive block within Y. The algorithms represent strings of …
☆ 保存　𝄞 引用　被引用次数：1801　相关文章　所有 9 个版本　≫

# Naive vs RK vs FA vs KMP

右移+1（依序扫描），再一一匹配判断

| b | a | c | a | c | a | b | a | c | a |
|---|---|---|---|---|---|---|---|---|---|

| a | c | a | c |
|---|---|---|---|

与 $P$ "无关"
（无需预处理）

**Naive**

$T$

| b | a | c | a | c | a | b | a | c | a |
|---|---|---|---|---|---|---|---|---|---|

| a | c | a | c |
|---|---|---|---|

$P$

**RK**

| a | c | a | c |
|---|---|---|---|

$P$

$p(P)$

**KMP**

**FA**

| b | a | c | a | c | a | b | a | c | a |
|---|---|---|---|---|---|---|---|---|---|

| a | c | a | c |
|---|---|---|---|

根据 $\delta(4, a) = 3$，即状态4时输入a得到状态3
（3个匹配），接着求 $\delta(3, b)$?

$\delta(P: Q, \Sigma)$，输入字符不匹配时，快速移动 $P$

| b | a | c | a | c | a | b | a | c | a |
|---|---|---|---|---|---|---|---|---|---|

| a | c | a | c |
|---|---|---|---|

根据 $\pi[4] = 2$，此例中是右移到 $P2$ 匹配
处，判断 $P[3]$ 是否与 $T$ 的下一个匹配?

$\pi(P)$，输入字符不匹配时，快速移动 $P$

# Naive vs RK vs FA vs KMP

| 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 | $\Theta(m)$ | $\Theta(n)$ |

| | 关键 | 特征 |
|---|---|---|
| FA | 求 $\delta(P: Q, \Sigma)$ | 输入字符$T[i]$不匹配时，快速移动 $P$ ，每个$T[i]$匹配一次 |
| KMP | 求 $\pi(P)$ | 输入字符$T[i]$不匹配时，快速移动 $P$ ，每个$T[i]$匹配一次 |

# 32.3 String matching with finite automata（有限自动机）

- **Many string-matching algorithms build a finite automaton (Machine) that scans the text $T$ for all occurrences of the pattern $P$.**

- **These string-matching automata are very efficient:**
  - **they examine each text character *exactly once*;**
  - **taking constant time per text character.**

- **These string-matching automata (Machine) are very efficient:**
    - **they examine each text character *exactly once* ;**
    - **taking constant time per text character.**
- **The matching time is** $\Theta(n)$**.**
    - **The preprocessing time (to build the automaton by pattern) can be large if** $\Sigma$ **is large.** （对西文文本来说，小写26，大写26，数字10，共62，再加上各种标点符号、或特殊符号、或希腊字母等，$\Sigma$约百余个字符，不算大；若中文，$\Sigma$可以很大）
    - **Section 32.4 describes a clever way around this problem.**

text $T$ | a | b | c | a | b | a | a | b | c | a | b | a | c |

pattern $P$ —— $s = 3$

$M$

```
input
text T  →  for a pattern P     →  output       →  EOT     ── yes ──→  stop
           build M                 matching          (EOF)
                                   case              ── no ──┐
```

# Finite automata (Machine)

- **A *finite automaton M* is a 5-tuple** $M = (Q, q_0, A, \Sigma, \delta)$ **, where**

  - $Q$ is a finite set of **states**,

  - $q_0 \in Q$ is the **start state**,

  - $A \subseteq Q$ is a distinguished set of **accepting states**,

  - $\Sigma$ is a finite **input alphabet**,

  - $\delta$ is a function from $Q \times \Sigma$ into $Q$, called the **transition function** of $M$.

- The $M$ begins in state $q_0$ and reads the characters of its input string one at a time. If the $M$ is in state $q$ and reads input $a$, it moves ("makes a transition") from state $q$ to $\delta(q, a)$. Whenever its current state $q$ is a member of $A$, the $M$ is said to have **accepted** the string read so far. An input that is not accepted is said to be **rejected**.

# Finite automata: an example

- Figure 32.6: A simple two-state finite automaton with state set $Q = \{0, 1\}$, start state $q_0 = 0$, and input alphabet $\Sigma = \{a, b\}$.

| state | input a | b |
|-------|---------|---|
| 0 | 1 | 0 |
| 1 | 0 | 0 |

(a)

(b)

(a) A tabular representation of the transition function $\delta$.

(b) An equivalent state-transition diagram.

a b a a a
<0, 1, 0, 1, 0, 1>

- **State 1 is the only accepting state** (shown blackened). Directed edges represent transitions. For example, the edge from state 1 to 0 labeled b indicates $\delta(1, b) = 0$. This automaton accepts those strings that end in an odd number of a's. For example, the sequence of states this automaton enters for input abaaa (including the start state) is <0, 1, 0, 1, 0, 1>, so it accepts this input. For input abbaa, the sequence of states is <0, 1, 0, 0, 1, 0>, so it rejects this input.

# Finite automata: *final-state function*

- A finite automaton $M$ induces a function $\Phi$, called the ***final-state function***, from $\Sigma^*$ to $Q$ such that $\Phi(\omega)$ is the state $M$ ends up in after scanning the string $\omega$. 终态函数: $M$ 读入字符串 $\omega$ 后得到状态 $\Phi(\omega)$

- Thus, $M$ accepts a string $\omega$ if and only if $\Phi(\omega) \in A$.

- The function $\Phi$ is defined by the recursive relation

$$\Phi(\varepsilon) = q_0 ,$$

$$\Phi(\omega a) = \delta(\Phi(\omega), a) \text{ for } w \in \Sigma^*, a \in \Sigma .$$

字符串匹配过程: $M$ 从头到尾依序扫描文本 $T$, 扫描到的字符串都是 $T$ 的前缀 $T_x$, 若有 $P$ 匹配, 则此时 $P$ 为 $T_x$ 的一个后缀。也就是, 求文本 $T$ 的前缀的 $P$ 后缀。

$\varepsilon$     $\omega$     $a$

$\longrightarrow M \longrightarrow q_0 \longrightarrow \cdots \longrightarrow M \longrightarrow q \longrightarrow M \longrightarrow q'$

$q_0 = \Phi(\varepsilon)$     $q = \Phi(\omega)$     $q' = \Phi(\omega a) = \delta(\Phi(\omega), a)$
$= \delta(q, a)$

# String-matching automata



- There is a string-matching automaton(**Machine**) for every pattern *P*

- This automaton must be constructed from the pattern in a preprocessing step before it can be used to search the text string.

- Figure illustrates this construction for the pattern *P* = ababaca.



(a)

$$\omega \qquad\qquad \omega$$

aba         aba<span style="color:red">a</span>

| | |           |

*P* : ababaca      ababaca

$\delta(2, a) = 3$      $\delta(3, a) = 1$

$q' = \Phi(\omega a) = \delta(\Phi(\omega), a) = \delta(q, a)$

含义：读入$\omega$时状态（与 *P* 中字符匹配个数）为$q$，再读入a时的状态是什么？

|  | input | | | *P* |
|---|---|---|---|---|
| state | 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)

$$M = (Q, q_0, A, \Sigma, \delta)$$

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

(c)

应用：读入$T_k$，看$T_k$的后缀跟*P*的前缀$P_p$的匹配情况，匹配数为$m$，则字符串匹配。

# String-matching automata: an example

(a) A state-transition diagram for the string-matching automaton that accepts all strings ending in the string **ababaca**. State 0 is the start state, and state 7 (shown blackened) is the only accepting state. A directed edge from state $i$ to state $j$ labeled $\alpha$ represents $\delta(i, \alpha) = j$. The right-going edges forming the "spine" of the automaton, shown heavy, correspond to successful matches between pattern and input characters.

The left-going edges correspond to failing matches.

$M = (Q, q_0, A, \Sigma, \delta)$

Some edges corresponding to failing matches are not shown; by convention, if a state $i$ has no outgoing edge labeled $\alpha$ for some $\alpha \in \Sigma$, then $\delta(i, \alpha) = 0$.



$M$ of $P = $ ababaca

(a diagram)

(a)

Figure 32.7

$M = (Q, q_0, A, \Sigma, \delta)$

(b) The corresponding transition function $\delta$ (a table), and the pattern string $P = $ ababaca. The entries corresponding to successful matches between pattern and input characters are shown shaded.

(c) The operation of the automaton on the text $T = $ abababacaba.

自动机 $M$ 处理 $T_i$ 后，其状态为 $\Phi(T_i)$

One occurrence of the pattern is found, ending in position 9.

$M$ of $P = $ ababaca
(a table)

| state | input a | b | c | P |
|---|---|---|---|---|
| 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 | b | a | c | a | b | a |
| state $\phi(T_i)$ | 0 | 1 | 2 | 3 | 4 | 5 | 4 | 5 | 6 | 7 | 2 | 3 |

(c)

## Kore: how to build $\delta(q, a)$ ?

The operation of the automaton on the text $T =$ ababababacaba.

自动机处理$T_i$后，其状态为$\Phi(T_i)$

$M = (Q, q_0, A, \Sigma, \delta)$



(a)

$\delta(0, a) = 1, \delta(3, b) = 4, \delta(4, c) = 0, \delta(5, b) = 4$ ?

$T$ | a b a b a b a c a b a | ...

$P$ | a b a b a c a



| | input | | | |
|---|---|---|---|---|
| state | a | b | c | P |
| 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 | b | a | c | a | b | a |
| state $\phi(T_i)$ | 0 | 1 | 2 | 3 | 4 | 5 | 4 | 5 | 6 | 7 | 2 | 3 |

(c)

The operation of the automaton on the text $T = $ abababacaba.

$M = (Q, q_0, A, \Sigma, \delta)$

自动机处理$T_i$后，其状态为$\Phi(T_i)$

$\delta(0, a) = 1$ ?



(a)

|       | input |   |   | P |
|-------|-------|---|---|---|
| state | 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)

$T$ a b a b a b a c a b a   ...

$P$ a b a b a c a

从空字符开始，从文本串 $T$ 的第一字符依序扫描，输入 a…时（…表示还有很多字符），$P$ 的第1个跟其匹配，即 $\delta(0, a) = 1$；【＊从T一个一个地扫描，跟KMP有相似性】

...

$T$ a ...

$P$ a b a b a c a

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

(c)

# String-matching automata: an example

The operation of the automaton on the text $T =$ ababababacaba.

$M = (Q, q_0, A, \Sigma, \delta)$

自动机处理$T_i$后，其状态为$\Phi(T_i)$

$\delta(0, \text{a}) = 1$, $\delta(3, \text{b}) = 4$ ?

(a)

$T$: a b a b a b a c a b a ...

$P$: a b a b a c a

| | input | | | |
|---|---|---|---|---|
| state | a | b | c | P |
| 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)

a b a ...

$T$: a b a b ...

$P$: a b a b

状态3时（有3个匹配），
输入b，即文本串 $T$ 为
abab…时（…表示还有很多
字符），$P$ 的前4个跟其匹
配，即 $\delta(3, \text{b}) = 4$

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

(c)

The operation of the automaton on the text $T =$ **ababababacaba**.

$M = (Q, q_0, A, \Sigma, \delta)$

自动机处理$T_i$后，其状态为$\Phi(T_i)$



(a)

$\delta(0,\text{a}) = 1$, $\delta(3,\text{b}) = 4$, **$\delta(4,\text{c}) = 0$** ?

$T$ | a b a b a b a c a b a | ... |

$P$ | a b a b a c a |

**a b a b ...**

$T$ **a b a b c ...**

$P$ **a b a b a**

**a b a b**

**a b a**

**......**

Naive: 状态4时（有4个匹配），输入c，即文本串 $T$ 为ababc…时，$P$ 的前5个跟其不匹配，即$\delta(4, c) \mathrel{!=} 5$；
把 $P$ 按字符右移1位（寻找新的可能匹配），$P$ 的前4个跟其不匹配，即$\delta(4, c) \mathrel{!=} 4$；
把 $P$ 按字符右移2位，$P$ 的前3个跟其不匹配，即$\delta(4, c) \mathrel{!=} 3$；以此类推。

|  | input | | | $P$ |
|---|---|---|---|---|
| state | 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)

The operation of the automaton on the text $T =$ **abababacaba**.

$M = (Q, q_0, A, \Sigma, \delta)$

自动机处理$T_i$后，其状态为$\Phi(T_i)$



(a)

$\delta(0, a) = 1, \delta(3, b) = 4, \delta(4, c) = 0, \boldsymbol{\delta(5, b) = 4}$ ?

| $T$ | a b a b a b a c a b a | ... |
|-----|-----------------------|-----|

| $P$ | a b a b a c a |
|-----|---------------|

每次把 $P$ 右移1位后，都从 $P$ 的第一个字符开始匹配，跟naive方法一样？

**肯定不用这样做！**

**快速右移到 $\delta(q, x)$ 处！**

**求 $\delta(q, x)$ 是关键！**

$T$   a b a b a **...**
    a b a b a b **...**
$P$   a b a b a c
     a b a b a
      a b a b

|       | input |   |   |     |
|-------|-------|---|---|-----|
| state | a     | b | c | $P$ |
| 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 | b | a | c | a | b | a |
| state $\phi(T_i)$ | 0 | 1 | 2 | 3 | 4 | 5 | 4 | 5 | 6 | **7** | 2 | 3 |

(c)

- Suffix function $\sigma$ corresponding to $P$ :

  A mapping from $\Sigma^*$ to $\{0, 1, \ldots, m\}$ such that $\sigma(x)$ is the length of the longest prefix of $P$ that is a suffix of $x$ ($x$ 的后缀且是$P$的最长前缀的长度) ：

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

- The suffix function $\sigma$ is well defined since the empty string $P_0 = \varepsilon$ is a suffix of every string. As examples,

  ◆ for the pattern $P = \text{ab}$, we have $\sigma(\varepsilon) = 0$, $\sigma(\text{ccaca}) = 1$, and $\sigma(\text{ccab}) = 2$.

- For a pattern $P$ of length $m$, we have $\sigma(x) = m$ if and only if $P \sqsupset x$.

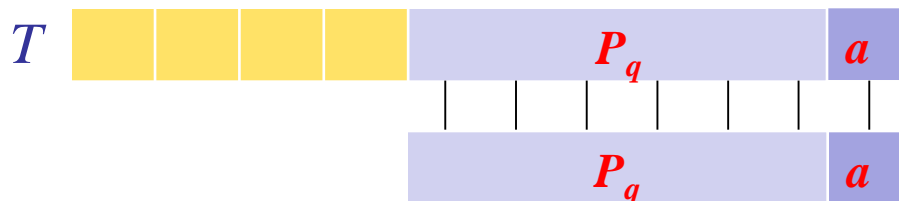- From the definition of the suffix function, if $x \sqsupset y$, then $\sigma(x) \le \sigma(y)$.

$M = (Q, q_0, A, \Sigma, \delta)$

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



We define the ***string-matching automaton*** that corresponds to a given pattern $P[1 .. m]$ as follows.

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

$$\delta(q, a) = \sigma(P_q a) \qquad\qquad (32.3)$$

- where, the state set $Q$ is $\{0, 1, \ldots, m\}$, the start state $q_0$ is state 0, and state $m$ is the only accepting state $A$.

$\delta(q, a) = \sigma(P_q a)$ 的定义合理，后面将证明，$\delta(q, a) = \sigma(P_q a) = \sigma(T_i a)$，即，扫描 $T_i$ 后，匹配为 $P_q$，接着读入 $a$，对 $T_i a$ 的匹配与对 $P_q a$ 的匹配是一样的(**Lemma 32.1**)。$P_q a$ 的长度比 $T_i a$ 短，处理起来（求 $\delta$）就简单得多。
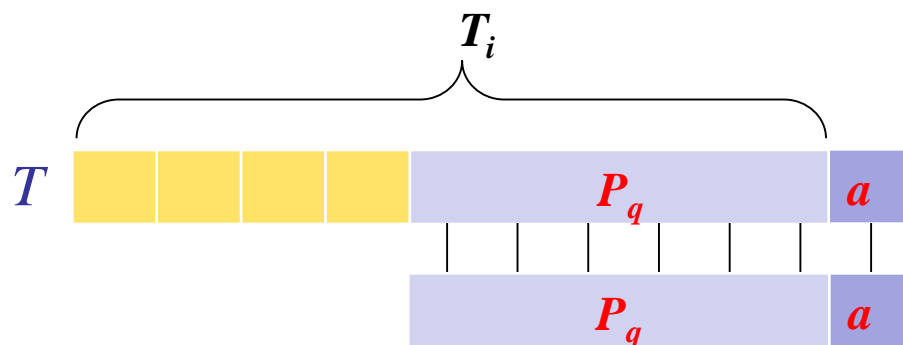
$T_i$

$M = (Q, q_0, A, \Sigma, \delta)$

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

We define the *machine* **M** :

$Q = \{0, 1, \ldots, m\}; q_0 = 0; A = \{m\};$

$$\delta(q, a) = \sigma(P_q a) \qquad \textbf{32.3)}$$



Intuitively, the machine **M** maintains an invariant:

$$\Phi(T_i) = \sigma(T_i) , \qquad ( \text{ where, } \Phi(T_i) = q = \sigma(T_i) ). \qquad \textbf{(32.4)}$$

自动机 $M$ 扫描字符串 $T$ 的过程中，扫描到前缀子串 $T_i$ 时状态为 $q$（为 $T_i$ 的后缀函数 $\sigma(T_i)$），接着扫描下一个字符 $T[i+1]$（记为 $a$），状态转移到 $\delta(q, a) = \sigma(P_q a)$，这就是扫描到前缀子串 $T_{i+1}$ 时状态（为 $T_{i+1}$ 的后缀函数 $\sigma(T_{i+1})$）

$$\Phi(T_{i+1}) = \Phi(T_i a) = \delta(\Phi(T_i), a) = \delta(q, a) = \sigma(P_q a) \overset{?}{=} \sigma(T_i a) = \sigma(T_{i+1}) \qquad \textbf{(32.4)*}$$

[ (32.3) maintains the invariant, or, it is rationale for defining (32.3). ]

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

We define the *machine* **M** :

$$Q = \{0, 1, \dots, m\};\ q_0 = 0;\ A = \{m\};\ \Sigma;$$

$$\delta(q, a) = \sigma(P_q a) \qquad \textbf{(32.3)}$$

$T_i$

$T$

$P_q$   $a$

$P_q$   $a$

- $\Phi(T_{i+1}) = \Phi(T_i a) = \delta(\Phi(T_i), a) = \delta(q, a) = \sigma(P_q a) \overset{?}{=} \sigma(T_i a) = \sigma(T_{i+1})$    **(32.4)**
  [ (32.3) maintains the invariant, or, it is rationale for defining (32.3). ]

- *Lemma 32.3:* $\sigma(T_i a) = \sigma(P_q a)$                        **(32.A)**

  **this lemma means definition (32.3) maintains the desired invariant (32.4).**

- **Compute: With (32.A), to compute $\sigma(T_i a)$, we can compute $\sigma(P_q a)$.**
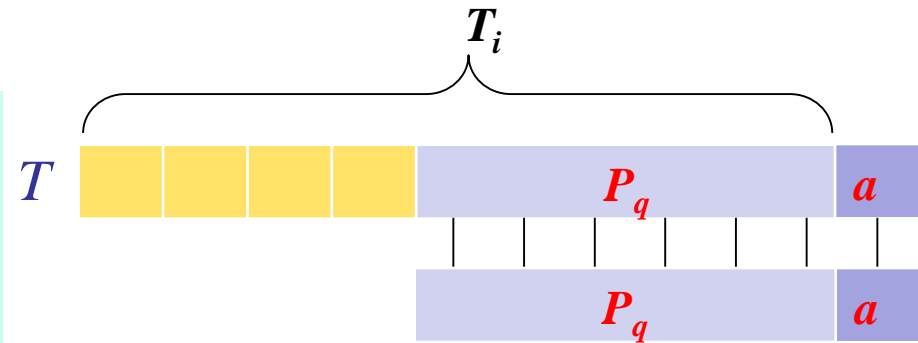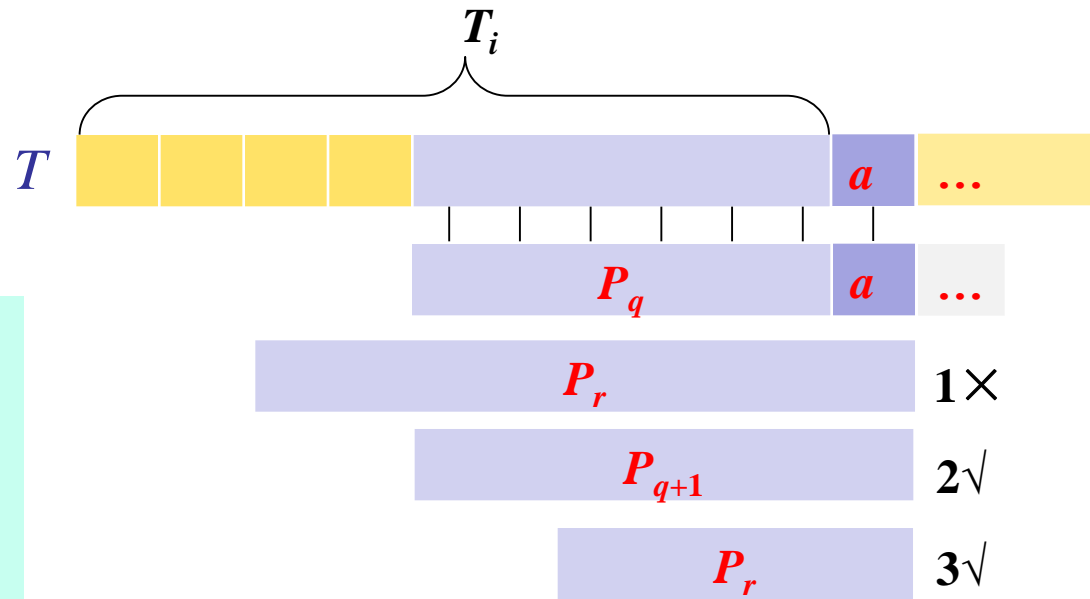
# String-matching automata

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

We define the *machine* **M** :

$Q = \{0, 1, \ldots, m\};$

$q_0 = 0; A = \{m\}; \Sigma;$

$\delta(q, a) = \sigma(P_q a).$ **(32.3)**



$T_i$

$T$

$a$ ...

$P_q$ $a$ ...

$P_r$ $1\times$

$P_{q+1}$ $2\checkmark$

$P_r$ $3\checkmark$

- $\Phi(T_{i+1}) = \Phi(T_i a) = \delta(\Phi(T_i), a) = \delta(q, a) = \sigma(P_q a) = \sigma(T_i a) = \sigma(T_{i+1})$ **(32.4)**
  [ (32.3) maintains the invariant, or, it is rationale for defining (32.3). ]

- *Lemma 32.3*: **If $\sigma(T_i) = \sigma(P_q) = q$, then $\sigma(T_i a) = \sigma(P_q a)$ .** **(32.A)**
  *Proof*
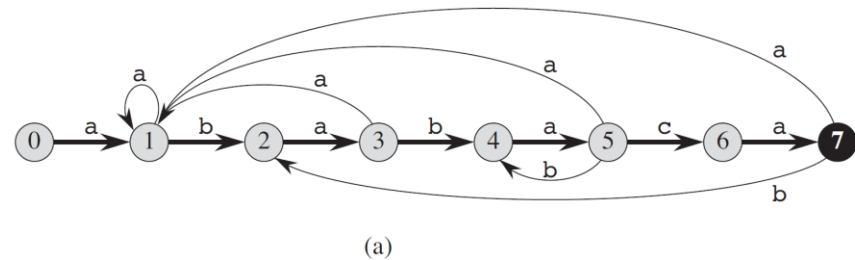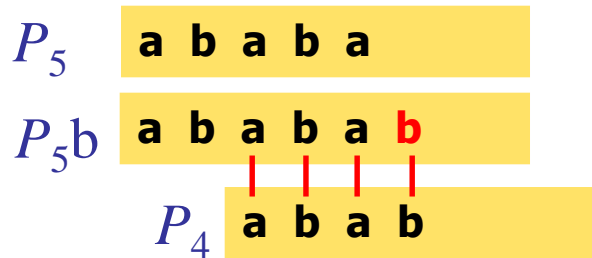  **Situation 1 is impossible （$\sigma(T_i a)>q+1$不可能）. if 1 满足，$\sigma(T_i) > q$ , 与假设矛盾。**
  **Apparently, if $P[q+1] = a$ , it is situation 2, $\sigma(T_i a) = \sigma(P_{q+1}) = q+1$; else , situation 3.**

- *Form 32.3* **and** *32.A* **shows the automaton is in state $\sigma(T_i)$ after scanning character $T[i]$. Since $\sigma(T_i) = m$ if and only if $P \sqsupset T_i$ , the machine is in the accepting state $m$ if and only if the pattern $P$ has just been scanned.**

# String-matching automata

For example, in the string-matching automaton of Figure 32.7, $\delta(5, \mathbf{b}) = 4$.

We make this transition because if the automaton reads a **b** in state $q = 5$, then $P_q\mathbf{b} = \mathbf{ababab}$, and then, $\delta(5, \mathbf{b}) = \sigma(\mathbf{ababab}) = 4$.

$P_5$   **a b a b a**

$P_5\mathbf{b}$   **a b a b a b**

$P_4$   **a b a b**

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

We define the *machine* **M** :

$Q = \{0, 1, \ldots, m\};$

$q_0 = 0; A = \{m\}; \Sigma;$

$\delta(q, a) = \sigma(P_q a).$     (32.3)



(a)

Figure 32.7

|  | input | | | |
|---|---|---|---|---|
| state | a | b | c | P |
| 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 | b | a | c | a | b | a |
| state $\phi(T_i)$ | 0 | 1 | 2 | 3 | 4 | 5 | 4 | 5 | 6 | **7** | 2 | 3 |

(c)

$\sigma(x) = \max \{k : P_k \sqsupset x\}.$
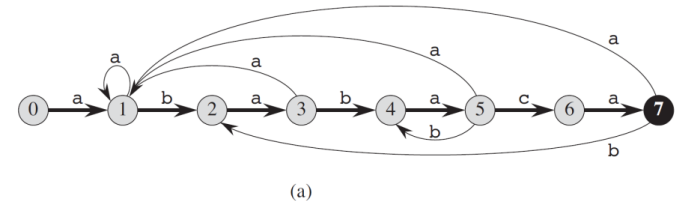
A string-matching automaton

$M = (Q, q_0, A, \Sigma, \delta)$ :

$\quad Q = \{0, 1, \ldots, m\}; \quad q_0 = 0; \quad A = \{m\};$

$\quad \delta(q, a) = \sigma(P_q a) = \sigma(T_{i-1} a) .$ $\qquad$ **(32.3)**



(a)

| state | input a | b | c | P |
|---|---|---|---|---|
| 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 | b | a | c | a | b | a |
| state $\phi(T_i)$ | 0 | 1 | 2 | 3 | 4 | 5 | 4 | 5 | 6 | **7** | 2 | 3 |

(c)

**FINITE-AUTOMATON-MATCHER(*T, δ, m*)**

1 $n \leftarrow length[T]$
2 $q \leftarrow 0$
3 **for** $i \leftarrow 1$ **to** $n$ $\quad$ // scan *T*
4 $\quad\quad a \leftarrow T[i]$
5 $\quad\quad q \leftarrow \delta(q, a)$
6 $\quad\quad$ **if** $q == m$
7 $\quad\quad\quad$ **print** "Pattern occurs with shift" $i - m$

**Running time ?**

◆ **The matching time is $\Theta(n)$.**

◆ **However, it does not include the preprocessing time required to compute the transition function $\delta$.**
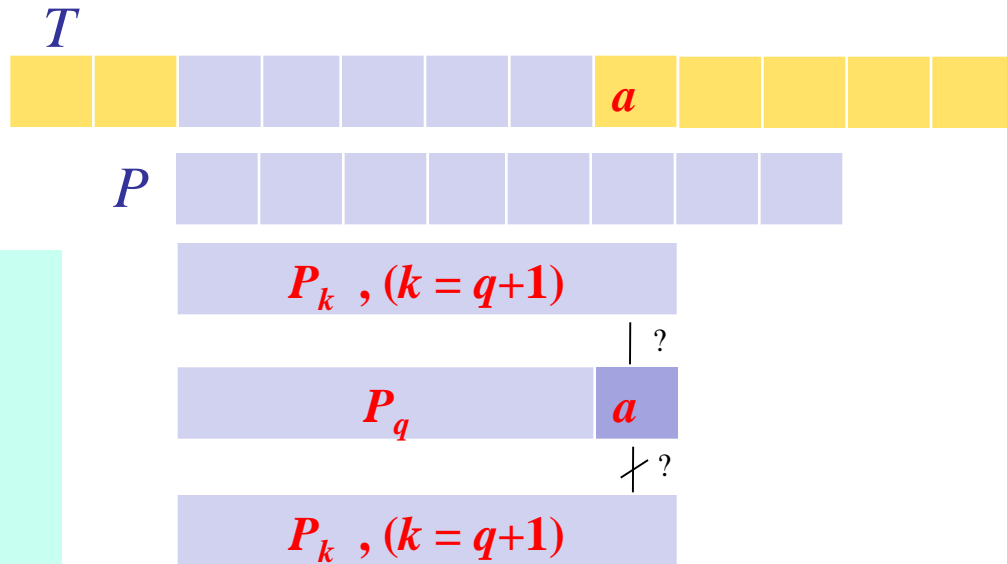
# Computing the transition function

$T$

$P$

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

A string-matching automaton
$M = (Q, q_0, A, \Sigma, \delta)$ :
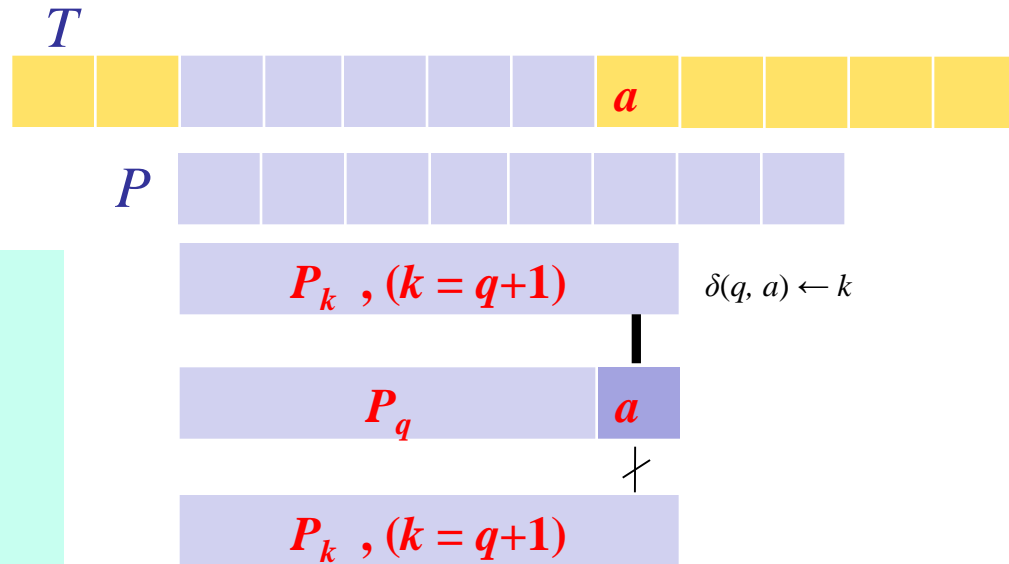  $Q = \{0, 1, \ldots, m\}$; $q_0 = 0$; $A = \{m\}$;
  $\delta(q, a) = \sigma(P_q a)$ .       **(32.3)**

$P_k$ , $(k = q+1)$

$|$ ?

$P_q$   $a$

$\not\vdash$ ?

$P_k$ , $(k = q+1)$

Computing $\delta$ from a given pattern $P[1 .. m]$ :

```
COMPUTE-TRANSITION-FUNCTION(P, Σ)
1  m ← length[P]
2  for q ← 0 to m
3      for each character a ∈ Σ
4          k ← min(m, q + 1)
5          while P_k !⊐ P_q a
6              k--
7          δ(q, a) ← k
8  return δ
```

$P_q$ 时（即$T$与$P$的前$q$个字符匹配时），输入第 $q+1$（即第$k$个）字符 $a$ 时：

1. $k = q+1$ （超过$m$时，取$m$，匹配数不大于$m$）

2. $P_k \sqsupset P_q a$ ?

# Computing the transition function



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

A string-matching automaton

$M = (Q, q_0, A, \Sigma, \delta)$ :

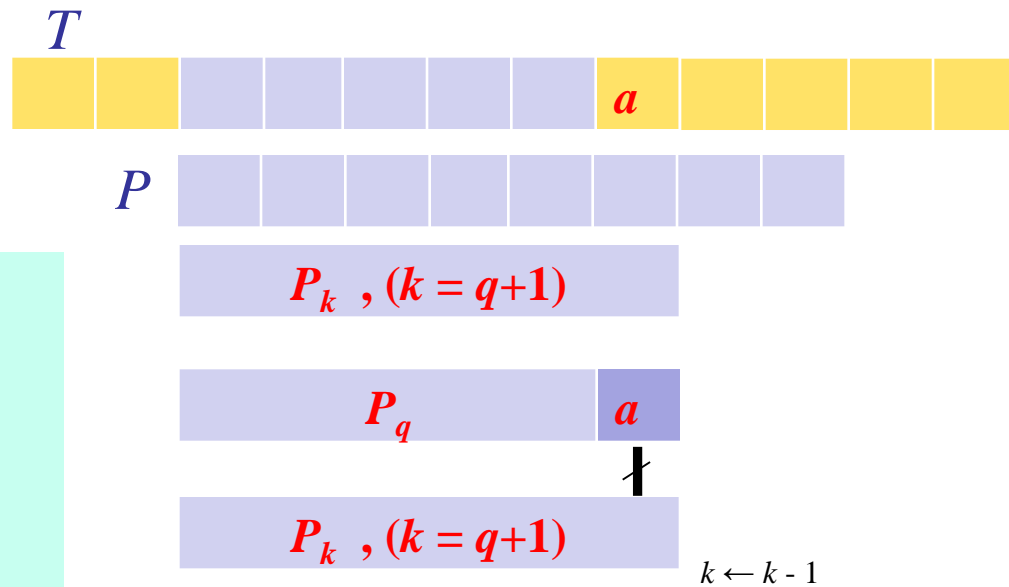  $Q = \{0, 1, \ldots, m\}$;  $q_0 = 0$;  $A = \{m\}$;

  $\delta(q, a) = \sigma(P_q a)$ .        **(32.3)**

Computing $\delta$ from a given pattern $P[1 .. m]$ :

```
COMPUTE-TRANSITION-FUNCTION(P, Σ)
1  m ← length[P]
2  for q ← 0 to m
3      for each character a ∈ Σ
4          k ← min(m, q + 1)
5          while P_k !⊐ P_q a
6              k--
7          δ(q, a) ← k
8  return δ
```

$P_q$ 时（即$T$与$P$的前$q$个字符匹配时），输入第 $q+1$（即第$k$个）字符 $a$ 时：

1. $k = q+1$（超过$m$时，取$m$，匹配数不大于$m$）

2. $P_k \sqsupset P_q a$ ?

3. 若2成立，则 $P_q a == P_k$，即，对在 $T$ 的继续扫描过程中，若扫描的下一个字符 $a==P[k]$，则匹配字符增加1

# Computing the transition function



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

A string-matching automaton

$M = (Q, q_0, A, \Sigma, \delta)$ :
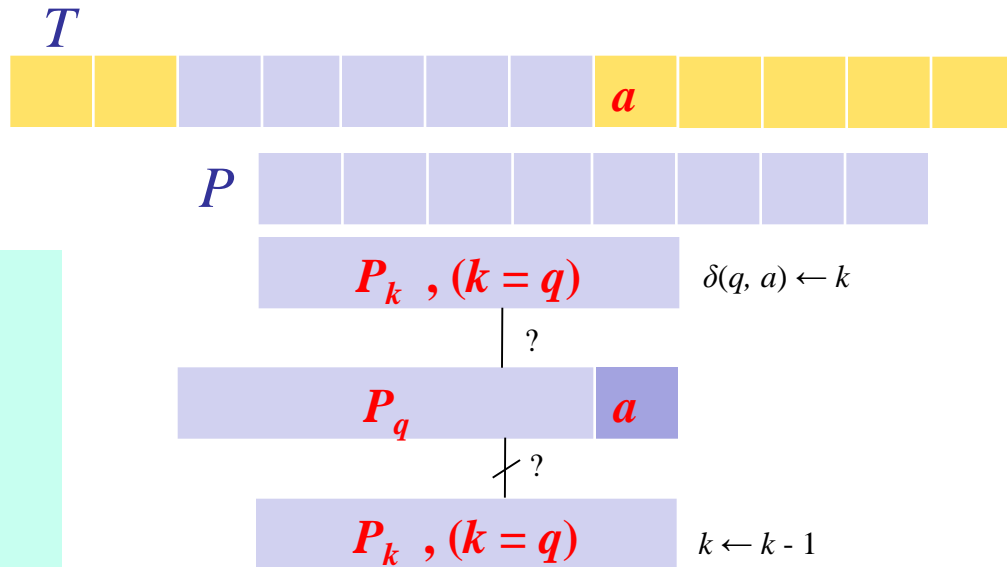
$Q = \{0, 1, \ldots, m\}$;  $q_0 = 0$;  $A = \{m\}$;

$\delta(q, a) = \sigma(P_q a)$ .               **(32.3)**

Computing $\delta$ from a given pattern $P[1 .. m]$ :

```
COMPUTE-TRANSITION-FUNCTION(P, Σ)
1  m ← length[P]
2  for q ← 0 to m
3      for each character a ∈ Σ
4          k ← min(m, q + 1)
5          while P_k !⊐ P_q a
6              k--
7          δ(q, a) ← k
8  return δ
```
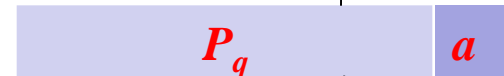
$P_q$ 时（即$T$与$P$的前$q$个字符匹配时），输入第 $q+1$（即第$k$个）字符 $a$ 时：
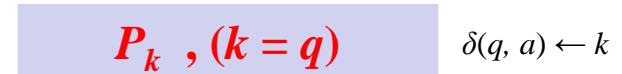
1. $k = q+1$（超过$m$时，取$m$，匹配数不大于$m$）

2. $P_k \sqsupset P_q a$ ?

3. 若2成立，则 $P_q a == P_k$，即，对在 $T$ 的继续扫描过程中，若扫描的下一个字符 $a == P[k]$，则匹配字符增加1

4. 若2不成立，即 $P_q a \mathrel{!=} P_k$，即，对在$T$的继续扫描过程中，若扫描的下一个字符$a \mathrel{!=} P[k]$，模版右移($k$--)，goto step 2

# Computing the transition function



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

A string-matching automaton

$M = (Q, q_0, A, \Sigma, \delta):$

$\quad Q = \{0, 1, \ldots, m\}; \ q_0 = 0; \ A = \{m\};$

$\quad \delta(q, a) = \sigma(P_q a).$　　　　(32.3)
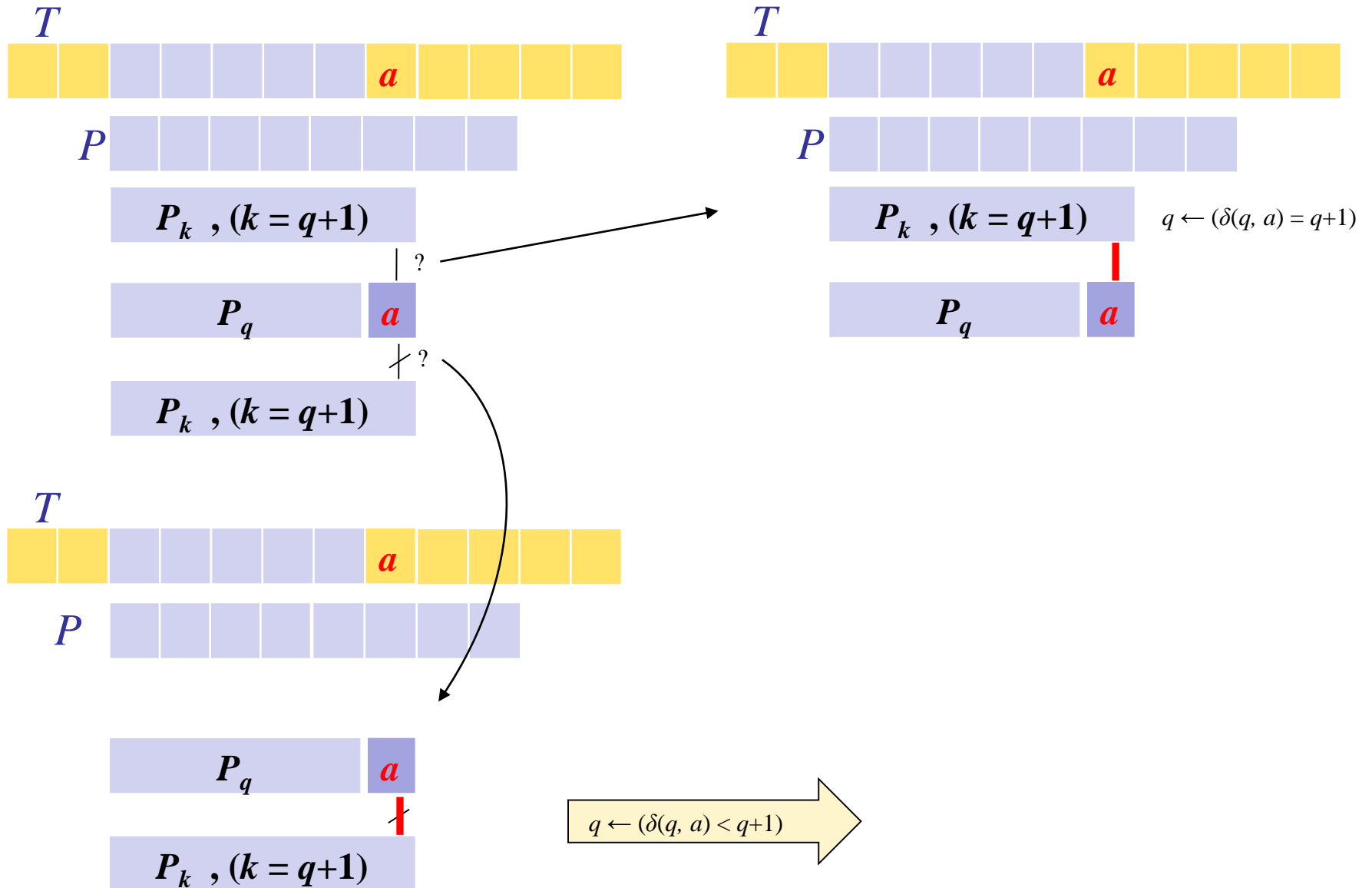
Computing $\delta$ from a given pattern $P[1 .. m]$ :
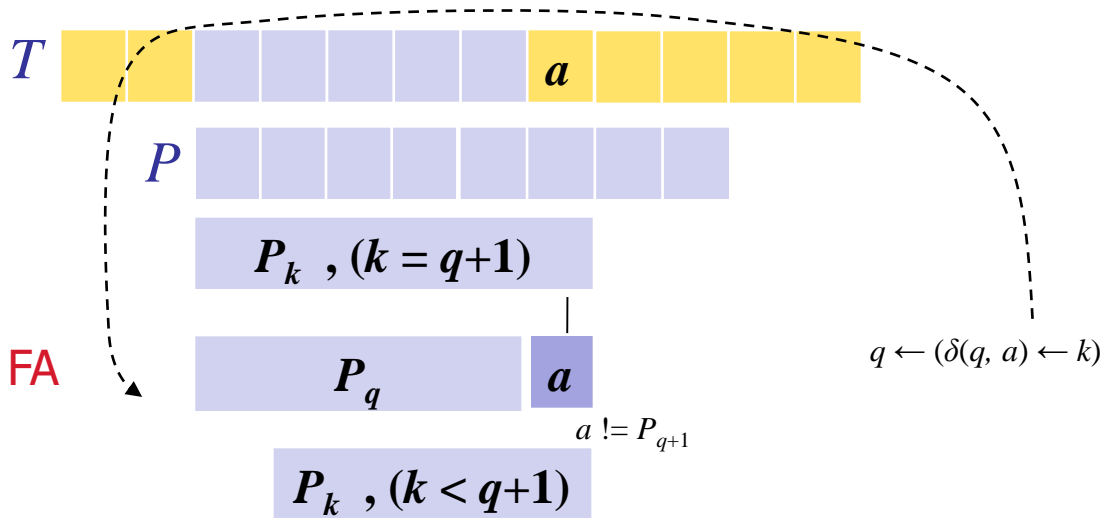
```
COMPUTE-TRANSITION-FUNCTION(P, Σ)
1  m ← length[P]
2  for q ← 0 to m
3      for each character a ∈ Σ
4          k ← min(m, q + 1)
5          while P_k !⊐ P_q a
6              k--
7          δ(q, a) ← k
8  return δ
```

$P_q$ 时（即$T$与$P$的前$q$个字符匹配时），输入第 $q+1$（即第$k$个）字符 $a$ 时：

1. $k = q+1$（超过$m$时，取$m$，匹配数不大于$m$）

2. $P_k \sqsupset P_q a$ ?

3. 若2成立，则 $P_q a == P_k$，即，对在 $T$ 的继续扫描过程中，若扫描的下一个字符 $a == P[k]$，则匹配字符增加1

4. 若2不成立，即 $P_q a \mathrel{!=} P_k$，即，对在$T$ 的继续扫描过程中，若扫描的下一个字符$a \mathrel{!=} P[k]$，模版右移($k$--)，goto step 2

# Computing the transition function

$T$

$P$

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

A string-matching automaton

$M = (Q, q_0, A, \Sigma, \delta)$ :

$\quad Q = \{0, 1, \ldots, m\}; \quad q_0 = 0; \quad A = \{m\};$

$\quad \delta(q, a) = \sigma(P_q a) .$        (32.3)

$P_k$ , $(k = q)$     $\delta(q, a) \leftarrow k$

?

$P_q$   $a$

?

$P_k$ , $(k = q)$     $k \leftarrow k - 1$

Computing $\delta$ from a given pattern $P[1 .. m]$ :

COMPUTE-TRANSITION-FUNCTION($P, \Sigma$)
1  $m \leftarrow length[P]$
2  **for** $q \leftarrow 0$ to $m$
3      **for** each character $a \in \Sigma$
4          $k \leftarrow \min(m, q + 1)$
5          **while** $P_k \; !\!\sqsupset P_q a$
6              $k$--
7          $\delta(q, a) \leftarrow k$
8  **return** $\delta$

- **Running time ?**

$T$

$P$

$P_k$ , $(k = q+1)$

$|$ ?

$P_q$   $a$

$\dashv$ ?

$P_k$ , $(k = q+1)$

$T$

$P$

$P_k$ , $(k = q+1)$   $q \leftarrow (\delta(q, a) = q+1)$

$P_q$   $a$

$T$

$P$

$P_q$   $a$

$q \leftarrow (\delta(q, a) < q+1)$

$P_k$ , $(k = q+1)$

# *32.4 The Knuth-Morris-Pratt algorithm

$T$

$P$

$P_k$ , $(k = q+1)$

FA

$P_q$    $a$

$a \,!\!= P_{q+1}$

$P_k$ , $(k < q+1)$

$q \leftarrow (\delta(q, a) \leftarrow k)$

根据输入的 $a$ ，然后查表 $\delta$ 知道需要转移的位置。

自动机构造完后，从 $\delta$ 已经知道输入 $a$ 后 $P$ 应快速右移多少（这种思想跟 KMP 相似，但 FA 的核心在于求 $\delta$ 有额外计算开销）。

$T$ 中一个字符扫描一次。

$T$

$P$

$P_k$ , $(k = q+1)$

$q$++ // 匹配数增加1

$P_q$    $a$

$\nmid$ ?

$P_k$ , $(k = q+1)$

?

$T$

$a$

$P$

**Naive shifting.    No!**

# *32.4 The Knuth-Morris-Pratt algorithm

$T$

$P$

$P_k$ , $(k = q+1)$

FA

$P_q$  $a$

$a \ne P_{q+1}$

$P_k$ , $(k < q+1)$

$q \leftarrow (\delta(q, a) \leftarrow k)$

根据输入的 $a$，然后查表 $\delta$ 知道需要转移的位置。

自动机构造完后，从 $\delta$ 已经知道输入 $a$ 后 $P$ 应快速右移多少（这种思想跟 KMP 相似，但 FA 的核心在于求 $\delta$ 有额外计算开销）。

$T$ 中一个字符扫描一次。

$T$

$P$

$P_k$ , $(k = q+1)$

$q{+}{+}$ // 匹配数增加1

KMP

$P_q$  $a$

$P_k$ , $(k = q+1)$

$k \leftarrow \pi[q]$

$T$

$P$

**Naive shifting.    No!**

$P_q$

预处理时已知

$P_k$ , $(k < q)$

$q \leftarrow k$ //$\pi[q]$

KMP本质：$P$的前缀$P_q$与$T$的匹配；$P_k$ $(k<q)$是 $P_q$ 的后缀（也是 $P_q$ 的前缀），最大的$k$是多少？$q \leftarrow \pi(q)$ 后，重新看是否 $a == P[q+1]$（重复执行此过程）

KMP 的前缀函数构造完后，不需要输入下一个 $a$ 就知道移动多少位置（从 $P_q$ 移动到 $P_k$）。$T$ 中一个字符可能扫描多次。

# *32.4 The Knuth-Morris-Pratt algorithm

KMP is a linear-time string-matching algorithm due to Knuth, Morris, and Pratt.

# *32.4 The Knuth-Morris-Pratt algorithm



(a)

$P_5 \sqsupset T_{s+5}$ , but,

$T[s+5+1] \mathrel{!=} P[5+1]$  ( $q = 5$ )

Naive shifting, $P_4 \sqsupset T_{s+1+4}$ ?      No!

# *32.4 The Knuth-Morris-Pratt algorithm



(a)

$P_5 \sqsupset T_{s+5}$ , but,

$T[s+5+1] != P[5+1]$   ( $q = 5$ )



?

**Naive shifting, $P_4 \sqsupset T_{s+1+4}$ ?      No!**
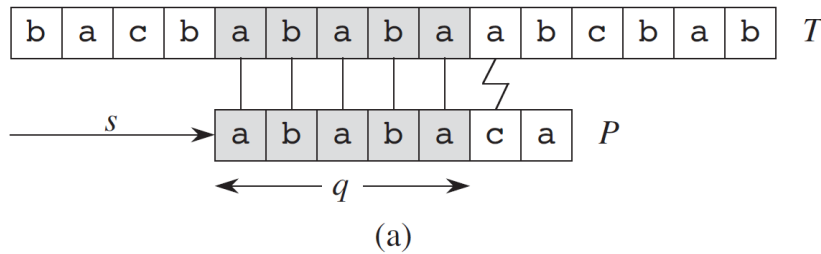


$s' = s + 2$

(b)

$P_q$

$P_k$

(c)

**We have already known that $P_k$ is the maximum suffix of $P_q$ , that is, $P_{k'} \sqsupset P_q$ ( $k' < q$, and $k = \max(k')$ ). For this example, $q$ is 5, $k$ is 3. So, we have $P_3 \sqsupset P_5 \sqsupset T_{s+5}$ , we just check whether ..**

$T[s+5+1] != P[3+1]$    ( $q \leftarrow k = 3$ )

**prefix function** for the pattern $P$ is the function..

$\pi : \{1, 2, \ldots, m\} \rightarrow \{0, 1, \ldots, m\text{-}1\}$ such that

$\pi[q] = \max\{k : k < q \text{ and } P_k \sqsupset P_q \}$

$P$ 的 前缀$P_q$ 的后缀是$P$ 的 前缀$P_k$ （最长的）

✔ $P$ 的 前缀$P_q$ 的 前缀$P_k$ 是 前缀$P_q$ 的后缀（最长的）

$P_q$ 是 $P$ 的前缀，

$P_k$ 是 $P_q$ 的前缀，且 $P_k$ 是 $P_q$ 的后缀，

最大的 $k$ 即为 $\pi[q]$

$P$ = ababaca
$P_5$ = ababa
$P_3$ = aba
$\Rightarrow$ $\pi[5] = 3$



(a)



(b)



(c)

$P_5 \sqsupset T_{s+5}$ , but,

$T[s+5+1] \neq P[5+1]$ （$q = 5$）

**We have already known that $P_k$ is the maximum suffix of $P_q$ , that is, $P_{k'} \sqsupset P_q$ ($k' < q$, and $k = \max(k')$ ). For this example, $q$ is 5, $k$ is 3. So, we have $P_3 \sqsupset P_5 \sqsupset T_{s+5}$ , we just check whether ..**

$T[s+5+1] \neq P[3+1]$ （$q \leftarrow k = 3$）

**prefix function** for the pattern $P$ is the function..

$\pi : \{1, 2, \ldots, m\} \rightarrow \{0, 1, \ldots, m\text{-}1\}$ such that

$\pi[q] = \max\{k : k < q \text{ and } P_k \sqsupset P_q \}$.

$P$

| a | b | a | b | a | c | a |
|---|---|---|---|---|---|---|

$P_5$    a b a b a c a

$P_3$    a b a b a c a     $\pi[5] = 3$

$P_1$    a b a b a c a     $\pi[3] = 1$
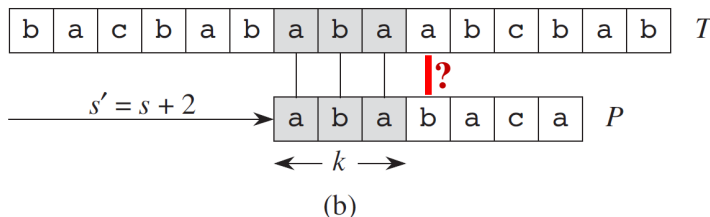
$P_0$    $\varepsilon$ a b a b a c a     $\pi[1] = 0$

| $i$ | 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)

$P$ 的 前缀$P_q$ 的 前缀$P_k$ 是 前缀$P_q$ 的后缀（最长的）

# *32.4 The Knuth-Morris-Pratt algorithm

KMP-MATCHER$(T, P)$

```
1   n = T.length
2   m = P.length
3   π = 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] ≠ T[i]
7           q = π[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 = π[q]                    // look for the next match
```
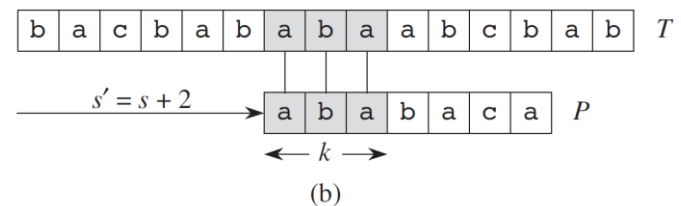
**we have  $P_q \sqsupset T[i\text{-}1]$,**

**check whether  $P[q+1] \ne T[i]$**

(a)

| $i$ | 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 |

(b)

# *32.4 The Knuth-Morris-Pratt algorithm

KMP-MATCHER($T$, $P$)

1. $n = T.length$
2. $m = P.length$
3. $\pi = $ COMPUTE-PREFIX-FUNCTION($P$)
4. $q = 0$
5. **for** $i = 1$ **to** $n$
6.     **while** $q > 0$ and $P[q + 1] \neq T[i]$
7.         $q = \pi[q]$
8.     **if** $P[q + 1] == T[i]$
9.         $q = q + 1$
10.     **if** $q == m$
11.         print "Pattern occurs with shift" $i - m$
12.         $q = \pi[q]$

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** $\pi$

**prefix function** for the pattern $P$ is the function..

$\pi : \{1, 2, \ldots, m\} \rightarrow \{0, 1, \ldots, m\text{-}1\}$ such that
$\pi[q] = \max\{k : k < q \text{ and } P_k \sqsupset P_q \}$.

we have $P_k \sqsupset P[q\text{-}1]$,
check if $P[k+1] \neq P[q]$

| $i$ | 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)

$P_5$   | a | b | a | b | a | c | a |

$P_3$   | a | b | a | b | a | c | a |   $\pi[5] = 3$

$P_1$   | a | b | a | b | a | c | a |   $\pi[3] = 1$

$P_0$   $\varepsilon$ | a | b | a | b | a | c | a |   $\pi[1] = 0$

(b)

# *32.4 KMP algorithm


(a)


(b)

```
COMPUTE-PREFIX-FUNCTION(P)
 1   m = P.length
 2   let π[1..m] be a new array
 3   π[1] = 0
 4   k = 0
 5   for q = 2 to m
 6       while k > 0 and P[k + 1] ≠ P[q]
 7           k = π[k]
 8       if P[k + 1] == P[q]
 9           k = k + 1
10       π[q] = k
11   return π
```

we have $P_k \sqsupset P[q\text{-}1]$,
check if $P[k+1] \neq P[q]$

$q = 1: k \leftarrow 0, \pi[1] \leftarrow 0$

$q = 2: \pi[2] \leftarrow 0$



$q = 3: P[1] == P[3]$
$k \leftarrow 1, \pi[3] \leftarrow 1$



$q = 4: P[2] == P[4]$
$k \leftarrow 2, \pi[4] \leftarrow 2$
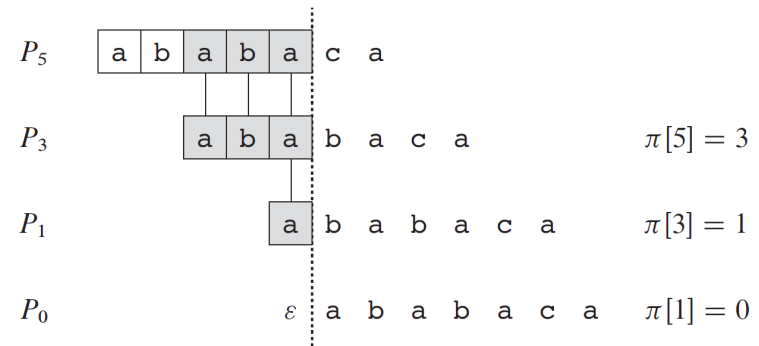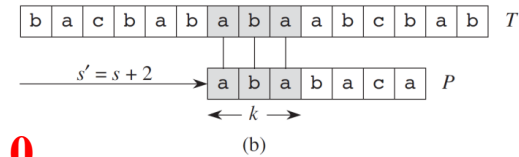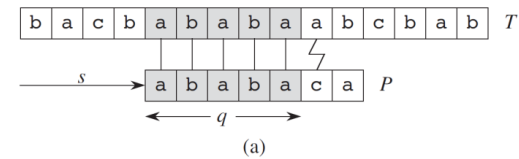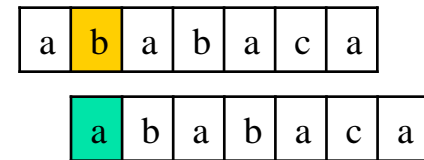


**prefix function** for the pattern $P$ is the function..

$\pi : \{1, 2, \ldots, m\} \rightarrow \{0, 1, \ldots, m\text{-}1\}$ such that
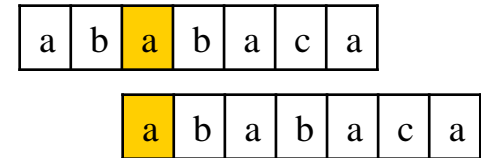$\pi[q] = \max\{k : k < q \text{ and } P_k \sqsupset P_q \}$.

| $i$ | 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 |

# *32.4 KMP algorithm

Running time?　　预习chapter17

Amortized analysis (accounting): $\Theta(n)$, $\Theta(m)$

```
KMP-MATCHER(T, P)
1   n = T.length
2   m = P.length
3   π = COMPUTE-PREFIX-FUNCTION(P)
4   q = 0
5   for i = 1 to n
6       while q > 0 and P[q + 1] ≠ T[i]
7           q = π[q]
8       if P[q + 1] == T[i]
9           q = q + 1
10      if q == m
11          print "Pattern occurs with shift" i − m
12          q = π[q]
```

```
COMPUTE-PREFIX-FUNCTION(P)
1   m = P.length
2   let π[1..m] be a new array
3   π[1] = 0
4   k = 0
5   for q = 2 to m
6       while k > 0 and P[k + 1] ≠ P[q]
7           k = π[k]
8       if P[k + 1] == P[q]
9           k = k + 1
10      π[q] = k
11  return π
```

**prefix function:**

$\pi[q] = \max\{k : k < q \text{ and } P_k \sqsupset P_q\}$

| $i$ | 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)

$P_5$ 　 a b a b a c a

$P_3$ 　 a b a b a c a 　 $\pi[5] = 3$

$P_1$ 　 a b a b a c a 　 $\pi[3] = 1$

$P_0$ 　 $\varepsilon$ a b a b a c a 　 $\pi[1] = 0$

(b)

KMP algorithm avoids computing the transition function $\delta$, and its matching time is $\Theta(n)$ using just an auxiliary function $\pi$, which we precompute from the pattern in time $\Theta(m)$ and store in an array $\pi[1 .. m]$.

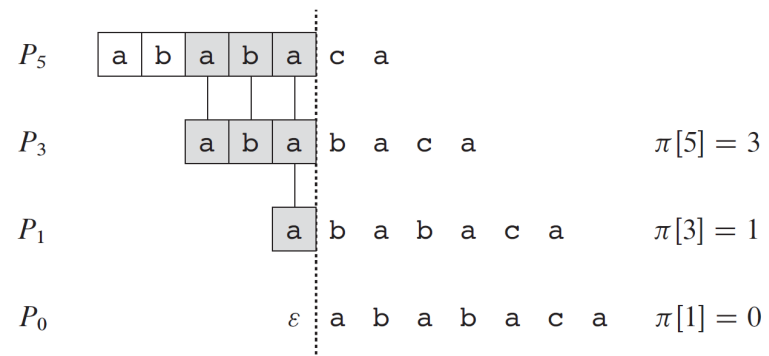**prefix function:**

$$\pi[q] = \max\{k : k < q \text{ and } P_k \sqsupset P_q \}$$

| $i$ | 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)

$P_5$    a b a b a c a

$P_3$    a b a b a c a      $\pi[5] = 3$

$P_1$    a b a b a c a      $\pi[3] = 1$

$P_0$    $\varepsilon$ a b a b a c a      $\pi[1] = 0$

(b)

# Exercises

**32.1-2**

**32.3-1**

**32.3-2**

**32.4-1**

Compute the prefix function $\pi$ for the pattern:
(1) **aabaaaabab**
(2) **abbbabbabbabbbabbabb**

**32.4-7**



| $i$ | 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)

$P_5$ : a b a b a c a

$P_3$ : a b a b a c a    $\pi[5] = 3$

$P_1$ : a b a b a c a    $\pi[3] = 1$

$P_0$ : $\varepsilon$ a b a b a c a    $\pi[1] = 0$

(b)