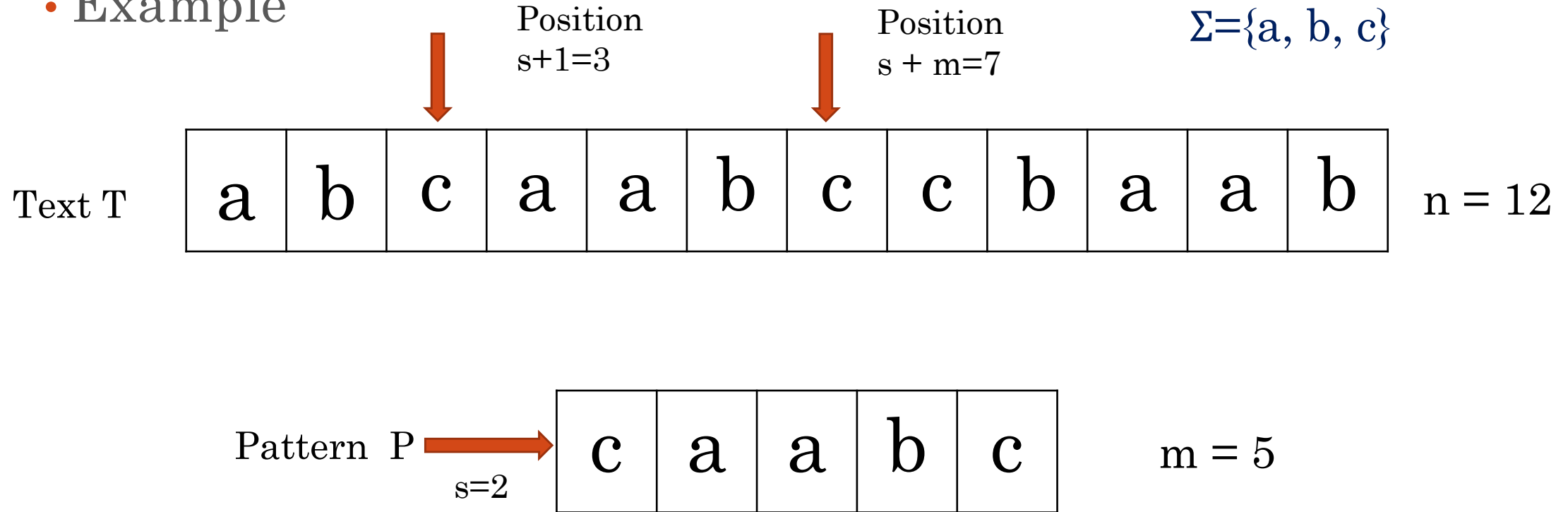# String Matching

# Introduction

- Naïve String Matching Algorithm
- String Matching with Finite Automata
- Knuth Morris Pratt Algorithm

# Naïve String Matching Algorithm

- String matching or pattern recognition is a problem for searching a pattern to be searched within a text under certain conditions and find out all occurrences of it.

- Pattern and text will be in form of an array of characters drawn from finite alphabet $\Sigma$.

- Pattern is denoted as $P[1...m]$ and Text as $T[1...n]$ where $m$ and $n$ are their respective length such that $n \geq m \geq 1$.

- If pattern P occurs in Text after $\boldsymbol{s}$ shifts then $\boldsymbol{P[1...m] = T[s+1...s+m]}$ where $\boldsymbol{n - m \geq s \geq 0}$.

- If P occurs after finite shift $\boldsymbol{s}$ in T, then we can say $\boldsymbol{s}$ is a valid shift, otherwise invalid shift

# Naïve String Matching Algorithm

- Example

Position
s+1=3

Position
s + m=7

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

Text T

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

n = 12

Pattern P

s=2

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

m = 5

# Naïve String Matching Algorithm

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" $s$

# String Matching Algorithm

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

String-matching algorithms, their preprocessing and matching times

# String Matching with Finite Automata

Finite Automata

A finite automaton $\mathbf{M}$ is a 5-tuple $(Q, \Sigma, \delta, s, F)$:

$\mathbf{Q}$: the finite set of states

$\boldsymbol{\Sigma}$ : the finite input alphabet

$\boldsymbol{\delta}$ : the "transition function of $\mathbf{M}$" from $Q \times \Sigma$ to $Q$
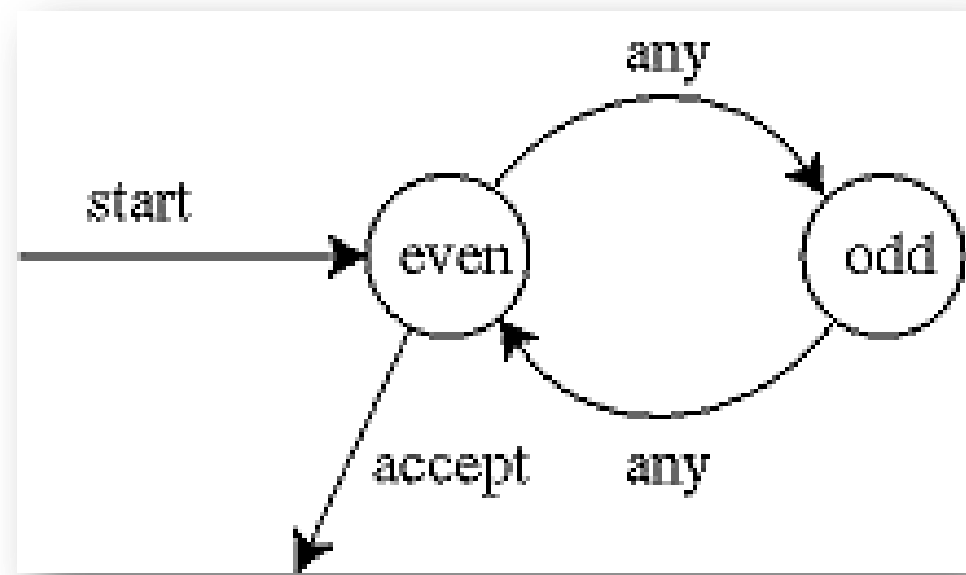
$s \in \mathbf{Q}$: the start state

$\mathbf{F} \subset \mathbf{Q}$: the set of final (accepting) states

# String Matching with Finite Automata

## How it works

A finite automaton accepts strings in a specific language. It begins in state $q_0$ and reads characters one at a time from the input string. It makes transitions from state **q to** $\delta$**(q,a)** based on these characters, and if when it reaches the end of the tape it is in one of the accept states, that string is accepted by the language.

# String Matching with Finite Automata

## How it works

- A finite automaton M induces a function $\phi$, called **the final-state function** , from $\Sigma^*$ to Q such that $\phi(w)$ is the state M ends up in after scanning the string $w$.

- Thus, M accepts a string $w$ if and only if $\phi(w)\epsilon$A.

- We define the function $\phi$ recursively, using the transition function:

$$\phi(\varepsilon) = q_0,$$
$$\phi(wa) = \delta(\phi(w), a) \ for \ w\epsilon\ \Sigma^*, a\ \epsilon\ \Sigma$$

# String Matching with Finite Automata

## String-Matching Automata

- For a given pattern P, we construct a string-matching automaton in a preprocessing step before using it to search the text string.

- In order to specify the string-matching automaton corresponding to a given pattern P [1…m], we first define an auxiliary function $\sigma$, called **the suffix function** corresponding to P.

- A suffix function w.r.t. pattern P[1..m], s , is a mapping from $\Sigma^*$ to {0,1,...,m} such that $\sigma(x)$ is the length of the longest prefix of P that is a suffix of x: $\sigma(x)=\max\{k: P_k \sqsupset x\}$.

# String Matching with Finite Automata

## The Suffix Function

In order to properly search for the string, the program must define a **suffix function (σ)** which checks to see how much of what it is reading matches the search string at any given moment.

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

$$P = \text{abaabc}$$

$$P_1 = \text{a}$$

$$P_2 = \text{ab}$$

$$P_3 = \text{aba}$$

$$P_4 = \text{abaa}$$

$$\sigma(\text{abbaba}) = \text{aba}$$

# String Matching with Finite Automata

## String-Matching Automata

- A suffix function is well defined since the empty string $P_0 = \varepsilon$ is a suffix of every string.

- For example, $P=ab$, $P_0 = \varepsilon$, $\sigma(\varepsilon)=0$, $\sigma(ccaca)=1$, $\sigma(ccab)=2$.

- For $P[1..m]$, $\sigma(x)=m$ if and only if $P \sqsupset x$ (* a valid shift *). The whole pattern is the suffix of x.

# String Matching with Finite Automata

## String-Matching Automata

- For any pattern P of length m, we can define its string-matching automata:

    $Q = \{0,\ldots,m\}$  (states)

    $q_0 = 0$  (start state)

    $F = \{m\}$  (accepting state)

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

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

- We define $\delta(q,a) = \sigma(P_q a)$ because we want to keep track of longest prefix of the pattern P that has matched the string T so far.

# String Matching with Finite Automata

## String-Matching Automata

- We consider the most recently read characters of T .

- In order, for a substring of T (let's say the substring ending at T[i]) to match some prefix $P_j$ of P, this prefix $P_j$ must be a suffix of $T_i$ .

- Suppose that q= $\phi(T_i)$, so that after reading $T_i$ , the automaton is in state q.

- We design the transition function $\delta$ so that this state number, q, tells us the length of the longest prefix of P that matches a suffix of $T_i$ . That is, in state q, $P_q \sqsupset T_i$ and q=$\sigma(T_i)$.

- Whenever q=m, all m characters of P match a suffix of Ti , and so we have found a match.

- Thus, since $\phi$ (Ti) and $\sigma$(Ti) both equal q, we shall see that the automaton maintains the following invariant:

$\phi$ (Ti) and $\sigma$(Ti)

# String Matching with Finite Automata

## Finite-Automaton-Matcher

The simple loop structure implies a running time for a string of length $n$ is O(n).

*However*: this is only the running time for the actual string matching. It does not include the time it takes to compute the transition function.

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     **if** $q = m$
6        **then** $s \leftarrow i - m$
7           print "Pattern occurs at shift" $s$

# String Matching with Finite Automata

## Computing the Transition Function

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

This procedure computes $\delta(q,a)$ according to its definition. The loop on line 2 cycles through all the states, while the nested loop on line 3 cycles through the alphabet. Thus, all state-character combinations are accounted for. Lines 4-7 set $\delta(q,a)$ to be the largest $k$ such that $P_k \sqsupset P_q a$.

# String Matching with Finite Automata

**Example :**

- **P= a b a b a c a**

- q= 3 (implies text is :a b a)        (step 2)

- $a \leftarrow \Sigma$                        (step 3)

   k = min(7+1, 3+2)=5, k-1=4,  (steps 4,5)

   $p_4 \sqsupset p_3.a$ ? No. k ← k-1=3        (step 5)

   $p_3 \sqsupset p_2.a$ ? Yes;  **δ(2,a) ← 3**    (steps 6,7)

- b← $\Sigma$                          (step 3)

   k = min(7+1, 3+2)=5, k-1=4,  (steps 4,5)

   $p_4 \sqsupset p_3.b$ ? Yes; **δ(3,b) = 4**    (steps 6,7)

- Similarly for c← $\Sigma$ ; **δ(3,c) = 0**

Compute-Transition-Function (P,Σ)                    1

$m$ ← length[P]                          2

For $q$ ← 0 to $m$                          3

    do for each character $a \in \Sigma$        4

        do k ← min($m$+1, $q$+2)                5

            repeat $k$ ← $k$-1                6

                until $P_k \sqsupset P_q a$            7

            δ($q,a$) ← $k$                    8

return δ                              9

# String Matching with Finite Automata

- This procedure builds $\delta(q,a)$ is a straight-forward way by definition. It considers all states q and all characters in $\Sigma$.

- For each combination, to find the the largest k such that $P_k \sqsupset P_q a$. The worst case time complexity is $O(m^3|\Sigma|)$

```
Compute-Transition-Function (P,Σ)

    m ← length[P]

    For q ← 0 to m

            do for each character a ∈ Σ

                    do k ← min(m+1, q+2)

                            repeat k ← k-1

                                    until P_k ⊃ P_q a

                            δ(q,a) ← k

    return δ
```

# KMP Algorithm

- KMP is the first linear time algorithm for string matching.

- Prevents re examination of previously matched characters.

- This algorithm avoids computing the transition function $\delta$ altogether, and its matching time is $\theta(n)$ using just an auxiliary function $\pi$.

- $\pi[q]$ (Prefix Table ot LPS Table) stores information that is needed to compute transition function $\delta(q,a)$ but that does not depend on a.

- array $\pi[q]$ has only m entries, whereas $\delta$ has $\theta(m|\sum|)$.

# KMP Algorithm

KMP-MATCHER$(T, P)$

1   $n = T.length$
2   $m = P.length$
3   $\pi = $ 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

# KMP Algorithm

COMPUTE-PREFIX-FUNCTION($P$)

1  $m = P.length$
2  let $\pi[1 \dots m]$ 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** $\pi$

The running time of COMPUTE-PREFIX-FUNCTION is $\Theta(m)$