

# String Matching Algorithms

Antonio Carzaniga

Faculty of Informatics  
Università della Svizzera italiana

December 22, 2011

- Problem definition
- Naïve algorithm
- Knuth-Morris-Pratt algorithm
- Boyer-Moore algorithm



- Given the text  
*Nel mezzo del cammin di nostra vita  
mi ritrovai per una selva oscura  
che la dritta via era smarrita...*

Find the string “trova”

- Given the text  
*Nel mezzo del cammin di nostra vita  
mi ritrovai per una selva oscura  
che la dritta via era smarrita...*

Find the string “trova”

- A more challenging example: How many times does the string “110011” appear in the following text

```
0011110101011010011000110101111011010111  
0110111001001010101011111011110110000101  
1011000010111111011110011000011111000100  
1001010010111011101011011110101001100101  
0010111001000011111110010011011101011010  
0110011011101001010010101000010100111110
```

- Given the text  
*Nel mezzo del cammin di nostra vita  
mi ritrovai per una selva oscura  
che la dritta via era smarrita...*

Find the string “trova”

- A more challenging example: How many times does the string “110011” appear in the following text

```
0011110101011010011000110101111011010111  
0110111001001010101011111011110110000101  
1011000010111111011110011000011111000100  
1001010010111011101011011110101001100101  
0010111001000011111110010011011101011010  
0110011011101001010010101000010100111110
```

# String Matching: Definitions

- Given a **text**  $T$ 
  - ▶  $T \in \Sigma^*$ : finite alphabet  $\Sigma$
  - ▶  $|T| = n$ : the length of  $T$  is  $n$

# String Matching: Definitions

- Given a **text**  $T$ 
  - ▶  $T \in \Sigma^*$ : finite alphabet  $\Sigma$
  - ▶  $|T| = n$ : the length of  $T$  is  $n$
  
- Given a **pattern**  $P$ 
  - ▶  $P \in \Sigma^*$ : same finite alphabet  $\Sigma$
  - ▶  $|P| = m$ : the length of  $P$  is  $m$

# String Matching: Definitions

- Given a **text**  $T$ 
  - ▶  $T \in \Sigma^*$ : finite alphabet  $\Sigma$
  - ▶  $|T| = n$ : the length of  $T$  is  $n$
  
- Given a **pattern**  $P$ 
  - ▶  $P \in \Sigma^*$ : same finite alphabet  $\Sigma$
  - ▶  $|P| = m$ : the length of  $P$  is  $m$
  
- Both  $T$  and  $P$  can be modeled as arrays
  - ▶  $T[1 \dots n]$  and  $P[1 \dots m]$

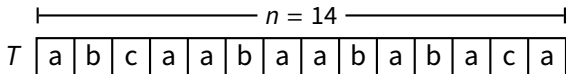
# String Matching: Definitions

- Given a **text**  $T$ 
  - ▶  $T \in \Sigma^*$ : finite alphabet  $\Sigma$
  - ▶  $|T| = n$ : the length of  $T$  is  $n$
  
- Given a **pattern**  $P$ 
  - ▶  $P \in \Sigma^*$ : same finite alphabet  $\Sigma$
  - ▶  $|P| = m$ : the length of  $P$  is  $m$
  
- Both  $T$  and  $P$  can be modeled as arrays
  - ▶  $T[1 \dots n]$  and  $P[1 \dots m]$
  
- Pattern  $P$  occurs with **shift**  $s$  in  $T$  iff
  - ▶  $0 \leq s \leq n - m$
  - ▶  $T[s + i] = P[i]$  for all positions  $1 \leq i \leq m$

- Problem: find all  $s$  such that
  - ▶  $0 \leq s \leq n - m$
  - ▶  $T[s + i] = P[i]$  for  $1 \leq i \leq m$

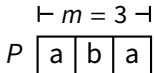
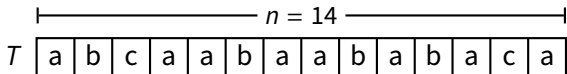
■ Problem: find all  $s$  such that

- ▶  $0 \leq s \leq n - m$
- ▶  $T[s + i] = P[i]$  for  $1 \leq i \leq m$



■ Problem: find all  $s$  such that

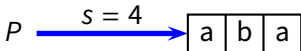
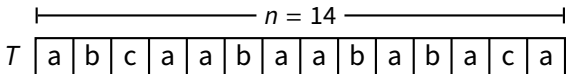
- ▶  $0 \leq s \leq n - m$
- ▶  $T[s + i] = P[i]$  for  $1 \leq i \leq m$



■ Result

■ Problem: find all  $s$  such that

- ▶  $0 \leq s \leq n - m$
- ▶  $T[s + i] = P[i]$  for  $1 \leq i \leq m$

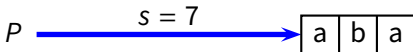
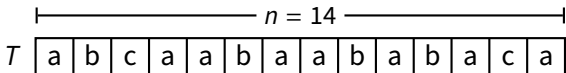


■ Result

$s = 4$

■ Problem: find all  $s$  such that

- ▶  $0 \leq s \leq n - m$
- ▶  $T[s + i] = P[i]$  for  $1 \leq i \leq m$



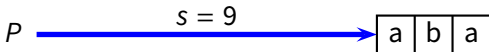
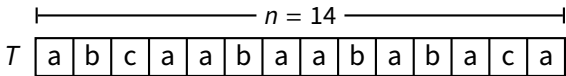
■ Result

$$s = 4$$

$$s = 7$$

■ Problem: find all  $s$  such that

- ▶  $0 \leq s \leq n - m$
- ▶  $T[s + i] = P[i]$  for  $1 \leq i \leq m$



■ Result

$$s = 4$$

$$s = 7$$

$$s = 9$$

# Naïve Algorithm

## Naïve Algorithm

- For each position  $s$  in  $0 \dots n - m$ , see if  $T[s + i] = P[i]$  for all  $1 \leq i \leq m$

- For each position  $s$  in  $0 \dots n - m$ , see if  $T[s + i] = P[i]$  for all  $1 \leq i \leq m$

```
NAIVE-STRING-MATCHING( $T, P$ )  
1   $n = \text{length}(T)$   
2   $m = \text{length}(P)$   
3  for  $s = 0$  to  $n - m$   
4      if SUBSTRING-AT( $T, P, s$ )  
5          OUTPUT( $s$ )
```

# Naïve Algorithm

- For each position  $s$  in  $0 \dots n - m$ , see if  $T[s + i] = P[i]$  for all  $1 \leq i \leq m$

## **NAIVE-STRING-MATCHING**( $T, P$ )

```
1  $n = \text{length}(T)$ 
2  $m = \text{length}(P)$ 
3 for  $s = 0$  to  $n - m$ 
4     if SUBSTRING-AT( $T, P, s$ )
5         OUTPUT( $s$ )
```

## **SUBSTRING-AT**( $T, P, s$ )

```
1 for  $i = 1$  to  $\text{length}(P)$ 
2     if  $T[s + i] \neq P[i]$ 
3         return FALSE
4 return TRUE
```

# Complexity of the Naïve Algorithm

# Complexity of the Naïve Algorithm

- Complexity of **NAIVE-STRING-MATCH** is  $O((n - m + 1)m)$

# Complexity of the Naïve Algorithm

- Complexity of **NAIVE-STRING-MATCH** is  $O((n - m + 1)m)$
- Worst case example

$$T = a^n, \quad P = a^m$$

i.e.,

$$T = \overbrace{aa \cdots a}^n, \quad P = \overbrace{aa \cdots a}^m$$

# Complexity of the Naïve Algorithm

- Complexity of **NAIVE-STRING-MATCH** is  $O((n - m + 1)m)$
- Worst case example

$$T = a^n, \quad P = a^m$$

i.e.,

$$T = \overbrace{aa \cdots a}^n, \quad P = \overbrace{aa \cdots a}^m$$

So,  $(n - m + 1)m$  is a tight bound, so the (worst-case) complexity of **NAIVE-STRING-MATCH** is

$$\Theta((n - m + 1)m)$$

# Improvement Strategy

- Observation

$T$ 

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

$P$ 

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

## ■ Observation

$T$ 

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

=

$P$ 

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

## ■ Observation

$T$ 

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

= =

$P$ 

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

## ■ Observation

$T$ 

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

= =  $\neq$

$P$ 

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

- Observation

$T$ 

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

= =  $\neq$

$P$ 

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

- What now?

# Improvement Strategy

## ■ Observation

$T$ 

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

= = ≠

$P$ 

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

## ■ What now?

- ▶ the naïve algorithm ***goes back to the second position in  $T$  and starts from the beginning of  $P$***

# Improvement Strategy

## ■ Observation

$T$ 

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

= = ≠

$P$ 

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

## ■ What now?

- ▶ the naïve algorithm **goes back to the second position in  $T$**  and **starts from the beginning of  $P$**
- ▶ can't we simply move along through  $T$ ?

## ■ Observation

$T$ 

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

= = ≠

$P$ 

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

## ■ What now?

- ▶ the naïve algorithm ***goes back to the second position in  $T$  and starts from the beginning of  $P$***
- ▶ can't we simply move along through  $T$ ?
- ▶ why?

# Improvement Strategy (2)

## Improvement Strategy (2)

- Here's a wrong but insightful strategy

## Improvement Strategy (2)

- Here's a wrong but insightful strategy

### **WRONG-STRING-MATCHING**( $T, P$ )

```
1   $n = \text{length}(T)$ 
2   $m = \text{length}(P)$ 
3   $q = 0$            // number of characters matched in  $P$ 
4   $s = 1$ 
5  while  $s \leq n$ 
6       $s = s + 1$ 
7      if  $T[s] == P[q + 1]$ 
8           $q = q + 1$ 
9          if  $q == m$ 
10             OUTPUT( $s - m$ )
11              $q = 0$ 
12     else  $q = 0$ 
```

## Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

## Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

*T*

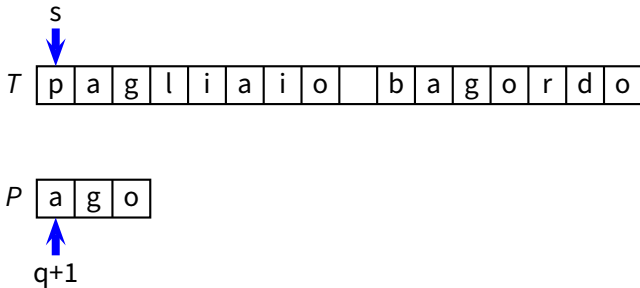
|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

*P*

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|

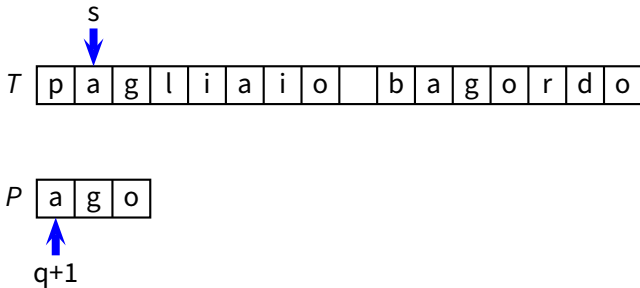
# Improvement Strategy (3)

## ■ Example run of **WRONG-STRING-MATCHING**



# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

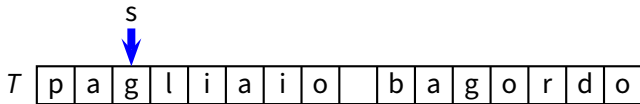


# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

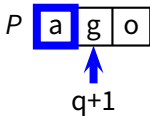
$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

A horizontal array of 16 cells representing the text T. The cells contain the characters 'p', 'a', 'g', 'l', 'i', 'a', 'i', 'o', a space, 'b', 'a', 'g', 'o', 'r', 'd', 'o'. A blue arrow labeled 's' points down to the third cell, which contains the character 'g'.

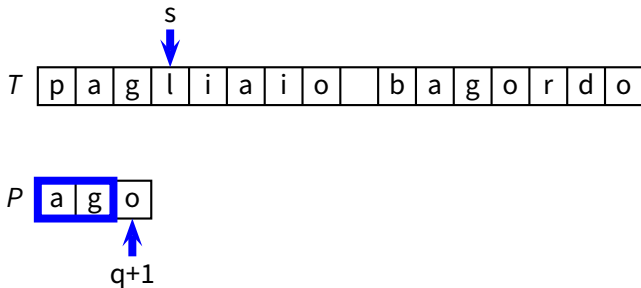
$P$ 

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|

A horizontal array of 3 cells representing the pattern P. The cells contain the characters 'a', 'g', 'o'. The first cell, containing 'a', is highlighted with a blue border. A blue arrow labeled 'q+1' points up to the second cell, which contains the character 'g'.

# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**




# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|



$P$ 

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|


  
 $q+1$

# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|



$P$ 

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|


  
 $q+1$

# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|



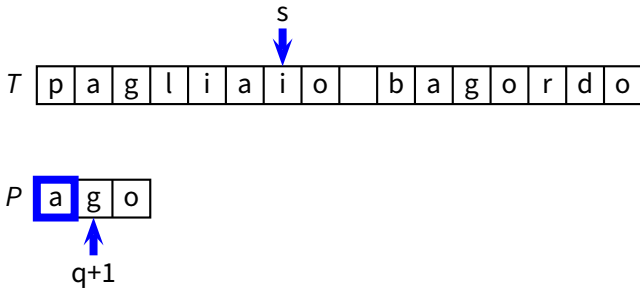
$P$ 

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|

  
 $q+1$

# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**




# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|



$P$ 

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|


  
 $q+1$

# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**


$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|



$P$ 


|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|



$q+1$

# Improvement Strategy (3)

## ■ Example run of **WRONG-STRING-MATCHING**

$T$  p a g l i a i o  s b a g o r d o

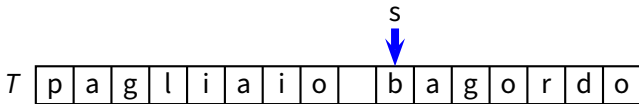
$P$  a g o

  
q+1

# Improvement Strategy (3)

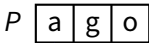
- Example run of **WRONG-STRING-MATCHING**

$T$  p a g l i a i o b a g o r d o



The diagram shows a horizontal array of 13 cells representing the string T. The characters are 'p', 'a', 'g', 'l', 'i', 'a', 'i', 'o', a space, 'b', 'a', 'g', 'o', 'r', 'd', 'o'. A blue arrow labeled 's' points down to the 'b' character, which is at the 10th position (index 9).

$P$  a g o



The diagram shows a horizontal array of 3 cells representing the string P. The characters are 'a', 'g', 'o'. A blue arrow labeled 'q+1' points up to the 'a' character, which is at the 0th position (index 0).


$q+1$

# Improvement Strategy (3)

## ■ Example run of **WRONG-STRING-MATCHING**

$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|



$P$ 

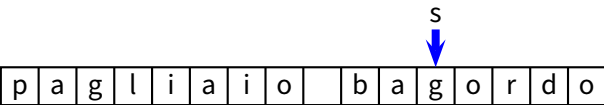
|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|

  
 $q+1$

# Improvement Strategy (3)

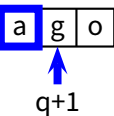
- Example run of **WRONG-STRING-MATCHING**

$T$  p a g l i a i o b a g o r d o



The diagram shows a horizontal array of 13 cells representing the text  $T$ . The cells contain the characters 'p', 'a', 'g', 'l', 'i', 'a', 'i', 'o', ' ', 'b', 'a', 'g', 'o', 'r', 'd', 'o'. A blue arrow labeled 's' points down to the third 'g' cell.

$P$  a g o

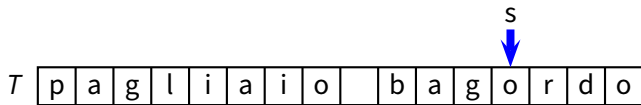


The diagram shows a horizontal array of 3 cells representing the pattern  $P$ . The cells contain the characters 'a', 'g', 'o'. A blue arrow labeled 'q+1' points up to the 'g' cell. The 'a' cell is highlighted with a blue border.

# Improvement Strategy (3)

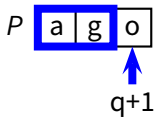
- Example run of **WRONG-STRING-MATCHING**

$T$  p a g l i a i o b a g o r d o



The diagram shows a horizontal array of 13 cells representing the text T. The characters are 'p', 'a', 'g', 'l', 'i', 'a', 'i', 'o', 'b', 'a', 'g', 'o', 'r', 'd', 'o'. A blue arrow labeled 's' points down to the 10th cell, which contains the character 'o'.

$P$  a g o



The diagram shows a horizontal array of 3 cells representing the pattern P. The characters are 'a', 'g', 'o'. A blue arrow labeled 'q+1' points up to the 3rd cell, which contains the character 'o'. The first two cells, 'a' and 'g', are enclosed in a blue rectangular box.

# Improvement Strategy (3)

## ■ Example run of **WRONG-STRING-MATCHING**

*T*

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

s



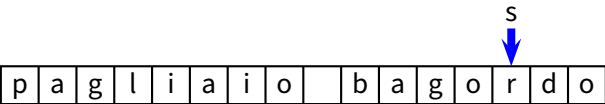
*P*

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|

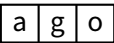
# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

$T$  p a g l i a i o b a g o r d o



$P$  a g o



$q+1$

Output: 10

# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

$s$   
↓

$P$ 

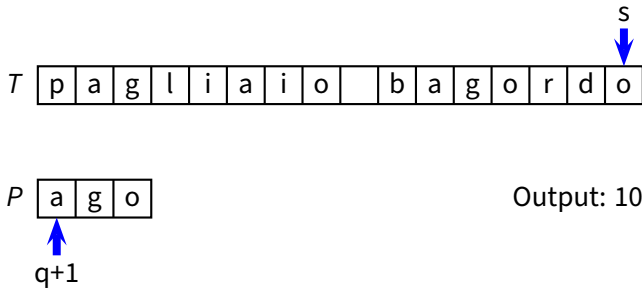
|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|

↑  
 $q+1$

Output: 10

# Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**



## Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

*T*

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

*P*

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|

Output: 10

- Done. Perfect!

## Improvement Strategy (3)

- Example run of **WRONG-STRING-MATCHING**

$T$ 

|   |   |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| p | a | g | l | i | a | i | o |  | b | a | g | o | r | d | o |
|---|---|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

$P$ 

|   |   |   |
|---|---|---|
| a | g | o |
|---|---|---|

Output: 10

- Done. Perfect!
- Complexity:  $\Theta(n)$

## Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?

## Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?

*T*

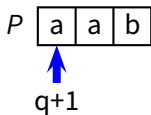
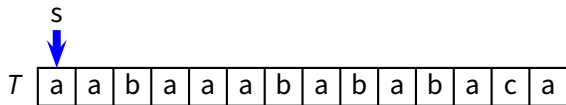
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | a | b | a | a | a | b | a | b | a | b | a | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

*P*

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

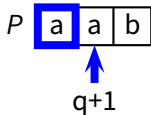
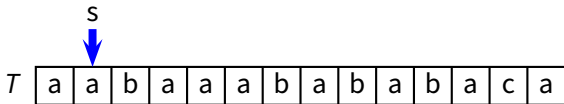
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



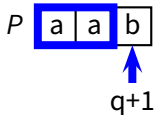
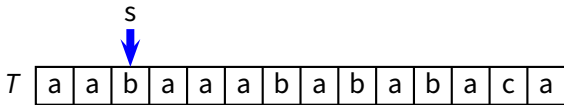
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



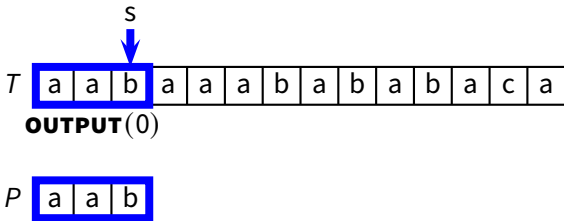
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



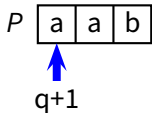
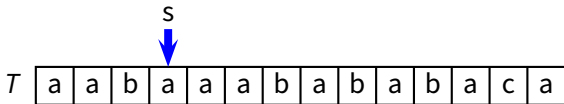
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



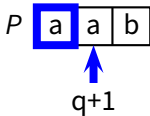
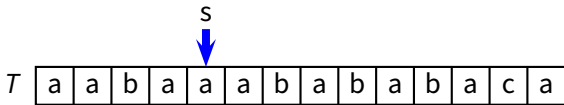
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



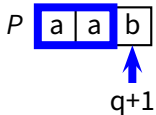
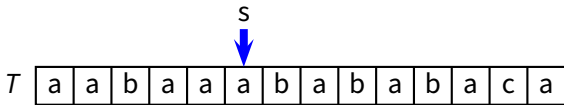
## Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



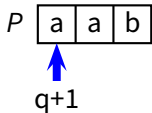
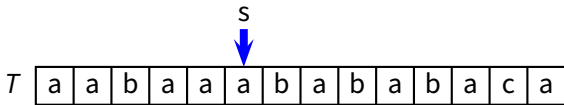
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



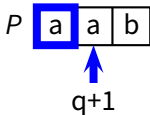
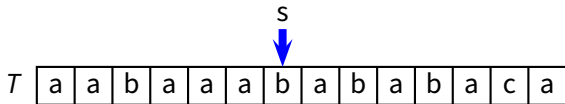
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



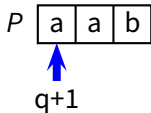
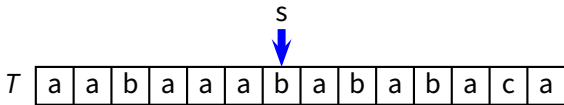
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



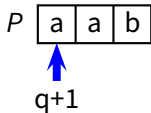
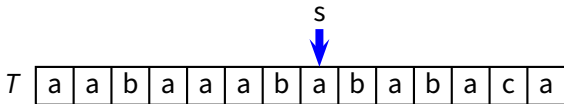
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



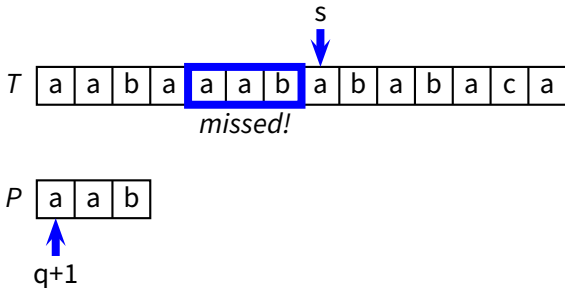
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



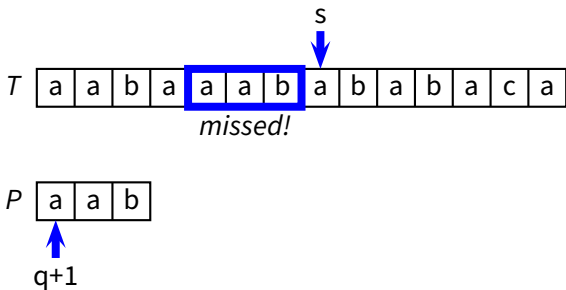
# Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



## Improvement Strategy (4)

- What is wrong with **WRONG-STRING-MATCHING**?



- So **WRONG-STRING-MATCHING** doesn't work, but it tells us something useful

## Improvement Strategy (5)

- Where did **WRONG-STRING-MATCHING** go wrong?

## Improvement Strategy (5)

- Where did **WRONG-STRING-MATCHING** go wrong?

*T*

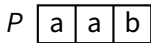
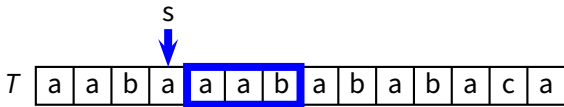
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | a | b | a | a | a | b | a | b | a | b | a | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

*P*

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

# Improvement Strategy (5)

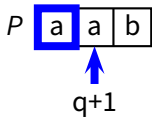
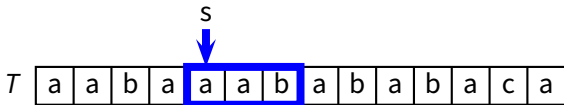
- Where did **WRONG-STRING-MATCHING** go wrong?



$q+1$

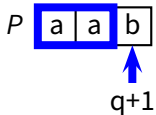
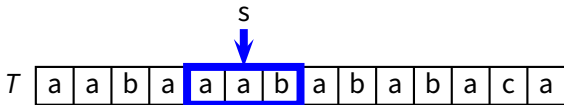
# Improvement Strategy (5)

- Where did **WRONG-STRING-MATCHING** go wrong?



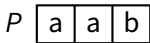
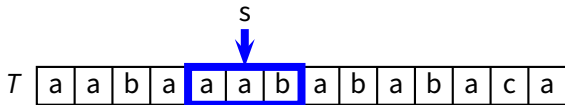
# Improvement Strategy (5)

- Where did **WRONG-STRING-MATCHING** go wrong?



# Improvement Strategy (5)

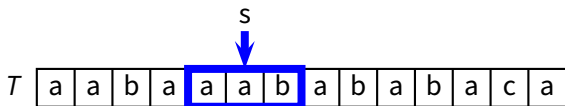
- Where did **WRONG-STRING-MATCHING** go wrong?



$q+1$

## Improvement Strategy (5)

- Where did **WRONG-STRING-MATCHING** go wrong?



- Wrong: by going all the way back to  $q = 0$  we throw away a good prefix of  $P$  that we already matched

## Improvement Strategy (6)

- Another example

# Improvement Strategy (6)

- Another example

$T$ 

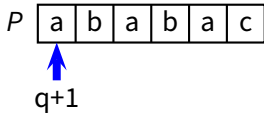
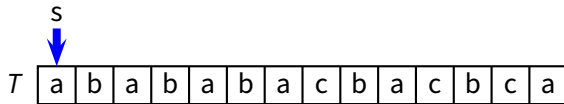
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | b | a | b | a | b | a | c | b | a | c | b | c | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$P$ 

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

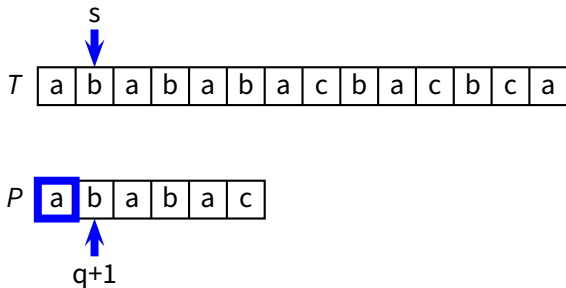
# Improvement Strategy (6)

- Another example



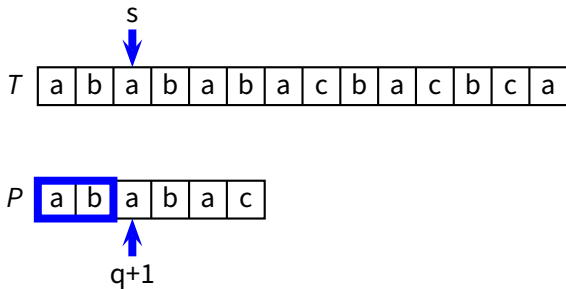
# Improvement Strategy (6)

- Another example



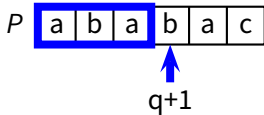
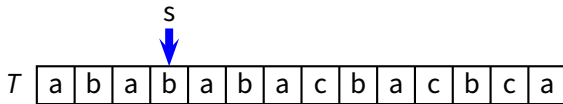
# Improvement Strategy (6)

- Another example



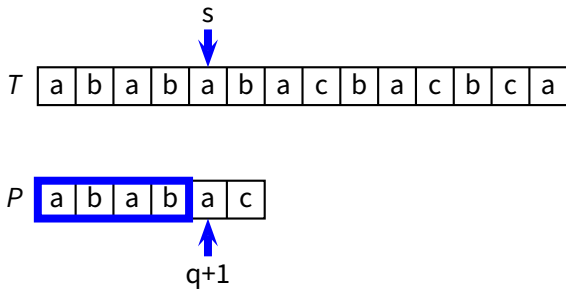
# Improvement Strategy (6)

- Another example



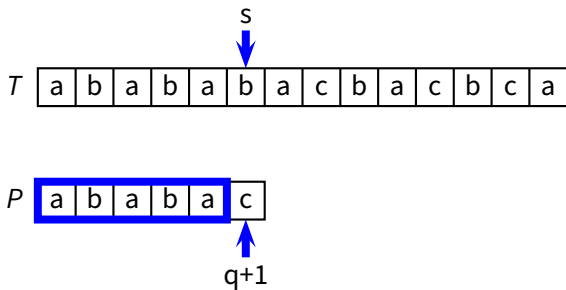
# Improvement Strategy (6)

- Another example



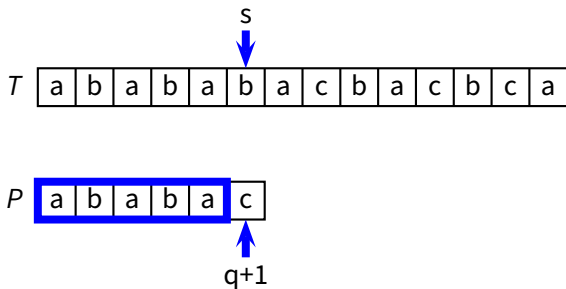
# Improvement Strategy (6)

- Another example



## Improvement Strategy (6)

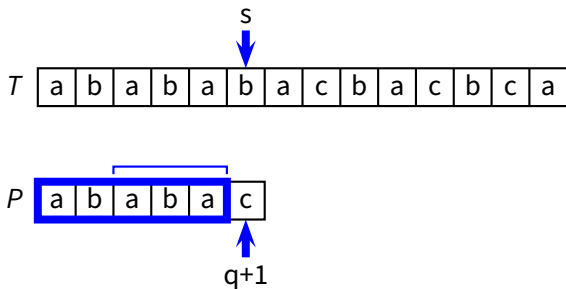
- Another example



- We have matched “ababa”

# Improvement Strategy (6)

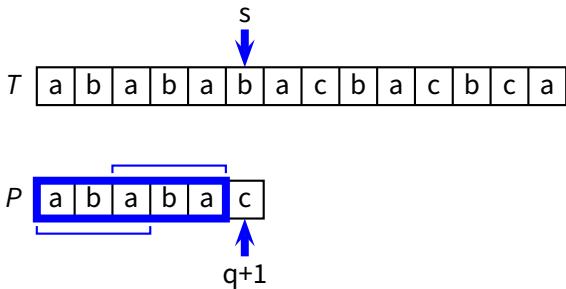
- Another example



- We have matched “ababa”
  - ▶ *suffix* “aba” can be *reused as a prefix*

# Improvement Strategy (6)

- Another example

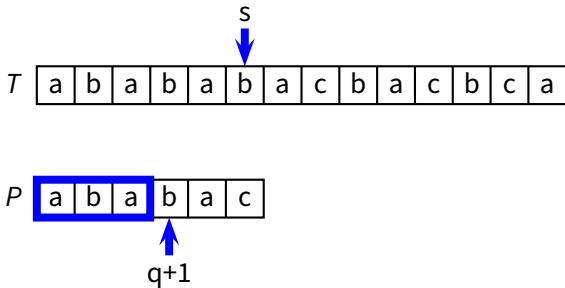


- We have matched “ababa”

- ▶ *suffix* “aba” can be *reused as a prefix*

# Improvement Strategy (6)

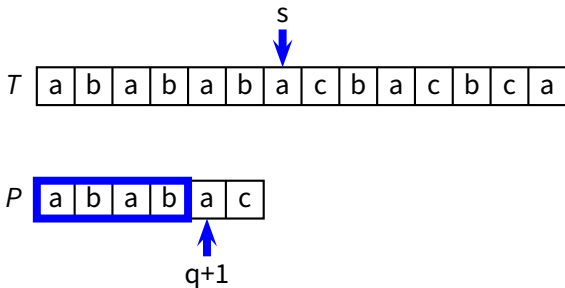
- Another example



- We have matched “ababa”
  - ▶ *suffix* “aba” can be *reused as a prefix*

# Improvement Strategy (6)

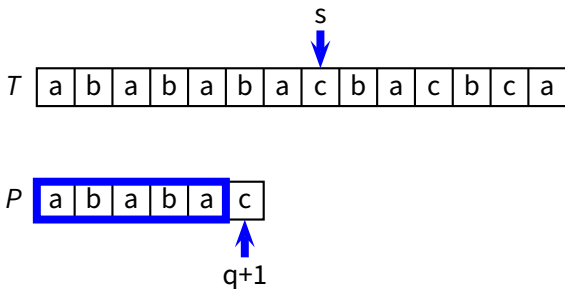
- Another example



- We have matched “ababa”
  - ▶ **suffix** “aba” can be **reused as a prefix**

## Improvement Strategy (6)

- Another example



- We have matched “ababa”
  - ▶ *suffix* “aba” can be *reused as a prefix*

# Improvement Strategy (6)

- Another example

$T$ 

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

  
**OUTPUT**(2)

$P$ 

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

- We have matched “ababa”
  - ▶ *suffix* “aba” can be *reused as a prefix*

- $P[1 \dots q]$  is the prefix of  $P$  matched so far

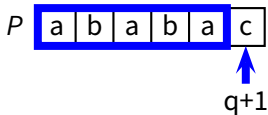
- $P[1 \dots q]$  is the prefix of  $P$  matched so far
- Find the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$ 
  - ▶ i.e., find  $0 \leq \pi < q$  such that  $P[q - \pi + 1 \dots q] = P[1 \dots \pi]$ 
    - ▶  $\pi = 0$  means that such a prefix does not exist

- $P[1 \dots q]$  is the prefix of  $P$  matched so far
- Find the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$ 
  - ▶ i.e., find  $0 \leq \pi < q$  such that  $P[q - \pi + 1 \dots q] = P[1 \dots \pi]$ 
    - ▶  $\pi = 0$  means that such a prefix does not exist

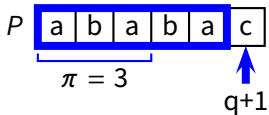
$P$ 

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

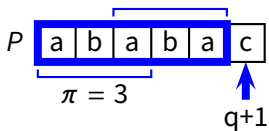
- $P[1 \dots q]$  is the prefix of  $P$  matched so far
- Find the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$ 
  - ▶ i.e., find  $0 \leq \pi < q$  such that  $P[q - \pi + 1 \dots q] = P[1 \dots \pi]$
  - ▶  $\pi = 0$  means that such a prefix does not exist



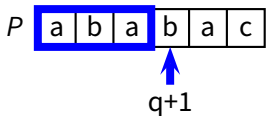
- $P[1 \dots q]$  is the prefix of  $P$  matched so far
- Find the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$ 
  - ▶ i.e., find  $0 \leq \pi < q$  such that  $P[q - \pi + 1 \dots q] = P[1 \dots \pi]$
  - ▶  $\pi = 0$  means that such a prefix does not exist



- $P[1 \dots q]$  is the prefix of  $P$  matched so far
- Find the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$ 
  - ▶ i.e., find  $0 \leq \pi < q$  such that  $P[q - \pi + 1 \dots q] = P[1 \dots \pi]$
  - ▶  $\pi = 0$  means that such a prefix does not exist

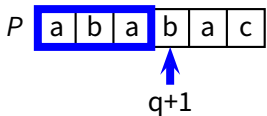


- $P[1 \dots q]$  is the prefix of  $P$  matched so far
- Find the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$ 
  - ▶ i.e., find  $0 \leq \pi < q$  such that  $P[q - \pi + 1 \dots q] = P[1 \dots \pi]$
  - ▶  $\pi = 0$  means that such a prefix does not exist



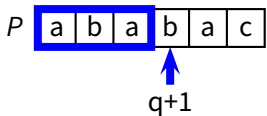
- Restart from  $q = \pi$

- $P[1 \dots q]$  is the prefix of  $P$  matched so far
- Find the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$ 
  - ▶ i.e., find  $0 \leq \pi < q$  such that  $P[q - \pi + 1 \dots q] = P[1 \dots \pi]$
  - ▶  $\pi = 0$  means that such a prefix does not exist



- Restart from  $q = \pi$
- Iterate as usual

- $P[1 \dots q]$  is the prefix of  $P$  matched so far
- Find the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$ 
  - ▶ i.e., find  $0 \leq \pi < q$  such that  $P[q - \pi + 1 \dots q] = P[1 \dots \pi]$
  - ▶  $\pi = 0$  means that such a prefix does not exist



- Restart from  $q = \pi$
- Iterate as usual
- In essence, this is the Knuth-Morris-Pratt algorithm

# The Prefix Function

- Given a pattern prefix  $P[1 \dots q]$ , the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$  depends only on  $P$  and  $q$

# The Prefix Function

- Given a pattern prefix  $P[1 \dots q]$ , the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$  depends only on  $P$  and  $q$
- This prefix is identified by its length  $\pi(q)$

# The Prefix Function

- Given a pattern prefix  $P[1 \dots q]$ , the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$  depends only on  $P$  and  $q$
- This prefix is identified by its length  $\pi(q)$
- Because  $\pi(q)$  depends only on  $P$  (and  $q$ ),  $\pi$  can be computed at the beginning by **PREFIX-FUNCTION**
  - ▶ we represent  $\pi$  as an array of length  $m$

# The Prefix Function

- Given a pattern prefix  $P[1 \dots q]$ , the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$  depends only on  $P$  and  $q$
- This prefix is identified by its length  $\pi(q)$
- Because  $\pi(q)$  depends only on  $P$  (and  $q$ ),  $\pi$  can be computed at the beginning by **PREFIX-FUNCTION**
  - ▶ we represent  $\pi$  as an array of length  $m$
- Example

$P$ 

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

# The Prefix Function

- Given a pattern prefix  $P[1 \dots q]$ , the longest prefix of  $P$  that is also a suffix of  $P[2 \dots q]$  depends only on  $P$  and  $q$
- This prefix is identified by its length  $\pi(q)$
- Because  $\pi(q)$  depends only on  $P$  (and  $q$ ),  $\pi$  can be computed at the beginning by **PREFIX-FUNCTION**
  - ▶ we represent  $\pi$  as an array of length  $m$
- Example

$P$ 

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

$\pi$ 

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 0 | 0 | 1 | 2 | 3 | 0 |
|---|---|---|---|---|---|

# The Knuth-Morris-Pratt Algorithm

## KMP-STRING-MATCHING( $T, P$ )

```
1   $n = \text{length}(T)$ 
2   $m = \text{length}(P)$ 
3   $\pi = \text{PREFIX-FUNCTION}(P)$ 
4   $q = 0$  // number of character matched
5  for  $i = 1$  to  $n$  // scan the text left-to-right
6      while  $q > 0$  and  $P[q + 1] \neq T[i]$ 
7           $q = \pi[q]$  // no match: go back using  $\pi$ 
8      if  $P[q + 1] == T[i]$ 
9           $q = q + 1$ 
10     if  $q == m$ 
11         OUTPUT( $i - m$ )
12          $q = \pi[q]$  // go back for the next match
```

# Prefix Function Algorithm

# Prefix Function Algorithm

- Computing the prefix function amounts to finding all the occurrences of a pattern  $P$  *in itself*

# Prefix Function Algorithm

- Computing the prefix function amounts to finding all the occurrences of a pattern  $P$  *in itself*
- In fact, **PREFIX-FUNCTION** is remarkably similar to **KMP-STRING-MATCHING**

## **PREFIX-FUNCTION**( $P$ )

```
1   $m = \text{length}(P)$ 
2   $\pi[1] = 0$ 
3   $k = 0$ 
4  for  $q = 2$  to  $m$ 
5      while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6           $k = \pi[k]$ 
7      if  $P[k + 1] == P[q]$ 
8           $k = k + 1$ 
9       $\pi[q] = k$ 
```

# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```

$P$ 

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

$\pi$ 

|  |  |  |  |  |  |
|--|--|--|--|--|--|
|  |  |  |  |  |  |
|--|--|--|--|--|--|

# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```

$P$ 

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

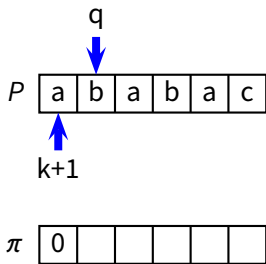
$\pi$ 

|   |  |  |  |  |  |
|---|--|--|--|--|--|
| 0 |  |  |  |  |  |
|---|--|--|--|--|--|

# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

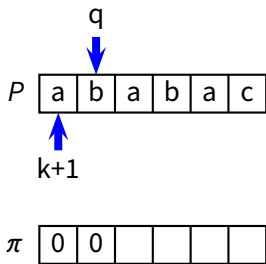
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

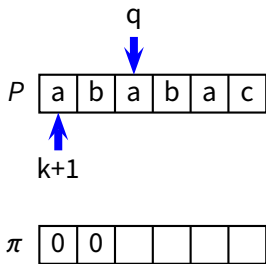
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

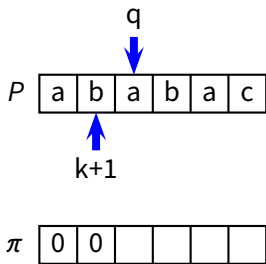
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

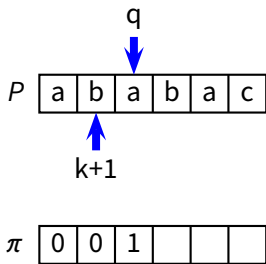
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

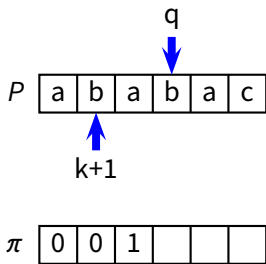
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

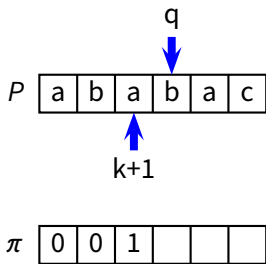
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

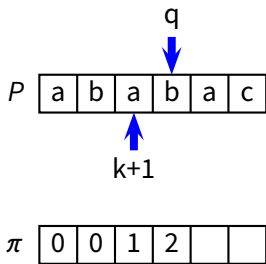
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

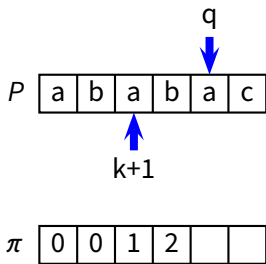
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

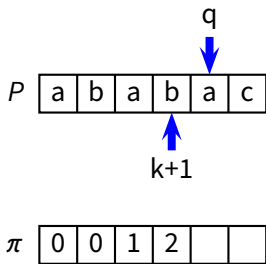
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

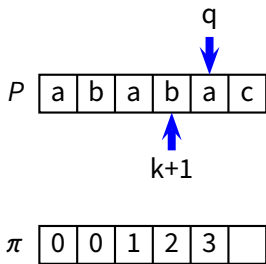
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

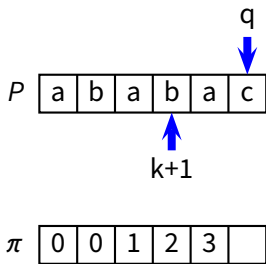
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

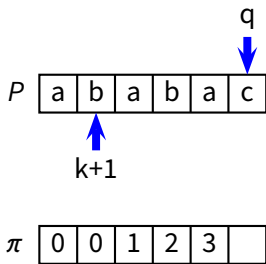
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

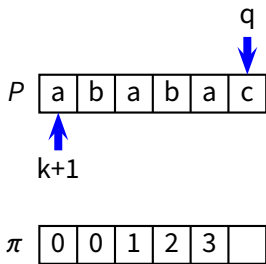
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

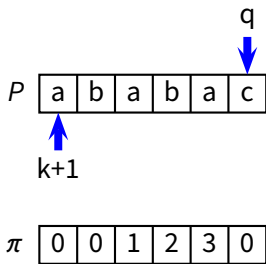
```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Prefix Function at Work

## PREFIX-FUNCTION( $P$ )

```
1  $m = \text{length}(P)$ 
2  $\pi[1] = 0$ 
3  $k = 0$ 
4 for  $q = 2$  to  $m$ 
5     while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
6          $k = \pi[k]$ 
7     if  $P[k + 1] == P[q]$ 
8          $k = k + 1$ 
9      $\pi[q] = k$ 
```



# Complexity of KMP

- $O(n)$  for the search phase

## Complexity of KMP

- $O(n)$  for the search phase
- $O(m)$  for the pre-processing of the pattern

# Complexity of KMP

- $O(n)$  for the search phase
- $O(m)$  for the pre-processing of the pattern
- The complexity analysis is non-trivial

# Complexity of KMP

- $O(n)$  for the search phase
- $O(m)$  for the pre-processing of the pattern
- The complexity analysis is non-trivial
- Can we do better?

- Knuth-Morris-Pratt is  $\Omega(n)$ 
  - ▶ KMP will *always* go through *at least*  $n$  character comparisons
  - ▶ it fixes our “wrong” algorithm in the case of *periodic* patterns and texts

- Knuth-Morris-Pratt is  $\Omega(n)$ 
  - ▶ KMP will *always* go through *at least*  $n$  character comparisons
  - ▶ it fixes our “wrong” algorithm in the case of *periodic* patterns and texts
- Perhaps there’s another algorithm that works better on the average case
  - ▶ e.g., in the absence of periodic patterns

## A New Strategy

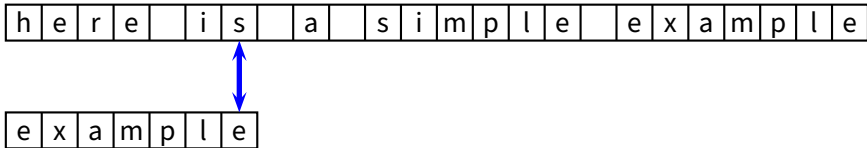
|   |   |   |   |  |   |   |  |   |  |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| h | e | r | e |  | i | s |  | a |  | s | i | m | p | l | e |  | e | x | a | m | p | l | e |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

## A New Strategy

|   |   |   |   |  |   |   |  |   |  |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| h | e | r | e |  | i | s |  | a |  | s | i | m | p | l | e |  | e | x | a | m | p | l | e |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| e | x | a | m | p | l | e |
|---|---|---|---|---|---|---|

## A New Strategy



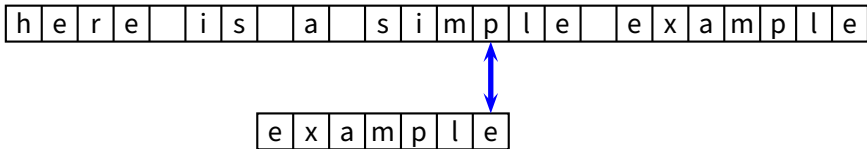
- We match the pattern right-to-left

|   |   |   |   |  |   |   |  |   |  |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| h | e | r | e |  | i | s |  | a |  | s | i | m | p | l | e |  | e | x | a | m | p | l | e |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| e | x | a | m | p | l | e |
|---|---|---|---|---|---|---|

- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern

## A New Strategy



- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern

# A New Strategy

h e r e   i s   a   s i m p l e   e x a m p l e

e x a m p l e

- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

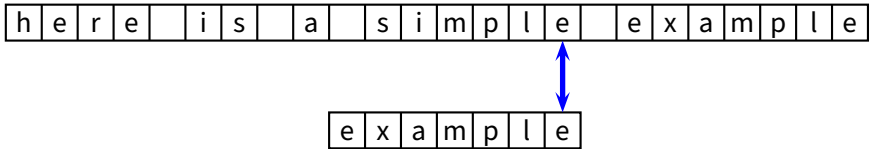
## A New Strategy

h e r e i s a s i m p l e e x a m p l e

e x a m p l e

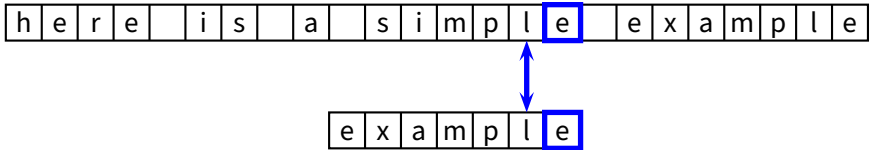
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

# A New Strategy



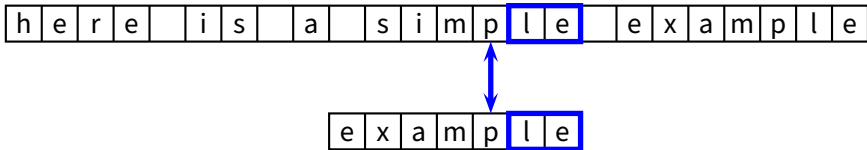
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

# A New Strategy



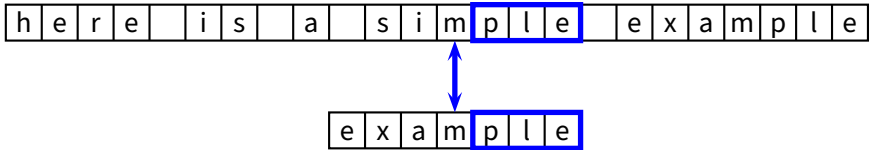
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

## A New Strategy



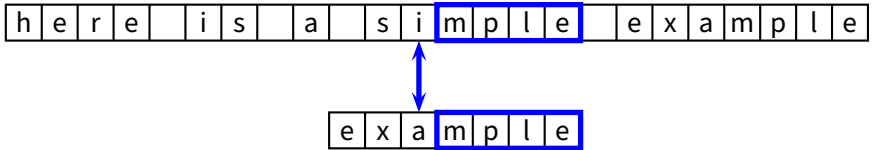
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

# A New Strategy



- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

# A New Strategy



- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

|   |   |   |   |  |   |   |  |   |  |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| h | e | r | e |  | i | s |  | a |  | s | i | m | p | l | e |  | e | x | a | m | p | l | e |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| e | x | a | m | p | l | e |
|---|---|---|---|---|---|---|

- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

|   |   |   |   |  |   |   |  |   |  |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| h | e | r | e |  | i | s |  | a |  | s | i | m | p | l | e |  | e | x | a | m | p | l | e |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| e | x | a | m | p | l | e |
|---|---|---|---|---|---|---|

- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

## A New Strategy

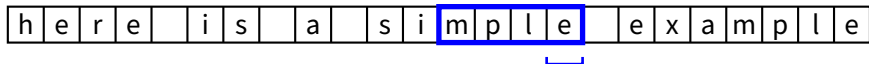
|   |   |   |   |  |   |   |  |   |  |   |   |   |   |   |   |  |   |   |   |   |   |   |   |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|
| h | e | r | e |  | i | s |  | a |  | s | i | m | p | l | e |  | e | x | a | m | p | l | e |
|---|---|---|---|--|---|---|--|---|--|---|---|---|---|---|---|--|---|---|---|---|---|---|---|

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| e | x | a | m | p | l | e |
|---|---|---|---|---|---|---|

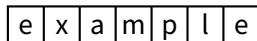
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

## A New Strategy

h e r e   i s   a   s i m p l e   e x a m p l e

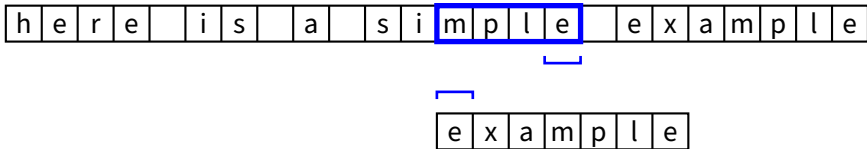


e x a m p l e



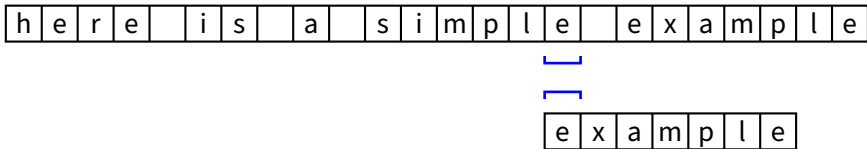
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

## A New Strategy



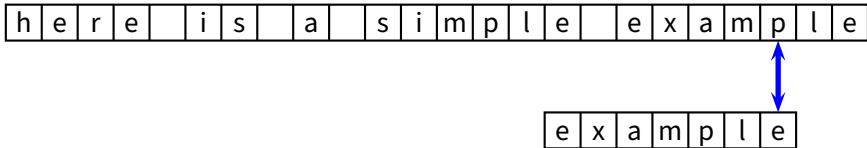
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$

# A New Strategy



- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

# A New Strategy



- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

# A New Strategy

h e r e   i s   a   s i m p l e   e x a m p l e

e x a m p l e

- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

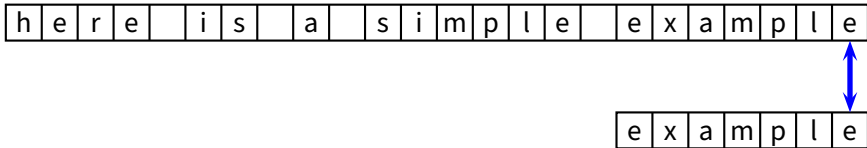
# A New Strategy

h e r e   i s   a   s i m p l e   e x a m p l e

e x a m p l e

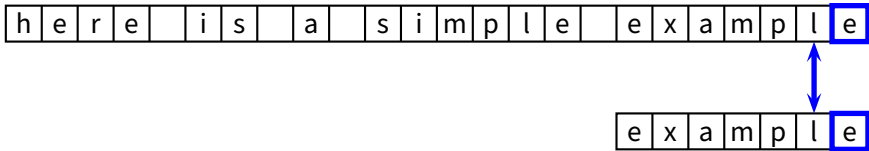
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

# A New Strategy



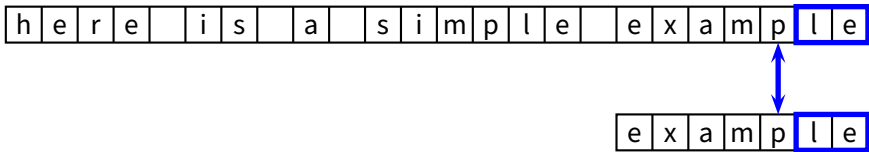
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

# A New Strategy



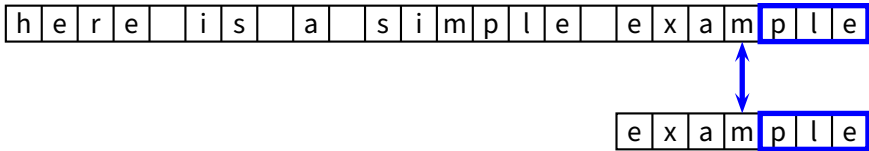
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

# A New Strategy



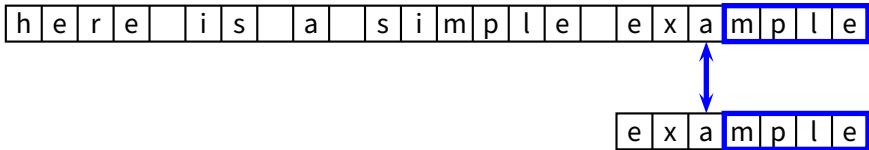
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

# A New Strategy



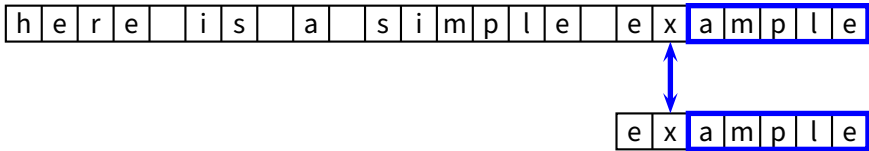
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

## A New Strategy



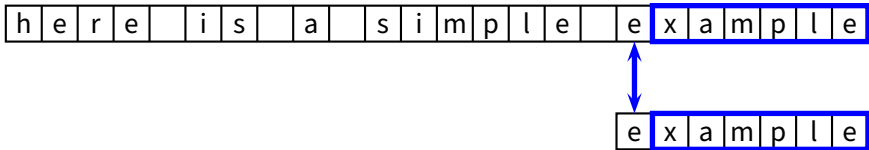
- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

## A New Strategy



- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

## A New Strategy



- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

# A New Strategy

h e r e   i s   a   s i m p l e   e x a m p l e

e x a m p l e

- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match

# A New Strategy

h e r e   i s   a   s i m p l e   e x a m p l e

e x a m p l e

- We match the pattern right-to-left
- If we find a bad character  $\alpha$  in the text, we can shift
  - ▶ so that the pattern skips  $\alpha$ , if  $\alpha$  is not in the pattern
  - ▶ so that the pattern lines up with the rightmost occurrence of  $\alpha$  in the pattern, if the pattern contains  $\alpha$
  - ▶ so that a pattern prefix lines up with a suffix of the current partial (or complete) match
- In essence, this is the Boyer-Moore algorithm

# Comments on Boyer-Moore

## Comments on Boyer-Moore

- Like KMP, Boyer-Moore includes a pre-processing phase

## Comments on Boyer-Moore

- Like KMP, Boyer-Moore includes a pre-processing phase
- The pre-processing is  $O(m)$

## Comments on Boyer-Moore

- Like KMP, Boyer-Moore includes a pre-processing phase
- The pre-processing is  $O(m)$
- The search phase is  $O(nm)$

## Comments on Boyer-Moore

- Like KMP, Boyer-Moore includes a pre-processing phase
- The pre-processing is  $O(m)$
- The search phase is  $O(nm)$
- The search phase can be as low as  $O(n/m)$  in common cases

## Comments on Boyer-Moore

- Like KMP, Boyer-Moore includes a pre-processing phase
- The pre-processing is  $O(m)$
- The search phase is  $O(nm)$
- The search phase can be as low as  $O(n/m)$  in common cases
- In practice, Boyer-Moore is the fastest string-matching algorithm for most applications