

# Analysis of Algorithms

## String Matching

### Basic String/Document Processing Algorithms

- String-driven information retrieval is at the core of many important computer applications, including:
  - Web searching and “surfing”: the Internet document formats HTML and XML are primarily text formats with added links to multimedia context
  - Searching for a certain DNA sequence in a genomic database, or searching for particular patterns in DNA sequences
- Document processing is one of the dominant functions of computers
  - editing, searching, transporting over the Internet, displaying, etc.

### Text Documents

- From the perspective of algorithm design, documents can be viewed as simple **character strings**, that is, they can be abstracted as a sequence of characters.
- At the heart of algorithms for searching and processing text are methods for dealing with character strings
- $A = \text{“CGTAAACTGCTTTAATCAAACGC”}$  DNA sequence
- $B = \text{“http://www.pucit.edu.pk”}$

## The String Matching Problem

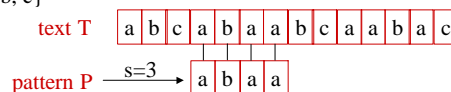
- Finding all occurrences of a pattern in a text
- We assume that
  - the text is an array  $T[1..n]$  of length  $n$  and the pattern is an array  $P[1..m]$  of length  $m$ .
  - the elements of  $T$  and  $P$  are characters drawn from a finite alphabet  $\Sigma$
- Ex.:  $\Sigma = \{0,1\}$  or  $\Sigma = \{a,b,\dots,z\}$
- The character arrays  $P$  and  $T$  are loosely called “strings of characters”

## The String Matching Problem

**Given:** Two strings  $T[1..n]$  and  $P[1..m]$  over alphabet  $\Sigma$ .

Want to find all occurrences of  $P[1..m]$  “the pattern” in  $T[1..n]$  “the text”.

**Example:**  $\Sigma = \{a, b, c\}$



**Terminology:**

- $P$  occurs with shift  $s$ .
- $P$  occurs beginning at position  $s+1$ .
- $s$  is a valid shift.

**Goal:** Find all valid shifts with which a given pattern  $P$  occurs in a text  $T$ .

### Notation and Terminology

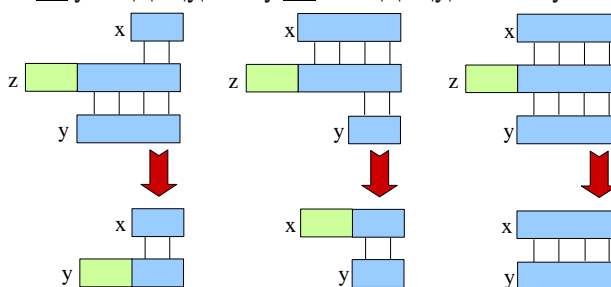
- An Alphabet  $\Sigma$  is a finite set of symbols
- $\Sigma^*$  denotes the set of all finite-length strings over  $\Sigma$ .
  - The empty string is denoted  $\epsilon$
- The length of a string  $x$  is denoted  $|x|$
- The concatenation of two strings  $x$  and  $y$  is denoted  $xy$ , and has length  $|x| + |y|$
- We say that string  $x$  is a **prefix** of string  $y$ , denoted  $x \hat{I} y$ , if  $y = xw$  for some string  $w \in \Sigma^*$ .

### Notation and Terminology

- Ex:  $vzk \hat{I} vzkavk$ ,  $vzk \hat{I} vzk$
- We say that string  $x$  is a **suffix** of string  $y$ , denoted  $x \hat{E} y$ , if  $y = wx$  for some string  $w \in \Sigma^*$ .
- $avk \hat{E} vzkavk$ ,  $avk \hat{E} avk$

### Lemma 32.1

**Lemma 32.1:** Suppose  $x \text{ suf } z$  and  $y \text{ suf } z$ . If  $|x| \leq |y|$  then  $x \text{ suf } y$ . If  $|x| \geq |y|$  then  $y \text{ suf } x$ . If  $|x| = |y|$  then  $x = y$



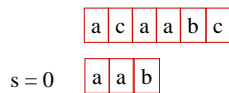
## Naïve Brute-Force Algorithm

```
Naïve(T, P)
  n := length[T];
  m := length[P];
  for s := 0 to n - m do
    if P[1..m] = T[s+1..s+m] then
      print "pattern occurs with shift s"
```

9

2/27/2003

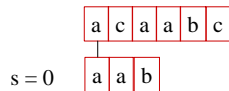
## Example



10

2/27/2003

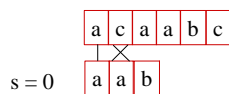
## Example



11

2/27/2003

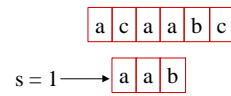
## Example



12

2/27/2003

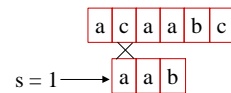
## Example



13

2/27/2003

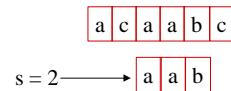
## Example



14

2/27/2003

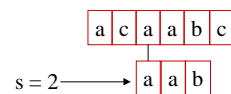
## Example



15

2/27/2003

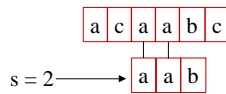
## Example



16

2/27/2003

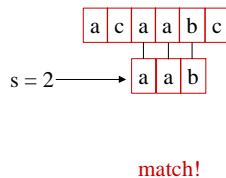
## Example



17

2/27/2003

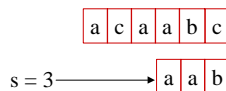
## Example



18

2/27/2003

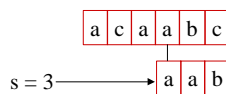
## Example



19

2/27/2003

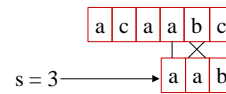
## Example



20

2/27/2003

## Example



Running time is  $\Theta((n - m + 1)m)$ .

21

2/27/2003

## String Matching with Finite Automata

- Many string matching algorithms build a finite automaton that scans  $T$  for all occurrences of  $P$
- String matching automata are very efficient, since they examine each character only once
- Running time (excluding the time to build the automaton) is  $\Theta(n)$
- The time to build the automaton can however be large, if  $\Sigma$  is large...

22

2/27/2003

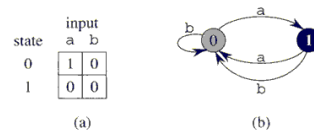
## Finite Automata

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

23

2/27/2003

## Example



**Figure 32.6** A simple two-state finite automaton with input alphabet  $\Sigma = \{a, b\}$ . (a) A tabular representation of the transition function. (b) A state-transition diagram. State 1 is the only accepting state. This automaton accepts those strings that end in an even number of 'a's. For example, the sequence of states this automaton goes through on the input "aaba" is  $(0, 1, 0, 1, 0, 1)$ , and so it accepts this input. It rejects the input "abab" because the sequence of states is  $(0, 1, 0, 0, 1, 0)$ , and so it rejects this input.

$Q = \{0, 1\}$   
 $q_0 = 0$   
 $A = \{1\}$   
 $\Sigma = \{a, b\}$

24

2/27/2003

## Final State Function

- The FA begins with state  $q_0$  and reads the characters of the input string one at a time.
- If the FA is in state  $q$  and reads input character  $c$  it makes a transition from state  $q$  to state  $\delta(q, c)$ .
- When the current state  $q \in A$ , the machine  $M$  is said to have **accepted** the string read so far. (An input that is not accepted is said to be **rejected**.)
- function  $\phi$ , called the **final state function** from  $\Sigma^*$  to  $Q$  such that  $\phi(w)$  is the state  $M$  ends up after scanning the string  $w$ , FA ends up in after scanning the string  $w$ .
- Thus  $M$  accepts a string if and only if  $\phi(w) \in A$ .

25

2/27/2003

## Final State Function

- Final state function**  $\phi(w)$  - recursive definition
- $\phi(\epsilon) = q_0$
- $\phi(wc) = \delta(\phi(w), c)$  for  $w \in \Sigma^*$ ,  $c \in \Sigma$

26

2/27/2003

## Suffix Function

- A string-matching automaton is constructed for pattern  $P$ , and then used to search the text string  $T$ .
- First, an auxiliary **suffix function**  $\sigma$  is defined for  $P[1..m]$ .
  - The suffix function is a mapping from  $\Sigma^*$  to  $\{0, 1, \dots, m\}$  such that  $\sigma(x)$  is the length of the longest prefix of  $P$  that is a suffix of  $x$ .
  - $\sigma(x) = \max \{k : P_k \supset x\} \ (x \Rightarrow T_i)$

27

2/27/2003

## Example

- $P = ab \Rightarrow \sigma(\epsilon) = 0$
- $P = ab \Rightarrow \sigma(ccaca) = 1$
- $P = ab \Rightarrow \sigma(ccab) = 2$
- Empty string  $P_0 = \epsilon$  is a suffix of every string.

28

2/27/2003

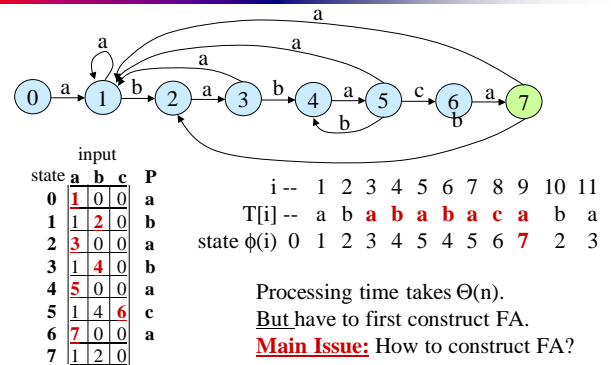
## String Matching Automata

- The string-matching automaton constructed for pattern  $P[1..m]$  is defined as follows:
  - The state set  $Q$  is  $\{0, 1, \dots, m\}$ . The start state  $q_0$  is state 0 and state  $m$  is the only accepting state.
  - The transition function  $\delta$ , for any state  $q$  and character  $c$ , is defined by
 
$$\delta(q, c) = \sigma(P_q c)$$
    - That is, in order to compute the length of the longest suffix of  $T_i c$  that is a prefix of  $P$ , we can compute the longest suffix of  $P_q c$  that is a prefix of  $P$ .

29

2/27/2003

## Finite Automata Algorithm



30

2/27/2003

## Finite-Automaton-Matcher

Finite-Automaton( $T, \delta, m$ )

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

$q \leftarrow 0$

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

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

**if**  $q = m$

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

- Running time is  $O(n)$

31

2/27/2003