# CMP(N)302: Design and Analysis of Algorithms

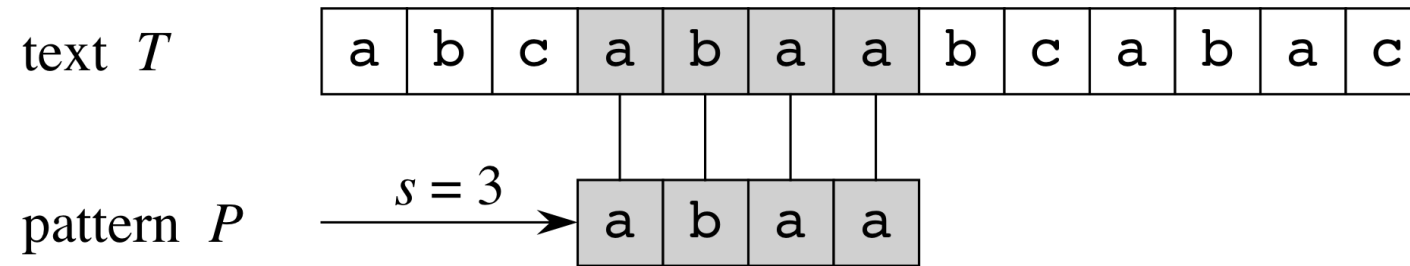## Lecture 11: String Matching

Ahmed Hamdy

Computer Engineering Department

Cairo University

Fall 2019

# String matching

- Simply put, find *all occurrences* of string called *pattern P* (of length $m$) inside another one called *text T* (of length $n$)



- Can be viewed as find the *shift s* ($0 \leq s \leq n - m$) by which the *P* appears in *T*

# String matching algorithms

- Performance of algorithms
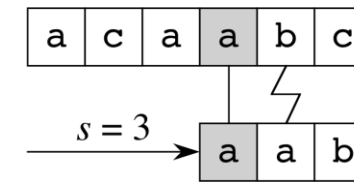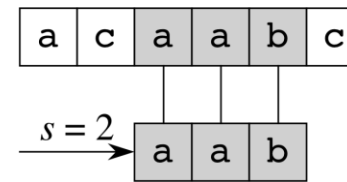
- $\Sigma$ denotes the alphabet
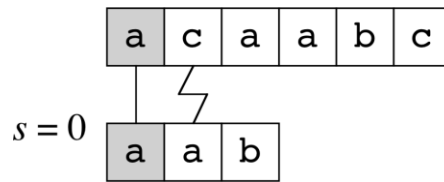
| Algorithm | Preprocessing time | Matching time |
|---|---|---|
| Naive | 0 | $O((n - m + 1)m)$ |
| Rabin-Karp | $\Theta(m)$ | $O((n - m + 1)m)$ |
| Finite automaton | $O(m\,\lvert\Sigma\rvert)$ | $\Theta(n)$ |
| Knuth-Morris-Pratt | $\Theta(m)$ | $\Theta(n)$ |

- When is each algorithm suitable??
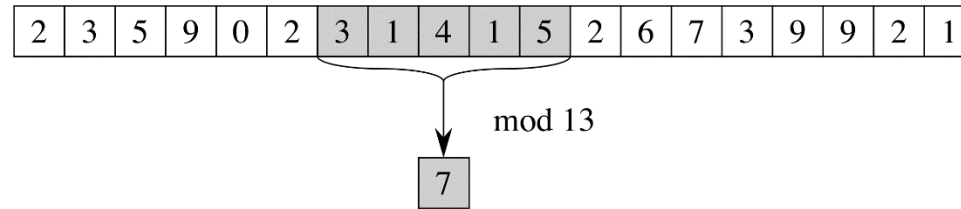
# Naive string-matching

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" $s$

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

$s = 0$    | a | a | b |

(a)

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

$s = 1$    | a | a | b |

(b)

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

$s = 2$    | a | a | b |

(c)

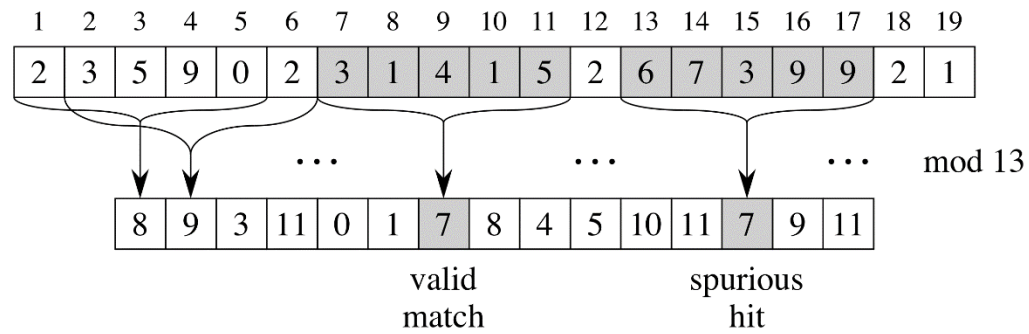| a | c | a | a | b | c |
|---|---|---|---|---|---|

$s = 3$    | a | a | b |

(d)

- Worst case running time $O((n - m + 1)m)$ which is $O(n^2)$ if $m = \lfloor n/2 \rfloor$

- Room for optimization where the algorithm does not make use of information from previous iteration
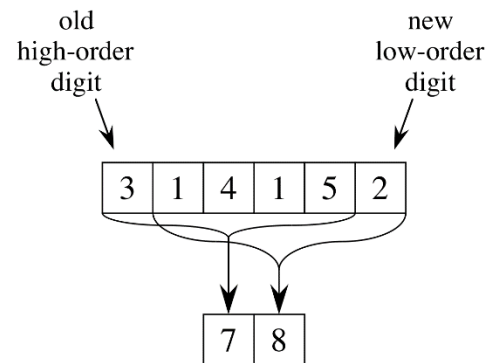
# Rabin-Karp algorithm

| 2 | 3 | 5 | 9 | 0 | 2 | 3 | 1 | 4 | 1 | 5 | 2 | 6 | 7 | 3 | 9 | 9 | 2 | 1 |

mod 13

| 7 |

(a)

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |

| 2 | 3 | 5 | 9 | 0 | 2 | 3 | 1 | 4 | 1 | 5 | 2 | 6 | 7 | 3 | 9 | 9 | 2 | 1 |

$\cdots$  $\cdots$  $\cdots$  mod 13

| 8 | 9 | 3 | 11 | 0 | 1 | 7 | 8 | 4 | 5 | 10 | 11 | 7 | 9 | 11 |

valid match

spurious hit

(b)

old high-order digit

new low-order digit

| 3 | 1 | 4 | 1 | 5 | 2 |

| 7 | 8 |

old high-order digit  shift  new low-order digit

$$14152 \equiv (31415 - 3 \cdot 10000) \cdot 10 + 2 \ (\mathrm{mod}\ 13)$$
$$\equiv (7 - 3 \cdot 3) \cdot 10 + 2 \ (\mathrm{mod}\ 13)$$
$$\equiv 8 \ (\mathrm{mod}\ 13)$$
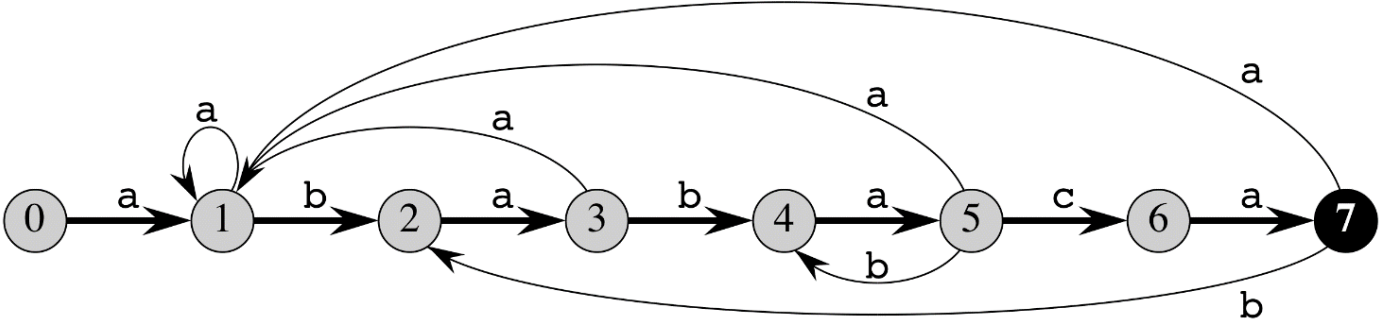
(c)

# Rabin-Karp algorithm

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

```
1   n = T.length
2   m = P.length
3   h = d^{m-1} mod q
4   p = 0
5   t_0 = 0
6   for i = 1 to m                         // preprocessing
7           p = (dp + P[i]) mod q
8           t_0 = (dt_0 + T[i]) mod 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]) mod q
```

- Though worst case is not better than the naïve algorithm, the average case is much better, typically $O((n - m + 1) + cm) = O(n + m)$

# Finite automata

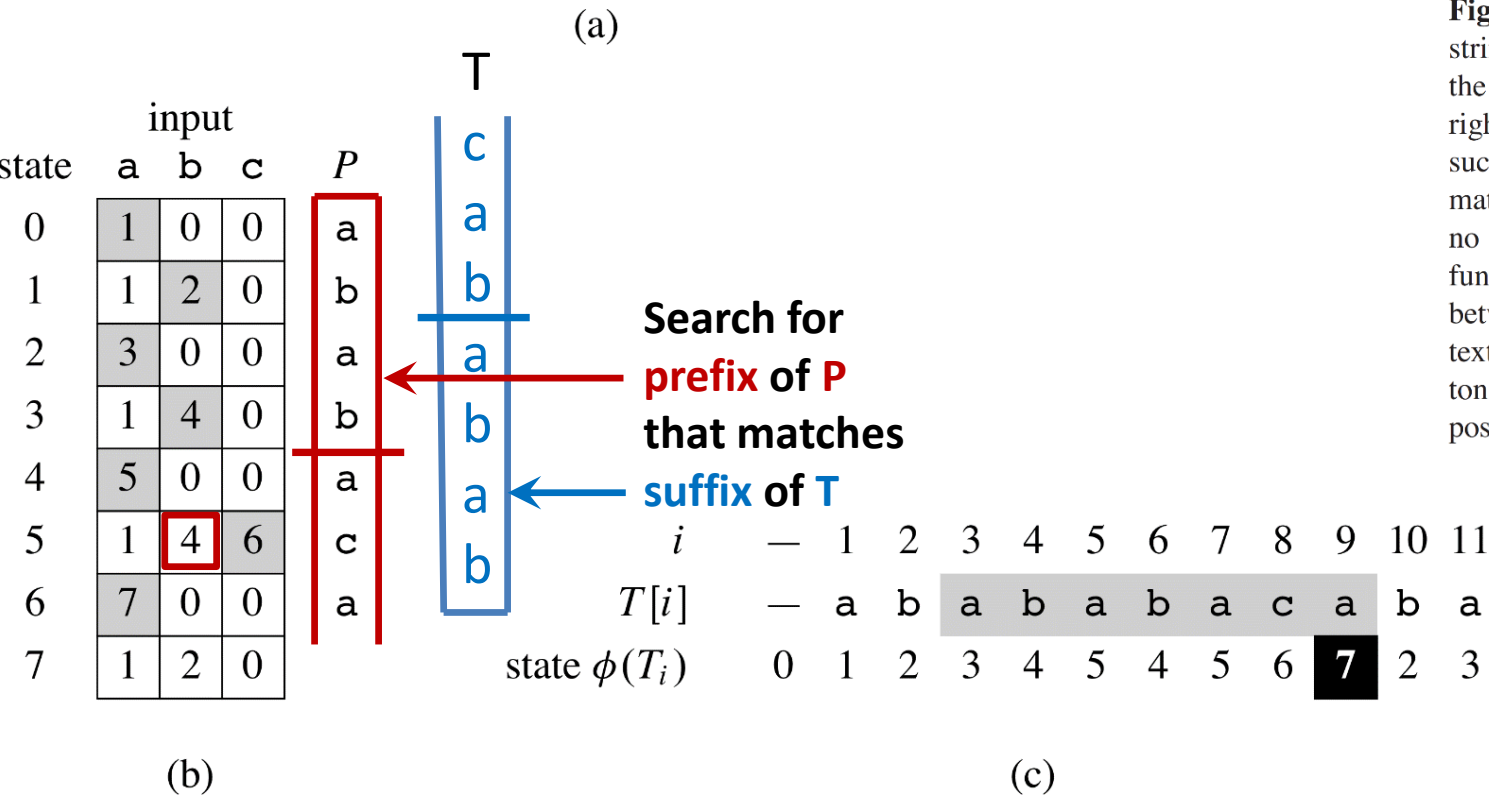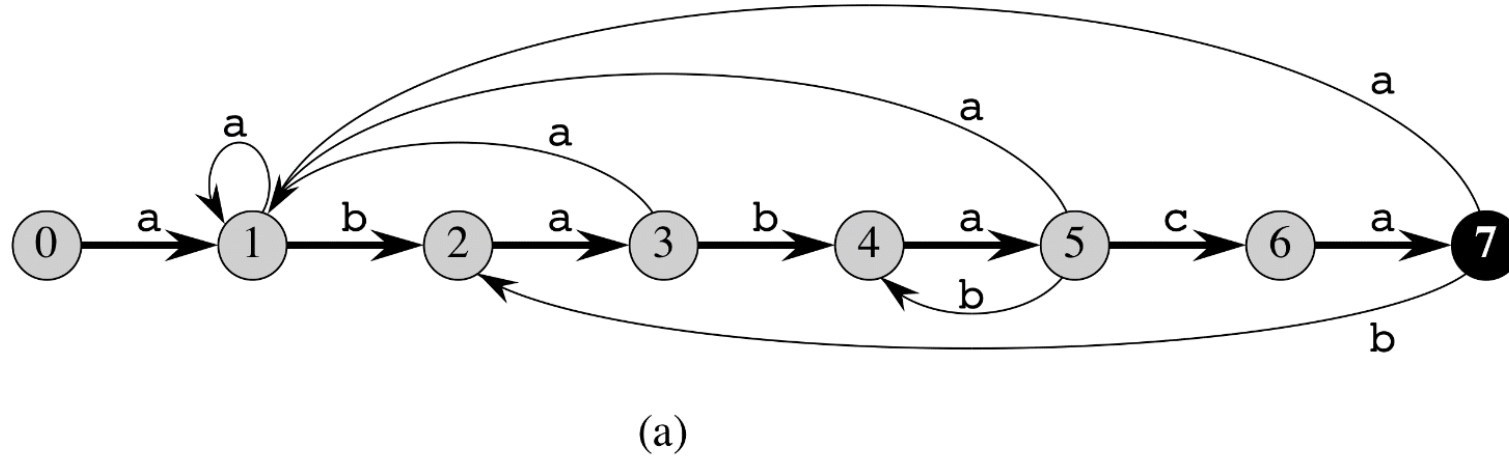- Complexity: $\Theta(n)$



(a)

**Figure 32.7** **(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 $a$ represents $\delta(i, a) = j$. The right-going edges forming the "spine" of the automaton, shown heavy in the figure, correspond to successful matches between pattern and input characters. The left-going edges correspond to failing matches. Some edges corresponding to failing matches are omitted; by convention, if a state $i$ has no outgoing edge labeled $a$ for some $a \in \Sigma$, then $\delta(i, a) = 0$. **(b)** The corresponding transition function $\delta$, 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. Under each text character $T[i]$ appears the state $\phi(T_i)$ that the automaton is in after processing the prefix $T_i$. The automaton finds one occurrence of the pattern, ending in position 9.

(b)

**Search for prefix of P that matches suffix of T**

(c)

# Finite Automata



(a)



(b)

When $q = 5$, initially $k = q + 2 - 1 = 6$

Search for
**prefix** of **P**
that matches
**suffix** of **T**

| $i$ | — | 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)

COMPUTE-TRANSITION-FUNCTION$(P, \Sigma)$

1  $m = P.length$
2  **for** $q = 0$ **to** $m$
3      **for** each character $a \in \Sigma$
4          $k = \min(m + 1, q + 2)$
5          **repeat**
6              $k = k - 1$
7          **until** $P_k \sqsupset P_q a$
8          $\delta(q, a) = k$
9  **return** $\delta$

# Finite Automata



(a)

COMPUTE-TRANSITION-FUNCTION$(P, \Sigma)$
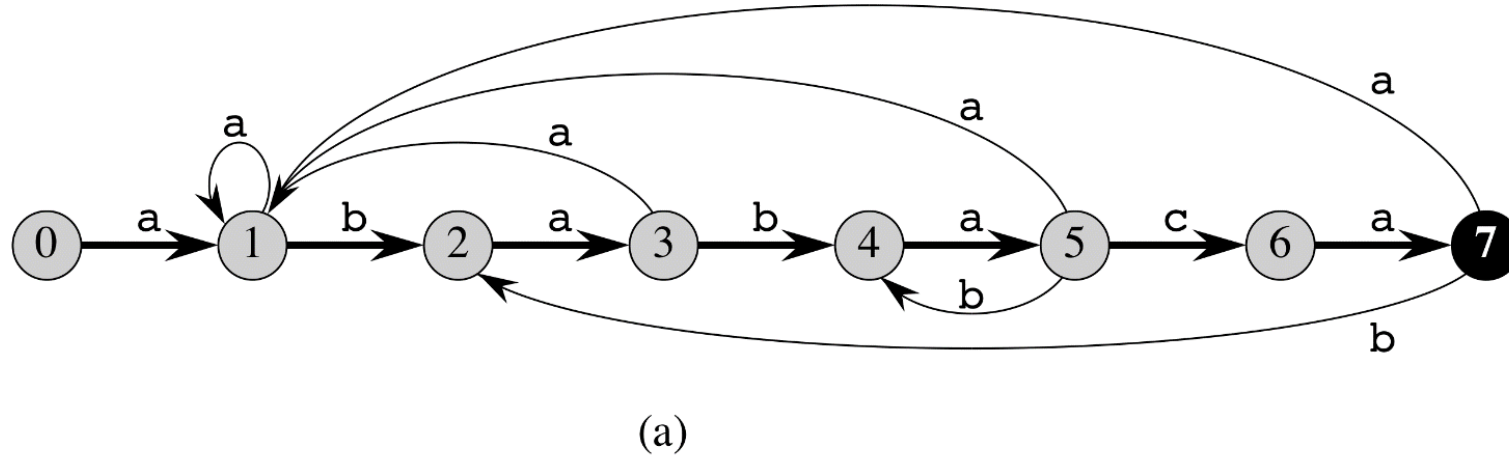
1  $m = P.length$
2  **for** $q = 0$ **to** $m$
3      **for** each character $a \in \Sigma$
4          $k = \min(m + 1, q + 2)$
5          **repeat**
6              $k = k - 1$
7          **until** $P_k \sqsupset P_q a$
8          $\delta(q, a) = k$
9  **return** $\delta$

|       | 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)

T: c a b a b a b

When $q = 5$, initially $k = k - 1 = 5$

Search for
**prefix** of **P**
that matches
**suffix** of **T**

| $i$          | — | 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 Automata



(a)

Compute-Transition-Function$(P, \Sigma)$
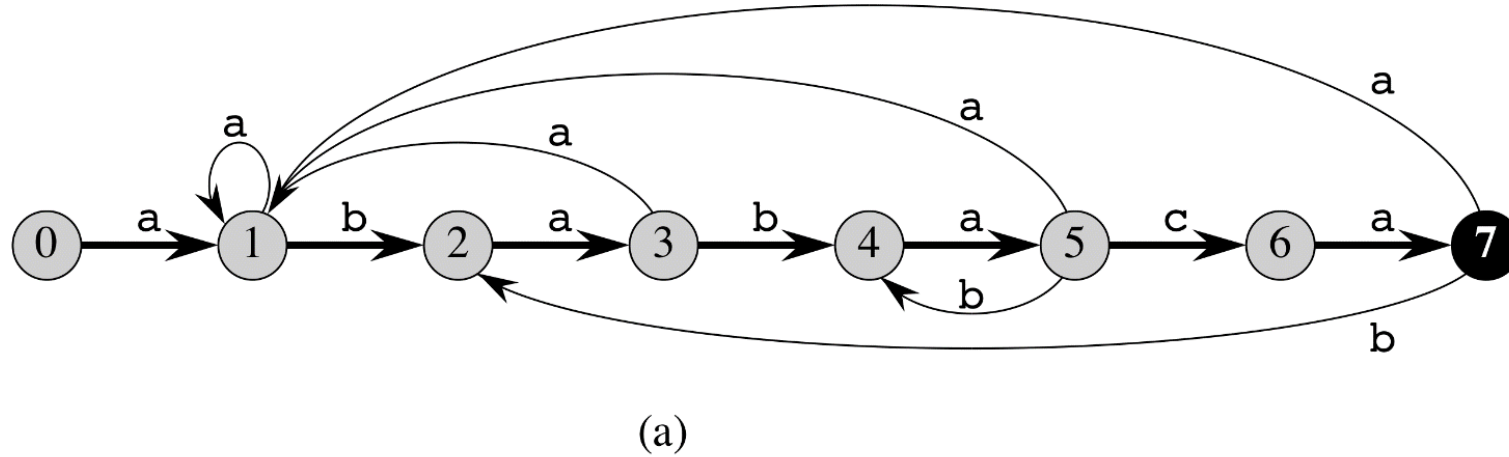
1  $m = P.length$
2  **for** $q = 0$ **to** $m$
3      **for** each character $a \in \Sigma$
4          $k = \min(m + 1, q + 2)$
5          **repeat**
6              $k = k - 1$
7          **until** $P_k \sqsupset P_q a$
8          $\delta(q, a) = k$
9  **return** $\delta$

| 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 | a |

input

T: c a b a b a b

When $q = 5$, initially $k = k - 1 = 4$

Search for **prefix** of **P** that matches **suffix** of **T**

| $i$ | — | 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 |

(b)                                     (c)

# Knuth-Morris-Pratt algorithm

| $i$ | 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)

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

$s \longrightarrow$   a b a b a c a   $P$

$\longleftarrow q \longrightarrow$

(a)

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

$s' = s + 2 \longrightarrow$   a b a b a c a   $P$

$\longleftarrow k \longrightarrow$

(b)

a b a b a   $P_q$

a b a   $P_k$

(c)

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

# Knuth-Morris-Pratt algorithm

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]   // While next character does not match, keep backtracking
7           k = π[k]                           through the π table
8       if P[k+1] == P[q]   // If next character matches,
9           k = k + 1        increment the matching size k
10      π[q] = k
11  return π
```

| $i$ | 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)