

# String Searching



# String Searching

- String searching is to find all the occurrences of a pattern in a text.
- **String Searching Domains**
  - Text editing
  - DNA (Deoxyribo Nucleic Acid) sequence matching
  - Web searching
  - Spelling and Grammar checking
  - Compression etc.

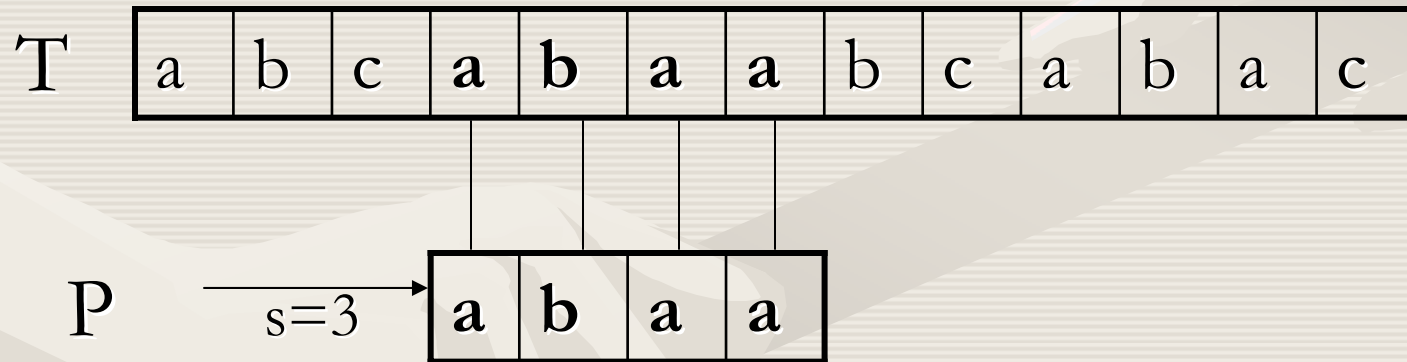
# Notations and Terminology

- $T[1 \dots n]$ ; array of text having length  $n$
- $P[1 \dots m]$ ; array of pattern to be searched in text  $T$  having length  $m$ , where  $m \leq n$
- $\Sigma$  is the finite set of alphabets.  
e.g.  $\Sigma = \{0, 1\}$  or  $\Sigma = \{a, b, \dots, z\}$ , etc.
- $P$  and  $T$  are often called strings of characters

# Notations and Terminology

- Valid & Invalid Shift

If  $P$  occurs in  $T$  beginning at position  $s+1$ , where  $0 \leq s \leq n - m$  and  $T[s+1.. s+m]$ , then we call ' $s$ ' a Valid Shift otherwise ' $s$ ' is called an Invalid Shift.



# Notations and Terminology

- $\Sigma^*$  (sigma star) denote the set of all finite length strings formed using characters from the alphabet set  $\Sigma$
- $\varepsilon$  denotes the Empty (zero length) String, also belongs to  $\Sigma^*$
- $|x|$  denotes the Length of string  $x$
- $xy$  denotes the Concatenation of two strings  $x$  and  $y$  having length  $|x| + |y|$

# Notations and Terminology

- **Prefix of a String**

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

- If  $\omega \sqsubseteq x$  then  $|\omega| \leq |x|$

- **Suffix of a String**

$\omega$  is a suffix of a string  $x$ , denoted by  $\omega \sqsupseteq x$ , if  $x = y\omega$  for some string  $y \in \Sigma^*$

- If  $\omega \sqsupseteq x$  then  $|\omega| \leq |x|$

# Notations and Terminology

- The empty string  $\epsilon$  is both a suffix and a prefix of every string
- For any strings  $x$  and  $y$  and any character  $a$ , we have  $x \sqsubset y$  if and only if  $xa \sqsubset ya$
- $\sqsubset$  and  $\sqsubseteq$  are transitive operations.

# Notations and Terminology

- **Lemma**

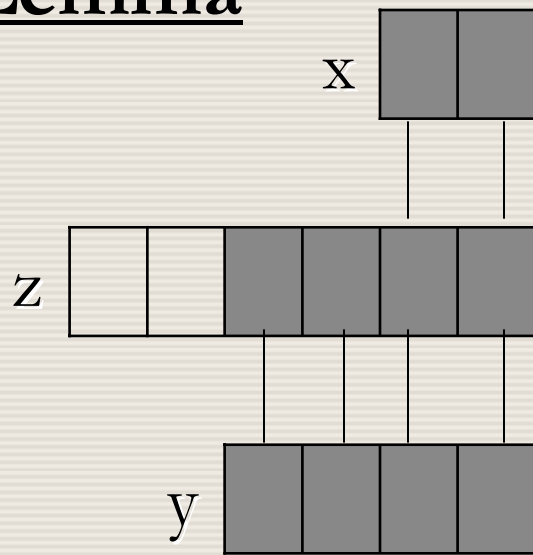
Suppose that  $x$ ,  $y$  and  $z$  are strings such that  $x \sqsubset z$  and  $y \sqsubset z$

- a) If  $|x| \leq |y|$  then  $x \sqsubset y$
- b) If  $|x| \geq |y|$  then  $y \sqsubset x$
- c) If  $|x| = |y|$  then  $x = y$

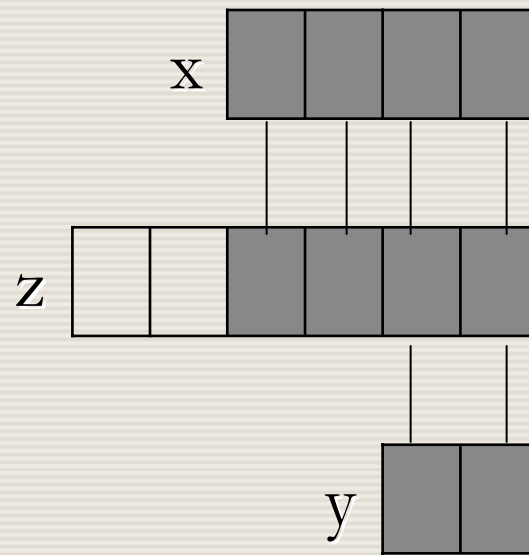


# Notations and Terminology

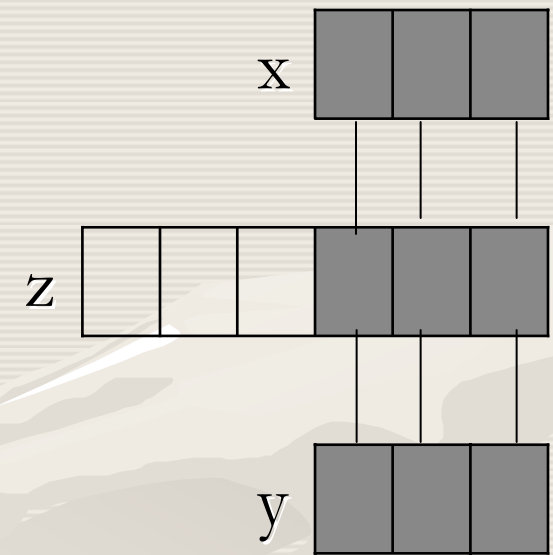
## Lemma



(a)



(b)



(c)

# Algorithm Categories

String Matching Algorithms



# Naïve String Matching Algorithm

- The naïve string-matching procedure can be interpreted graphically as sliding a “template” containing the pattern over the text.
- Noting for which shifts all of the characters on the template equal the corresponding characters in the text.
- The whole string is scanned character by character without any optimization or pre-processing.
- Also called **Brute Force Pattern Matching**.

# Naïve String Matching Algorithm

## NAÏVE-STRING-MATCHER (T, P)

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

$m \leftarrow \text{length}[P]$

**for**  $s \leftarrow 0$  **to**  $n - m$

**do if**  $P[1 .. m] = T[s+1 .. s+m]$

**then print** “Pattern occurs at shift”  $s$

# Execution of Naïve Algorithm

- $T = \text{abababacaba}$  (Text)
- $P = \text{ababaca}$  (Pattern)
- $n = 11$  (Length of Text)
- $m = 7$  (Length of Pattern)

# Execution of Naïve Algorithm

$n-m = 4$ , so loop will execute  $(0 - 4)$  5 times

	1	2	3	4	5	6	7	8	9	10	11
<b>T</b>	a	b	a	b	a	b	a	c	a	b	a

#	s	$P[1..m] = T[s+1..s+m]$	Output
1	0	ababaca = abababa ( <b>false</b> )	
2	1	ababaca = bababac ( <b>false</b> )	
3	2	ababaca = ababaca ( <b>true</b> )	Pattern Occurs at Shift 2
4	3	ababaca = babacab ( <b>false</b> )	
5	4	ababaca = abacaba ( <b>false</b> )	

# Properties of Naïve Algorithm

- No Pre-processing is involved
- Total running time is  $O((n-m+1)m)$
- No track of previously read characters
- May be suitable for smaller texts
- Performance decreases as the text grows larger
- No character set dependence
- No optimization has been performed

# Finite Automata

- Build and use a finite-automaton that scans the text string  $T$  for all occurrences of the pattern  $P$ .
- String automata are efficient: it examines each text character exactly once, taking constant time per text character. But the time to build the automaton, however, can be large if  $\sum$  is large.
- The automaton must be constructed from the pattern in a preprocessing step before it can be used to search the text string.



# Finite Automata

- A **Finite Automata** **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 finite **input alphabet (character set)**
- $\delta$  is a function from  $Q \times \Sigma$  into  $Q$ , called the **transition function of M**.
- The finite automaton begins in state  $q_0$  and reads the characters of its input string one at a time. If the automaton is in state  $q$  and reads input character  $a$ , it moves (“makes a transition”) from state  $q$  to state  $\delta(q, a)$

# Finite Automata

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

$m \leftarrow \text{length}[P]$

**for**  $q \leftarrow 0$  **to**  $m$

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

**do**  $k \leftarrow \min(m+1, q+2)$

**repeat**

$k \leftarrow k - 1$

**until**  $P_k \supseteq P_q a$

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

**return**  $\delta$

# Finite Automata

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

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

$q \leftarrow 0$

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

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

**if**  $q = m$

**then** print “Pattern occurs with shift”  $i - m$

# Finite Automata

## Computing the Transition Function

- $P = \text{ababaca}$  (Pattern)
- $m = 7$  (Length of Pattern)
- $\Sigma = \{a, b, c\}$  (Character Set)
- At start computing the transition function which gives us the state transitions depending upon the input character and current state.

# Finite Automata

## Computing the Transition Function

$m = 7$ , so the outer loop will execute

$(0 - 7)$  8 times.

$$P_x = P[1 .. x]$$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$P_k \supset P_q a$	$\delta(q, a)$
0	a	$k = \min(7+1, 0+2) = 2$	$2-1=1$	$a \supset a$ ( <b>T</b> )	$\delta(0, a) = 1$
	b	$k = \min(7+1, 0+2) = 2$	$2-1=1$	$a \supset b$ ( <b>F</b> )	$\delta(0, b) = 0$
			$1-1=0$	$\epsilon \supset b$ ( <b>T</b> )	
	c	$k = \min(7+1, 0+2) = 2$	$2-1=1$	$a \supset c$ ( <b>F</b> )	$\delta(0, c) = 0$
			$1-1=0$	$\epsilon \supset c$ ( <b>T</b> )	
1	a	$k = \min(7+1, 1+2) = 3$	$3-1=2$	$ab \supset aa$ ( <b>F</b> )	$\delta(1, a) = 1$
			$2-1=1$	$a \supset aa$ ( <b>T</b> )	

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$P_x = P[1 .. x]$$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$P_k \supset P_q a$	$\delta(q, a)$
1	b	$k = \min(7+1, 1+2) = 3$	$3-1=2$	$ab \supset ab$ ( <b>T</b> )	$\delta(1, b) = 2$
	c	$k = \min(7+1, 1+2) = 3$	$3-1=2$	$ab \supset ac$ ( <b>F</b> )	
			$2-1=1$	$a \supset ac$ ( <b>F</b> )	
			$1-1=0$	$\epsilon \supset ac$ ( <b>T</b> )	$\delta(1, c) = 0$
2	a	$k = \min(7+1, 2+2) = 4$	$4-1=3$	$aba \supset aba$ ( <b>T</b> )	$\delta(2, a) = 3$
	b	$k = \min(7+1, 2+2) = 4$	$4-1=3$	$aba \supset abb$ ( <b>F</b> )	
			$3-1=2$	$ab \supset abb$ ( <b>F</b> )	
			$2-1=1$	$a \supset abb$ ( <b>F</b> )	
			$1-1=0$	$\epsilon \supset abb$ ( <b>T</b> )	$\delta(2, b) = 0$

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$P_x = P[1 .. x]$$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$P_k \supset P_q a$	$\delta(q, a)$
2	c	$k = \min(7+1, 2+2) = 4$	$4-1=3$ $3-1=2$ $2-1=1$ $1-1=0$	$aba \supset abc$ ( <b>F</b> ) $ab \supset abc$ ( <b>F</b> ) $a \supset abc$ ( <b>F</b> ) $\epsilon \supset abc$ ( <b>T</b> )	$\delta(2, c) = 0$
3	a	$k = \min(7+1, 3+2) = 5$	$5-1=4$ $4-1=3$ $3-1=2$ $2-1=1$	$abab \supset abaa$ ( <b>F</b> ) $aba \supset abaa$ ( <b>F</b> ) $ab \supset abaa$ ( <b>F</b> ) $a \supset abaa$ ( <b>T</b> )	$\delta(3, a) = 1$
	b	$k = \min(7+1, 3+2) = 5$	$5-1=4$	$abab \supset abab$ ( <b>T</b> )	$\delta(3, b) = 4$

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$\mathbf{P}_x = \mathbf{P}[1 \dots x]$$

	1	2	3	4	5	6	7
$\mathbf{P}$	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$\mathbf{P}_k \supset \mathbf{P}_q a$	$\delta(q, a)$
3	c	$k = \min(7+1, 3+2) = 5$	$5-1=4$ $4-1=3$ $3-1=2$ $2-1=1$ $1-1=0$	$abab \supset abac$ ( <b>F</b> ) $aba \supset abac$ ( <b>F</b> ) $ab \supset abac$ ( <b>F</b> ) $a \supset abac$ ( <b>F</b> ) $\epsilon \supset abac$ ( <b>T</b> )	$\delta(3, c) = 0$
4	a	$k = \min(7+1, 4+2) = 6$	$6-1=5$	$ababa \supset ababa$ ( <b>T</b> )	$\delta(4, a) = 5$
	b	$k = \min(7+1, 4+2) = 6$	$6-1=5$ $5-1=4$ $4-1=3$	$ababa \supset ababb$ ( <b>F</b> ) $abab \supset ababb$ ( <b>F</b> ) $aba \supset ababb$ ( <b>F</b> )	



# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$P_x = P[1 .. x]$$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$P_k \supset P_q a$	$\delta(q, a)$
			$3-1=2$	$ab \supset ababb \text{ (F)}$	
			$2-1=1$	$a \supset ababb \text{ (F)}$	
			$1-1=0$	$\epsilon \supset ababb \text{ (F)}$	$\delta(4, b) = 0$
	c	$k = \min(7+1, 4+2) = 6$	$6-1=5$	$ababa \supset ababc \text{ (F)}$	
			$5-1=4$	$abab \supset ababc \text{ (F)}$	
			$4-1=3$	$aba \supset ababc \text{ (F)}$	
			$3-1=2$	$ab \supset ababc \text{ (F)}$	
			$2-1=1$	$a \supset ababc \text{ (F)}$	
			$1-1=0$	$\epsilon \supset ababc \text{ (F)}$	$\delta(4, c) = 0$

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$P_x = P[1 .. x]$$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$P_k \supset P_q a$	$\delta(q, a)$
5	a	$k = \min(7+1, 5+2) = 7$	$7-1=6$	$ababac \supset ababaa$ ( <b>F</b> )	$\delta(5, a) = 1$
			$6-1=5$	$ababa \supset ababaa$ ( <b>F</b> )	
			$5-1=4$	$abab \supset ababaa$ ( <b>F</b> )	
			$4-1=3$	$aba \supset ababaa$ ( <b>F</b> )	
			$3-1=2$	$ab \supset ababaa$ ( <b>F</b> )	
			$2-1=1$	$a \supset ababaa$ ( <b>T</b> )	
	b	$k = \min(7+1, 5+2) = 7$	$7-1=6$	$ababac \supset ababab$ ( <b>F</b> )	$\delta(5, b) = 4$
			$6-1=5$	$ababa \supset ababab$ ( <b>F</b> )	
			$5-1=4$	$abab \supset ababab$ ( <b>T</b> )	

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$\mathbf{P}_x = \mathbf{P}[1 \dots x]$$

	1	2	3	4	5	6	7
$\mathbf{P}$	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$\mathbf{P}_k \supset \mathbf{P}_q a$	$\delta(q, a)$
5	c	$k = \min(7+1, 5+2) = 7$	$7-1=6$	$\text{ababac} \supset \text{ababac} \text{ (T)}$	$\delta(5, c) = 6$
6	a	$k = \min(7+1, 6+2) = 8$	$8-1=7$	$\text{ababaca} \supset \text{ababaca} \text{ (T)}$	$\delta(6, a) = 7$
	b	$k = \min(7+1, 6+2) = 8$	$8-1=7$	$\text{ababaca} \supset \text{ababacb} \text{ (F)}$	
			$7-1=6$	$\text{ababac} \supset \text{ababacb} \text{ (F)}$	
			$6-1=5$	$\text{ababa} \supset \text{ababacb} \text{ (F)}$	
			$5-1=4$	$\text{abab} \supset \text{ababacb} \text{ (F)}$	
			$4-1=3$	$\text{aba} \supset \text{ababacb} \text{ (F)}$	
			$3-1=2$	$\text{ab} \supset \text{ababacb} \text{ (F)}$	
			$2-1=1$	$\text{a} \supset \text{ababacb} \text{ (F)}$	

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$P_x = P[1 .. x]$$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$P_k \supset P_q a$	$\delta(q, a)$
			$1-1=0$	$\epsilon \supset ababacb$ ( <b>T</b> )	$\delta(6, b)=0$
	c	$k = \min(7+1, 6+2)=8$	$8-1=7$	$ababaca \supset ababacc$ ( <b>F</b> )	
			$7-1=6$	$ababac \supset ababacc$ ( <b>F</b> )	
			$6-1=5$	$ababa \supset ababacc$ ( <b>F</b> )	
			$5-1=4$	$abab \supset ababacc$ ( <b>F</b> )	
			$4-1=3$	$aba \supset ababacc$ ( <b>F</b> )	
			$3-1=2$	$ab \supset ababacc$ ( <b>F</b> )	
			$2-1=1$	$a \supset ababacc$ ( <b>F</b> )	
			$1-1=0$	$\epsilon \supset ababacc$ ( <b>T</b> )	$\delta(6, c)=0$

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$P_x = P[1 .. x]$$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$P_k \supset P_q a$	$\delta(q, a)$
7	a	$k = \min(7+1, 7+2) = 8$	$8-1=7$	ababaca $\supset$ ababacaa ( <b>F</b> )	$\delta(7, a) = 1$
			$7-1=6$	ababac $\supset$ ababacaa ( <b>F</b> )	
			$6-1=5$	ababa $\supset$ ababacaa ( <b>F</b> )	
			$5-1=4$	abab $\supset$ ababacaa ( <b>F</b> )	
			$4-1=3$	aba $\supset$ ababacaa ( <b>F</b> )	
			$3-1=2$	ab $\supset$ ababacaa ( <b>F</b> )	
			$2-1=1$	a $\supset$ ababacaa ( <b>T</b> )	
	b	$k = \min(7+1, 7+2) = 8$	$8-1=7$	ababaca $\supset$ ababacab ( <b>F</b> )	
			$7-1=6$	ababac $\supset$ ababacab ( <b>F</b> )	

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$P_x = P[1 .. x]$$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$P_k \supset P_q a$	$\delta(q, a)$
			$6-1=5$	ababa $\supset$ ababacab ( <b>F</b> )	$\delta(7, b)=2$
			$5-1=4$	abab $\supset$ ababacab ( <b>F</b> )	
			$4-1=3$	aba $\supset$ ababacab ( <b>F</b> )	
			$3-1=2$	ab $\supset$ ababacab ( <b>T</b> )	
	c	$k = \min(7+1, 7+2)=8$	$8-1=7$	ababaca $\supset$ ababacac ( <b>F</b> )	
			$7-1=6$	ababac $\supset$ ababacac ( <b>F</b> )	
			$6-1=5$	ababa $\supset$ ababacac ( <b>F</b> )	
			$5-1=4$	abab $\supset$ ababacac ( <b>F</b> )	
			$4-1=3$	aba $\supset$ ababacac ( <b>F</b> )	

# Finite Automata

$m = 7$ , so the outer loop will execute  
(0 – 7) 8 times.

$$\mathbf{P}_x = \mathbf{P}[1 \dots x]$$

	1	2	3	4	5	6	7
$\mathbf{P}$	a	b	a	b	a	c	a

q	a	$k = \min(m+1, q+2)$	$k = k-1$	$\mathbf{P}_k \supset \mathbf{P}_q a$	$\delta(q, a)$
			$3-1=2$	$ab \supset ababacac$ ( <b>F</b> )	
			$2-1=1$	$a \supset ababacac$ ( <b>F</b> )	
			$1-1=0$	$\epsilon \supset ababacac$ ( <b>T</b> )	$\delta(7, c)=0$

# Finite Automata

This matrix is showing the results of execution of transition function. These results will be used to execute the Finite-Automaton-Matcher to find any valid shifts.

e.g.  $\delta(3, b)=4$ ,  $\delta(5, c)=6$   
etc.

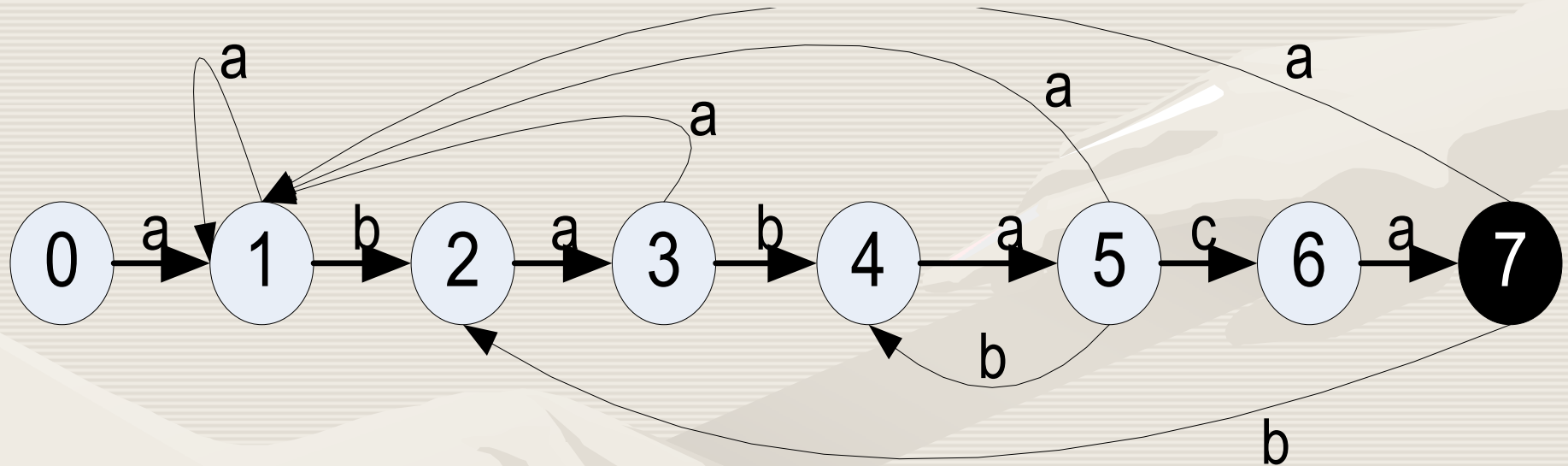
state	input			P
	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	



# Finite Automata

## State Transition Diagram

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 is the only accepting state. A directed edge from state  $i$  to state  $j$  labeled  $a$  represents  $\delta(i, a) = j$ . Some edges corresponding to failing matches are not shown; by convention, if a state  $i$  has no outgoing edge labeled  $a$  for some  $a \in \Sigma$ , then  $\delta(i, a) = 0$ .

# Finite Automata

## Executing Finite-Automaton-Matcher

- $T = \text{abababacaba}$  (Text)
- $n = 11$  (Length of Text)
- $P = \text{ababaca}$  (Pattern)
- $m = 7$  (Length of Pattern)
- $\Sigma = \{a, b, c\}$  (Character Set)
- $q = 0$  (Start State initially 0)

# Finite Automata

$n = 11$ , so loop will  
execute  $(1 - 11)$  11 times

	1	2	3	4	5	6	7	8	9	10	11
T	a	b	a	b	a	b	a	c	a	b	a

i	$q = \delta(q, T[i])$	$q = m$	Output
1	$q = \delta(0, a) = 1$	$1 = 7$ ( <b>False</b> )	
2	$q = \delta(1, b) = 2$	$2 = 7$ ( <b>False</b> )	
3	$q = \delta(2, a) = 3$	$3 = 7$ ( <b>False</b> )	
4	$q = \delta(3, b) = 4$	$4 = 7$ ( <b>False</b> )	
5	$q = \delta(4, a) = 5$	$5 = 7$ ( <b>False</b> )	
6	$q = \delta(5, b) = 4$	$4 = 7$ ( <b>False</b> )	
7	$q = \delta(4, a) = 5$	$5 = 7$ ( <b>False</b> )	
8	$q = \delta(5, c) = 6$	$6 = 7$ ( <b>False</b> )	

# Finite Automata

$n = 11$ , so loop will  
execute  $(1 - 11)$  11 times

	1	2	3	4	5	6	7	8	9	10	11
<b>T</b>	a	b	a	b	a	b	a	c	a	b	a

<b>i</b>	<b><math>q = \delta(q, T[i])</math></b>	<b><math>q = m</math></b>	<b>Output</b>
9	$q = \delta(6, a) = 7$	$7 = 7$ ( <b>True</b> )	Pattern occurs with shift $(9 - 7) = 2$
10	$q = \delta(7, b) = 2$	$2 = 7$ ( <b>False</b> )	
11	$q = \delta(2, a) = 3$	$3 = 7$ ( <b>False</b> )	

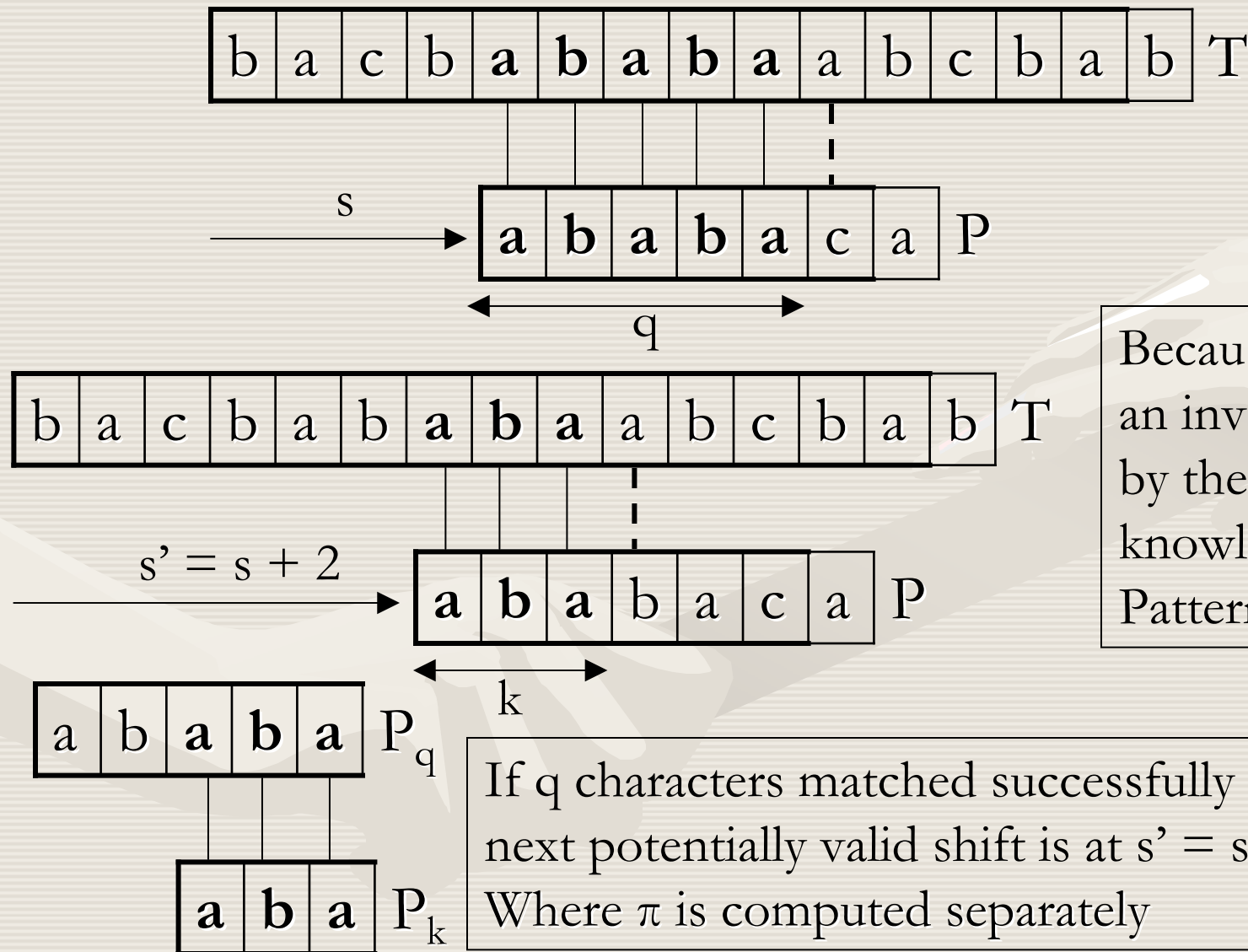
# Properties of Finite Automata Algorithm

- Pre-processing is required
- Total running time is  $O(m|\Sigma|) + \Theta(n)$
- Heavily dependent on Character Set
- Larger the Character Set more time it takes
- Performance decreases with the increase in the size of character set
- Un-necessarily stores all the valid and invalid state transitions

# Knuth-Morris-Pratt Algorithm

- This algorithm avoids the computation of the transition function  $\delta$  altogether, and its matching time is  $\Theta(n)$  using just an auxiliary (prefix) function  $\pi[1 \dots m]$  pre-computed from the pattern in time  $\Theta(m)$ .
- The prefix function  $\pi$  for a pattern encapsulates knowledge about how the pattern matches against shifts to itself.
- This information can be used to avoid testing useless shifts in the naïve pattern-matching algorithm or to avoid the pre-computation of  $\delta$  for a string-matching automaton.
- The information that  $q$  characters have matched successfully determines the corresponding text characters. Knowing these  $q$  text characters allows us to determine immediately that certain shifts are invalid.

# Knuth-Morris-Pratt Algorithm



Because  $s + 1$  is an invalid shift, by the knowledge of Pattern

If  $q$  characters matched successfully at shift  $s$ , the next potentially valid shift is at  $s' = s + (q - \pi[q])$ , Where  $\pi$  is computed separately

# Knuth-Morris-Pratt Algorithm

## COMPUTE-PREFIX-FUNCTION (P)

$m \leftarrow \text{length}[P]$

$\pi[1] \leftarrow 0$

$k \leftarrow 0$

**for**  $q \leftarrow 2$  **to**  $m$  **do**

**while**  $k > 0$  **and**  $P[k + 1] \neq P[q]$  **do**

$k \leftarrow \pi[k]$

**if**  $P[k + 1] = P[q]$  **then**

$k \leftarrow k + 1$

$\pi[q] \leftarrow k$

**return**  $\pi$



# Knuth-Morris-Pratt Algorithm

## KMP-MATCHER (T, P)

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

$m \leftarrow \text{length}[P]$

$\pi \leftarrow \text{COMPUTE-PREFIX-FUNCTION}(P)$

$q \leftarrow 0$  ▷ Number of characters matched

**for**  $i \leftarrow 1$  **to**  $n$  **do** ▷ Scan text from left to right

**while**  $q > 0$  **and**  $P[q + 1] \neq T[i]$  **do**

$q \leftarrow \pi[q]$  ▷ Next character does not match

**if**  $P[q + 1] = T[i]$  **then**

$q \leftarrow q + 1$  ▷ Next character matches

**if**  $q = m$  **then** ▷ Is all of P matched?

        print “Pattern occurs with shift”  $i - m$

$q \leftarrow \pi[q]$  ▷ Look for the next match

# Knuth-Morris-Pratt Algorithm

- Computing the Prefix Function

- $P = \text{ababaca}$  (Pattern)
- $m = 7$  (Length of Pattern)
- $\pi[1] = 0$  (Initial state is 0)
- $k = 0$
- First, computing the prefix function ( $\pi$ ) which gives us the knowledge about how the pattern matches against shifts with itself. This information can be used to avoid testing useless shifts.

# Knuth-Morris-Pratt Algorithm

## Computing the Prefix Function

$m = 7$ , so the outer loop will execute

$(2 - 7)$  6 times. Initially  $\pi[1] = 0, k = 0$

	1	2	3	4	5	6	7
<b>P</b>	a	b	a	b	a	c	a

q	$k > 0$ AND $P[k+1] \neq P[q]$	$k = \pi[k]$	if $P[k+1] = P[q]$	$k = k + 1$	$\pi[q] = k$
2	$0 > 0$ AND $a \neq b$ ( <b>False</b> )		$a = b$ ( <b>False</b> )		$\pi[2] = 0$
3	$0 > 0$ AND $a \neq a$ ( <b>False</b> )		$a = a$ ( <b>True</b> )	$0 + 1 = 1$	$\pi[3] = 1$
4	$1 > 0$ AND $b \neq b$ ( <b>False</b> )		$b = b$ ( <b>True</b> )	$1 + 1 = 2$	$\pi[4] = 2$
5	$2 > 0$ AND $a \neq a$ ( <b>False</b> )		$a = a$ ( <b>True</b> )	$2 + 1 = 3$	$\pi[5] = 3$
6	$3 > 0$ AND $b \neq c$ ( <b>True</b> )	1			
	$1 > 0$ AND $b \neq c$ ( <b>True</b> )	0			
	$0 > 0$ AND $a \neq c$ ( <b>False</b> )		$a = c$ ( <b>False</b> )		$\pi[6] = 0$
7	$0 > 0$ AND $a \neq a$ ( <b>False</b> )		$a = a$ ( <b>True</b> )	$0 + 1 = 1$	$\pi[7] = 1$

# Knuth-Morris-Pratt Algorithm

## Computing the Prefix Function

This matrix is showing the results of execution of prefix function. These results will be used to execute the KMP-Matcher to find any valid shifts.

P

1	2	3	4	5	6	7
a	b	a	b	a	c	a

State Matrix

i	$\pi[i]$
1	0
2	0
3	1
4	2
5	3
6	0
7	1

# Knuth-Morris-Pratt Algorithm

## Executing KMP-Matcher

- $T = \text{abababacaba}$  (Text)
- $n = 11$  (Length of Text)
- $P = \text{ababaca}$  (Pattern)
- $m = 7$  (Length of Pattern)
- $\pi$  (Computed Prefix Function for P)
- $q = 0$  (Number of Characters Matched)

# Knuth-Morris-Pratt Algorithm

	1	2	3	4	5	6	7		1	2	3	4	5	6	7	8	9	10	11
<b>P</b>	a	b	a	b	a	c	a		a	b	a	b	a	b	a	c	a	b	a
n = 11, so outer loop will execute (1 – 11) 11 times									<b>T</b>										

i	q>0 And P[q+1]≠T[i]	q= π[q]	P[q+1]= T[i]	q=q+1	q=m	q=π[q]	Output
1	0>0 AND a≠a (F)		a=a (T)	0+1=1	1=7 (F)		
2	1>0 AND b≠b (F)		b=b (T)	1+1=2	2=7 (F)		
3	2>0 AND a≠a (F)		a=a (T)	2+1=3	3=7 (F)		
4	3>0 AND b≠b (F)		b=b (T)	3+1=4	4=7 (F)		
5	4>0 AND a≠a (F)		a=a (T)	4+1=5	5=7 (F)		
6	5>0 AND c≠b (T)	3					
	3>0 AND b≠b (F)		b=b (T)	3+1=4	4=7 (F)		

# Knuth-Morris-Pratt Algorithm

	1	2	3	4	5	6	7		1	2	3	4	5	6	7	8	9	10	11
<b>P</b>	a	b	a	b	a	c	a		a	b	a	b	a	b	a	c	a	b	a

n = 11, so outer loop will execute (1 – 11) 11 times

i	q>0 And P[q+1]≠T[i]	q= π[q]	P[q+1]= T[i]	q=q+1	q=m	q= π[q]	Output
7	4>0 AND a≠a (F)		a=a (T)	4+1=5	5=7 (F)		
8	5>0 AND c≠c (F)		c=c (T)	5+1=6	6=7 (F)		
9	6>0 AND a≠a (F)		a=a (T)	6+1=7	7=7 (T)	1	Pattern Occurs with Shift (9-7) <b>2</b>
10	1>0 AND b≠b (F)		b=b (T)	1+1=2	2=7 (F)		
11	2>0 AND a≠a (F)		a=a (T)	2+1=3	3=7 (F)		

# Properties of Knuth-Morris-Pratt Algorithm

- Pre-processing is required
- Total running time is  $\Theta(m) + \Theta(n)$
- No dependency on Character Set
- Linear execution time
- Efficiently use the information about the Pattern and the characters read so far.
- Most efficient of the four algorithms.

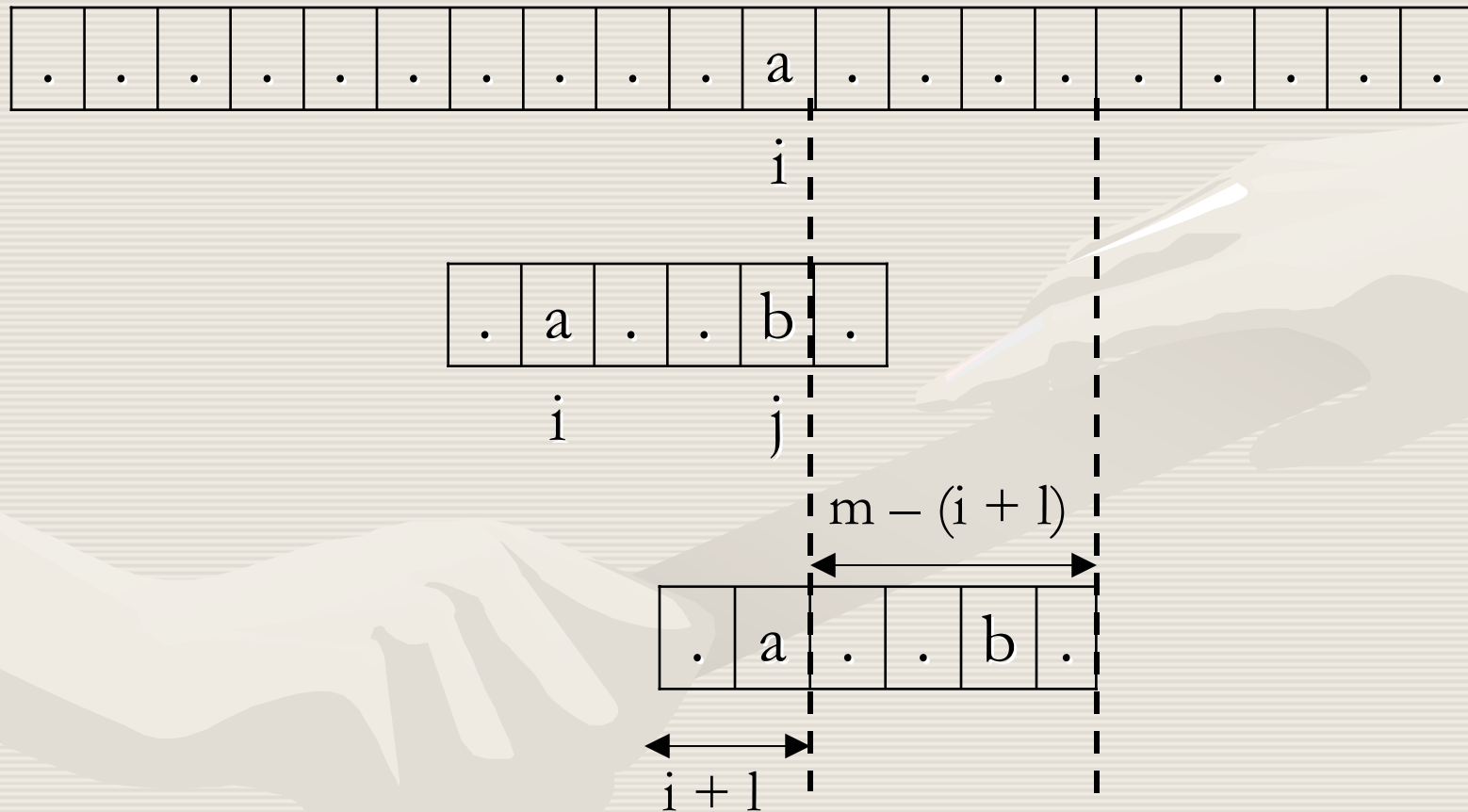


# Boyer Moore Algorithm

- It works fastest when alphabet is moderately sized and pattern is relatively long.
- Pattern is matched from right to left so in case of a mismatch we can skip a considerable number of characters.
- **Character Jump Heuristic**
  - During the testing of possible placement of  $P$  against  $T$ , a mismatch of text character  $T[i]=c$  with the corresponding pattern character  $P[j]$  is handled as follows.
  - If  $c$  is not contained anywhere in  $P$ , then shift  $P$  completely past  $T[i]$  (for it cannot match any character in  $P$ ). Otherwise shift  $P$  until an occurrence of character  $c$  in  $P$  gets aligned with  $T[i]$

# Boyer Moore Algorithm

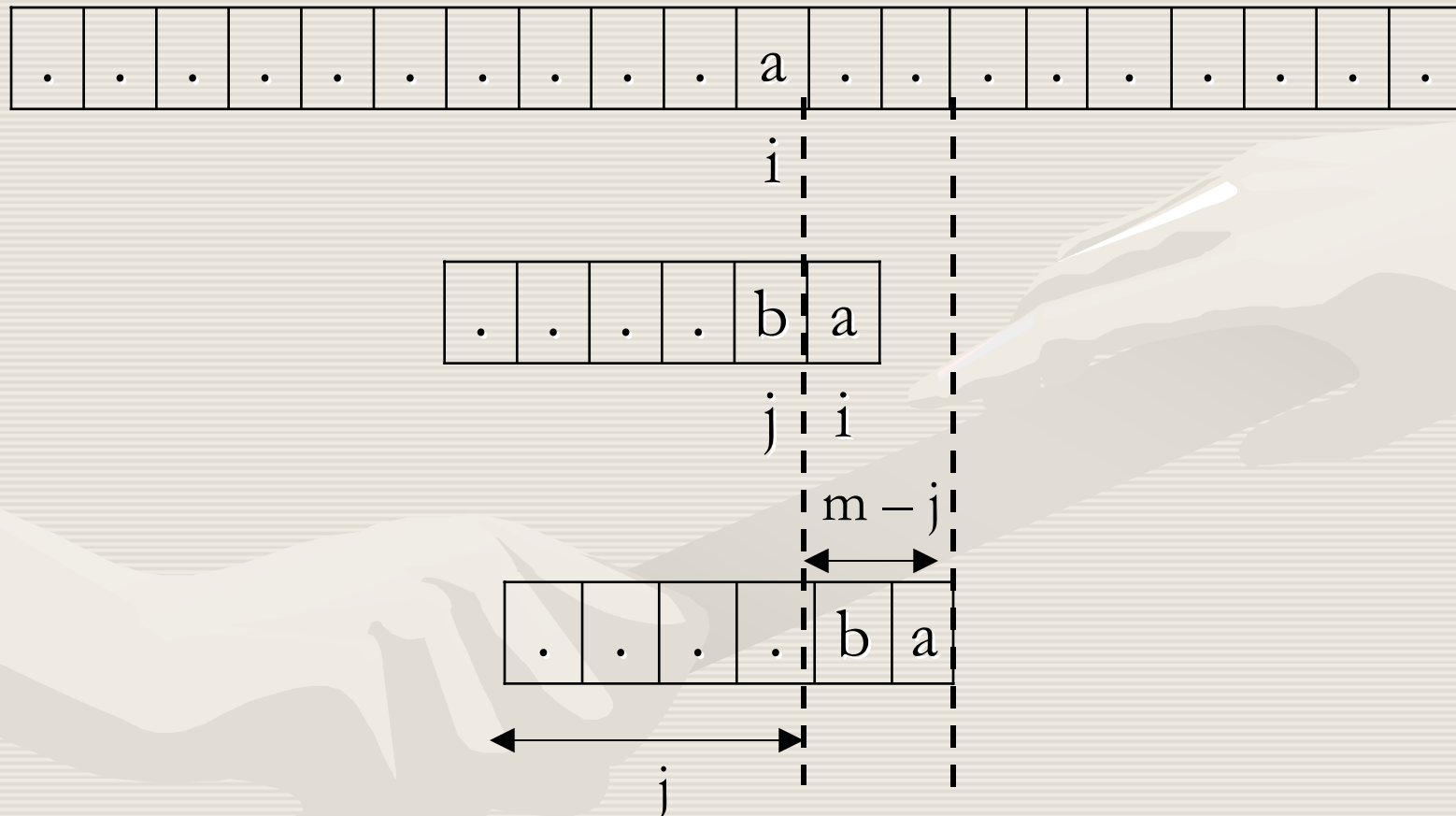
- Illustration of Jump in Boyer Moore Algorithm, where  $l$  denotes  $\text{Last}(P, T[i])$  (to be discussed shortly)



- Case 1: if  $i + 1 \leq j$ , then we shift the pattern by  $j - 1$  units**

# Boyer Moore Algorithm

- Illustration of Jump in Boyer Moore Algorithm, where  $i$  denotes  $\text{Last}(P, T[i])$  (to be discussed shortly)



- Case 2: if  $j < i + 1$ , then we shift the pattern by one unit

# Boyer Moore Algorithm

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

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

Since  $i+1 \leq j$  so shift by  $j - 1$  characters

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

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

Since  $j < i+1$  so shift by 1 character

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

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

Since  $i+1 \leq j$  so shift by  $j - 1$  characters

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

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

No match so shift by  $m$  characters

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

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

Since  $i+1 \leq j$  so shift by  $j-1$  characters

# Boyer Moore Algorithm

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

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

.....

# Boyer Moore Algorithm

## BOYER-MOORE-MATCHER (T, P)

$m \leftarrow \text{length}[P]$

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

$i \leftarrow m - 1$

$j \leftarrow m - 1$

**repeat**

**if**  $P[j] = T[i]$  **then**

**if**  $j = 0$  **then**

            print "Pattern Occurs with Shift"  $i$

$i \leftarrow i + m$

$j \leftarrow m - 1$

**else**

$i \leftarrow i - 1$

$j \leftarrow j - 1$

**else**

$i \leftarrow i + m - \min(j, 1 + \text{last}(P, T[i]))$

$j \leftarrow m - 1$

**until**  $i > n - 1$

# Boyer Moore Algorithm

**LAST(P, c)**

$i \longleftarrow \text{length}[P] - 1$

**do while**  $i \geq 0$

**if**  $P[i] = c$  **then**

**return**  $i$

$i \longleftarrow i - 1$

**return** -1

Since index start from 0 therefore

e.g. Last ( abxyabax, b) will return 5

Last ( abxyabax, a) will return 6

Last ( abxyabax, c) will return -1

- This function takes the Pattern P and a character c as input and return the index of the first character from the right side of pattern, if exists otherwise returns -1.

# Boyer Moore Algorithm

- Executing Boyer-Moore-Matcher
- $T = \text{abababacaba}$  (Text)
- $n = 11$  (Length of Text)
- $P = \text{ababaca}$  (Pattern)
- $m = 7$  (Length of Pattern)
- $i=6$  ( $m-1$ )
- $j=6$  ( $m-1$ )



# Boyer Moore Algorithm

	0	1	2	3	4	5	6		0	1	2	3	4	5	6	7	8	9	10
<b>P</b>	a	b	a	b	a	c	a	<b>T</b>	a	b	a	b	a	b	a	c	a	b	a

$n=11$  so loop will execute until  $i > 10$  ( $i > n-1$ ), Initially  $i=6$ , and  $j=6$  because  $i=j=m-1$

$P[j]=T[i]$	$j=0$	$i=i-1$	$j=j-1$	$i = i + m - \min(j, 1 + \text{last}(P, T[i]))$	$j = m-1$	Result
$a=a$ ( <b>T</b> )	F	$6-1=5$	$6-1=5$			
$c=b$ ( <b>F</b> )				$=5+7 - \min(5, 1 + \text{last}(P, b))$ $=12 - \min(5, 1+3)$ $=12 - 4$ $=8$	$7-1=6$	
$a=a$ ( <b>T</b> )	F	$8-1=7$	$6-1=5$			
$c=c$ ( <b>T</b> )	F	$7-1=6$	$5-1=4$			
$a=a$ ( <b>T</b> )	F	$6-1=5$	$4-1=3$			57

# Boyer Moore Algorithm

	0	1	2	3	4	5	6		0	1	2	3	4	5	6	7	8	9	10
<b>P</b>	a	b	a	b	a	c	a	<b>T</b>	a	b	a	b	a	b	a	c	a	b	a

$n=11$  so loop will continue to execute while  $i \leq 10$  (i.e. until  $i > 10$  ( $i > n-1$ ))

$P[j] = T[i]$	$j=0$	$i=i-1$	$j=j-1$	$i = i + m - \min(j, 1 + \text{last}(P, T[i]))$	$j = m-1$	Result
$b=b$ ( <b>T</b> )	F	$5-1=4$	$3-1=2$			
$a=a$ ( <b>T</b> )	F	$4-1=3$	$2-1=1$			
$b=b$ ( <b>T</b> )	F	$3-1=2$	$1-1=0$			
$a=a$ ( <b>T</b> )	T	<p>“Pattern occurs with shift” 2 and due to match <math>i</math> and <math>j</math> will be</p> <p><math>i = i + m = 2 + 7 = 9</math></p> <p><math>j = m - 1 = 7 - 1 = 6</math></p>				
$a=b$ ( <b>F</b> )				$= 9 + 7 - \min(6, 1 + \text{last}(P, b))$ $= 16 - \min(6, 1 + 3)$ $= 16 - 4 = 12$	$7-1=6$	

# Boyer Moore Algorithm

	0	1	2	3	4	5	6		0	1	2	3	4	5	6	7	8	9	10
<b>P</b>	a	b	a	b	a	c	a	<b>T</b>	a	b	a	b	a	b	a	c	a	b	a

$n=11$  so loop will continue to execute while  $i \leq 10$  (i.e. until  $i > 10$  ( $i > n-1$ ))

<b>P[j]= T[i]</b>	<b>j=0</b>	<b>i=i-1</b>	<b>j=j-1</b>	<b>i = i + m - min (j, 1 + last(P, T[i]) )</b>	<b>j= m-1</b>	<b>Result</b>
Loop will terminate since i becomes 12 ( $>10$ ) and because there are no more characters left in the string which can be completely matched with the pattern						

# Properties of Boyer-Moore Algorithm

- No Pre-processing is required
- Total running time is  $O(nm + |\Sigma|)$  (worst case)
- No track of previously read characters
- Scans the text from right to left
- No character set dependence
- Makes use of the pattern and index-character for Character Jumps (Character-Jump Heuristic)
- Good performance since a large portion of text is skipped

# Algorithm Time Comparison

Algorithm	Pre-Processing Time	Matching Time	Total Running Time
Knuth-Morris-Pratt	$\Theta(m)$	$\Theta(n)$	$\Theta(m) + \Theta(n)$
Finite Automaton	$O(m  \Sigma )$	$\Theta(n)$	$O(m  \Sigma ) + \Theta(n)$
Naïve	0	$O((n-m+1)m)$	$O((n-m+1)m)$
Boyer-Moore	0	$O(nm +  \Sigma )$	$O(nm +  \Sigma )$