



ALBERT-LUDWIGS-  
UNIVERSITÄT FREIBURG

# Algorithm Theory

**13 Text Search - Knuth, Morris, Pratt, Boyer,  
Moore**

**Christian Schindelhauer**

Albert-Ludwigs-Universität Freiburg  
Institut für Informatik  
Rechnernetze und Telematik  
Wintersemester 2007/08



# Text Search

- ▶ **Scenarios**
- ▶ **Static texts**
  - Literature databases
  - Library systems
  - Gene databases
  - World Wide Web
- ▶ **Dynamic texts**
  - Text editors
  - Symbol manipulators

# Text Search

Data type **string**

- array of character
- file of character
- list of character

Operations: (T, P of type **string**)

**length:** length ()

***i*-th character :** T [*i*]

**concatenation:** cat (T, P) T.P

# Problem Definition

## Given:

Text  $t_1 t_2 \dots t_n \in \Sigma^n$

pattern  $p_1 p_2 \dots p_m \in \Sigma^m$

## Goal:

Find one or all occurrences of the pattern in the text, i.e. positions  $i$  ( $0 \leq i \leq n - m$ ) such that

$$p_1 = t_{i+1}$$

$$p_2 = t_{i+2}$$

$$\vdots$$

$$p_m = t_{i+m}$$

# Problem Definition

**text:**  $t_1 \ t_2 \ \dots \ t_{i+1} \ \dots \ t_{i+m} \ \dots \ t_n$

**pattern:**  $p_1 \ \dots \ p_m$

**Possible Running Time :**

**1. Worst Solution:**

# possible alignments:  $n - m + 1$

# pattern positions:  $m$

$\rightarrow O(n m)$

**2. Best possible solution:**

At least one comparison per  $m$  consecutive text positions:

$\rightarrow \Omega ( m + n/m )$

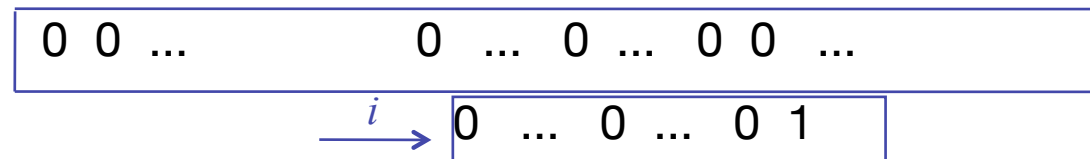
# Naive Method

For each possible position  $0 \leq i \leq n - m$  check at most  $m$  characters pairs. Whenever a mismatch occurs, shift to the next position

```
textsearchbf := proc (T :: string, P :: string)
# Input: text T, pattern P
# Output: list L of positions i, at which P occurs in T
  n := length (T); m := length (P);
  L := [];
  for i from 0 to n-m do
    j := 1;
    while j ≤ m and T[i+j] = P[j]
      do j := j+1 od;
    if j = m+1 then L := [L [], i] fi;
  od;
  RETURN (L)
end;
```

# Naive Method

Running time:



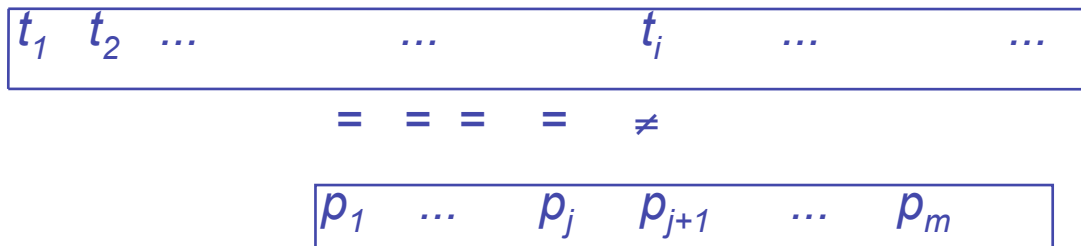
Worst Case:  $\Omega(m n)$

For real world examples, mismatches occur usually very often.

→ Running time  $\sim c n$

# Knuth-Morris-Pratt Algorithm (KMP)

Let  $t_i$  and  $p_{j+1}$  be the characters to be compared:



If, for a certain alignment, the first mismatch occurs for characters  $t_i$  and  $p_{j+1}$  then

- the last  $j$  characters compared in  $T$  equal the first  $j$  characters of  $P$
- $t_i \neq p_{j+1}$



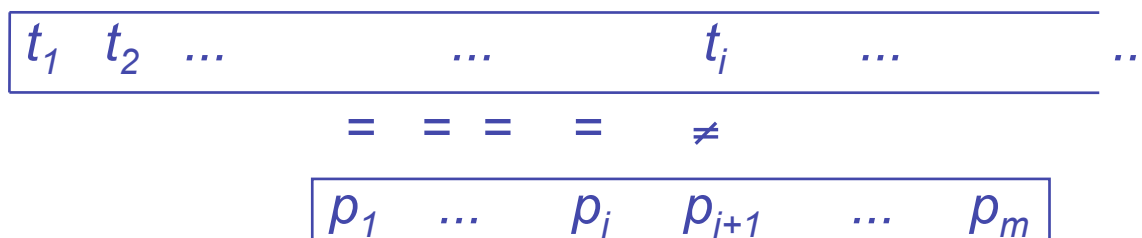
# Knuth-Morris-Pratt Algorithm (KMP)

## Idea:

Find  $j' = \text{next}[j] < j$  such that  $t_j$  can then be compared to  $p_{j'+1}$

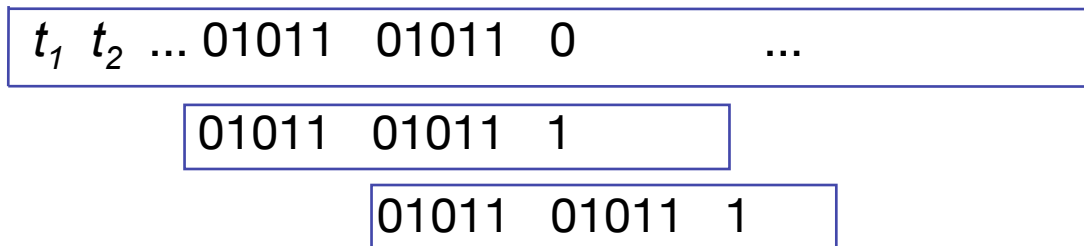
Find maximum  $j' < j$  such that  $P_{1\dots j} = P_{j'+1\dots j}$

Find the longest prefix of  $P$  which is a proper suffix of  $P_{1\dots j}$



# Knuth-Morris-Pratt Algorithm (KMP)

Example for determining next[j]:



next[j] = length of the longest prefix of  $P$  which is a proper  
suffix of  $P_{1 \dots j}$

# Knuth-Morris-Pratt Algorithm (KMP)

⇒ for  $P = 0101101011$ , next = [0,0,1,2,0,1,2,3,4,5] :

1	2	3	4	5	6	7	8	9	10
0	1	0	1	1	0	1	0	1	1
		0							
		0	1						
					0				
					0	1			
					0	1	0		
					0	1	0	1	
					0	1	0	1	1

# Knuth-Morris-Pratt Algorithm (KMP)

```
KMP := proc (T :: string, P :: string)
# Input: text T, pattern P
# Output: list L of positions i at which P occurs in T
  n := length (T); m := length(P);
  L := []; next := KMPnext(P);
  j := 0;
  for i from 1 to n do
    while j>0 and T[i] <> P[j+1] do j := next [j] od;
    if T[i] = P[j+1] then j := j+1 fi;
    if j = m then L := [L[], i-m] ;
      j := next [j]
    fi;
  od;
  RETURN (L);
end;
```

# Knuth-Morris-Pratt Algorithm (KMP)

Pattern: abrakadabra, next = [0,0,0,1,0,1,0,1,2,3,4]

```
a b r a k a d a b r a b r a b a b r a k ...  
| | | | | | | | | |  
a b r a k a d a b r a  
next[11] = 4
```

```
a b r a k a d a b r a b r a b a b r a k ...  
      | | | | |  
      a b r a k a d a b r a  
next[4] = 1
```

# Knuth-Morris-Pratt Algorithm (KMP)

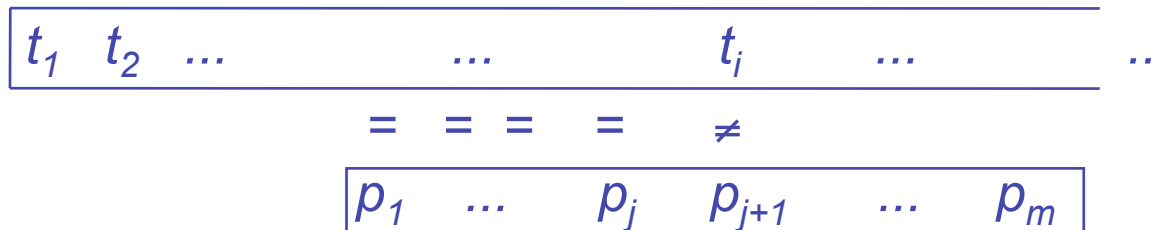
```
a b r a k a d a b r a b r a b a b r a k ...
      - | | | |
      a b r a k
      next [4] = 1

a b r a k a d a b r a b r a b a b r a k ...
      - | |
      a b r a k
      next[2] = 0

a b r a k a d a b r a b r a b a b r a k ...
      | | | | |
      a b r a k
```

# Knuth-Morris-Pratt Algorithm (KMP)

Correctness:



When starting the for-loop:

$$P_{1..j} = T_{i-j..i-1} \text{ und } j \neq m$$

if  $j = 0$ : we are located at the first character of  $P$

if  $j \neq 0$ :  $P$  can be shifted while  $j > 0$  and  $t_i \neq p_{j+1}$

# Knuth-Morris-Pratt Algorithm (KMP)

If  $T[i] = P[j+1]$ ,  $j$  and  $i$  can be increased (at the end of the loop)

If  $P$  has been compared completely ( $j = m$ ) an occurrence of  $P$  in  $T$  has been found and we can shift to the next position.



# Knuth-Morris-Pratt Algorithm (KMP)

## Running time:

- text pointer  $i$  will be never decreased
- text pointer  $i$  and pattern pointer  $j$  are always incremented together:  $\text{next}[j] < j$ ;  
 $j$  can be decreased only as many times as it has been increased

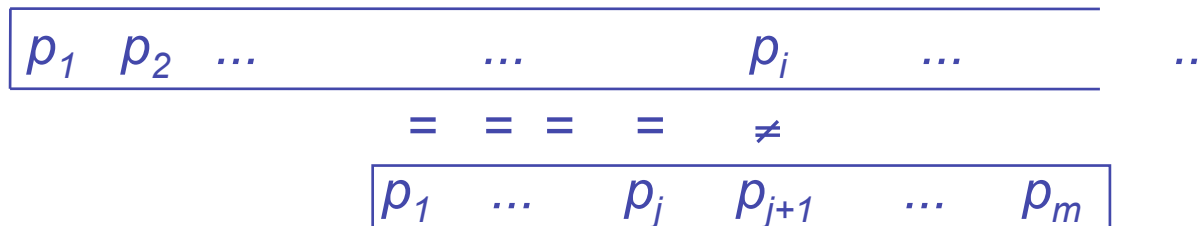
If the next-array is known, KMP runs in time  $O(n)$ .

# Computation of Array *next*

$next[i]$  = length of the longest prefix of  $P$  which is a proper suffix of  $P_{1\dots i}$

$next[1] = 0$

Let  $next[i-1] = j$ :



# Computation of Array *next*

Consider the following two cases

1)  $p_i = p_{j+1} \rightarrow \text{next}[i] = j + 1$

2)  $p_i \neq p_{j+1} \rightarrow$  replace  $j$  by  $\text{next}[j]$  until  $p_i = p_{j+1}$  or  $j = 0$ .

if  $p_i = p_{j+1}$  then  $\text{next}[i] = j + 1$  otherwise  $\text{next}[i] = 0$ .

# Computation of Array *next*

```
KMPnext := proc (P :: string)
#Input   : pattern P
#Output  : next-array for P
  m := length (P);
  next := array (1..m);
  next [1] := 0;
  j := 0;
  for i from 2 to m do
    while j > 0 and P[i] <> P[j+1]
      do j := next [j] od;
    if P[i] = P[j+1] then j := j+1 fi;
    next [i] := j
  od;
  RETURN (next);
end;
```

# Run Time of KMP

The KMP-algorithm runs in time  $O(n+m)$ .

Can text search be any faster?

# The Boyer-Moore-Algorithm (BM)

**Idea:** For any alignment of the pattern with the text, scan the characters from right to left rather than from left to right

**Example:**

```
he said abrakadabut but
 |  |  |  |  |  |  |  |  |  |
but |  |  |  |  |  |  |  |  |
    |  |  |  |  |  |  |  |  |
      |  |  |  |  |  |  |  |  |
        |  |  |  |  |  |  |  |
          |  |  |  |  |  |  |  |
            |  |  |  |  |  |  |
              |  |  |  |  |  |
                |  |  |  |  |
                  |  |  |  |
                    |  |  |
                     |  |
                      |
```

**Large jumps and few comparisons**

**Desired running time:  $O(m+n/m)$**

# BM – Last-Occurrence Function

For  $c \in \Sigma$  and pattern  $P$  let

$\delta(c) :=$  index of the right-most occurrence of  $c$  in  $P$

$$= \max \{j \mid p_j = c\}$$

$$= \begin{cases} 0 & \text{if } c \notin P \\ j & \text{if } c = p \text{ and } c \neq p, \text{ for } j < k \leq m \end{cases}$$

How hard is it to compute all values of  $\delta$ ?

Let  $|\Sigma| = l$ : Running time:  $O(m+l)$

# BM – Last-Occurrence Function

Let

$c$  = the character causing the mismatch

$j$  = the index of the current character in the pattern ( $c \neq p_j$ )

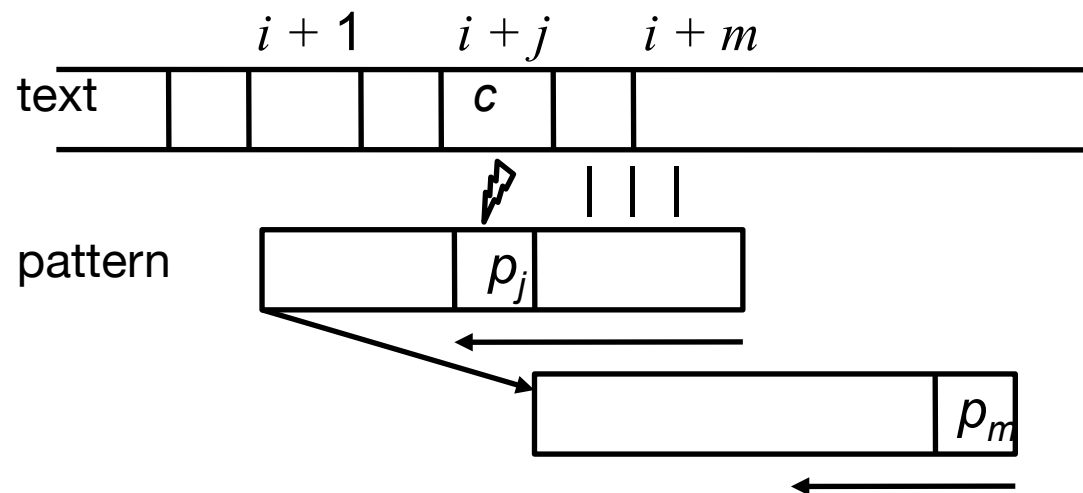


# BM – Last-Occurrence Function

## Computation of the pattern shift

**Case 1**  $c$  does not occur in the pattern  $P$  ( $\delta(c) = 0$ )

Shift the pattern  $j$  characters to the right

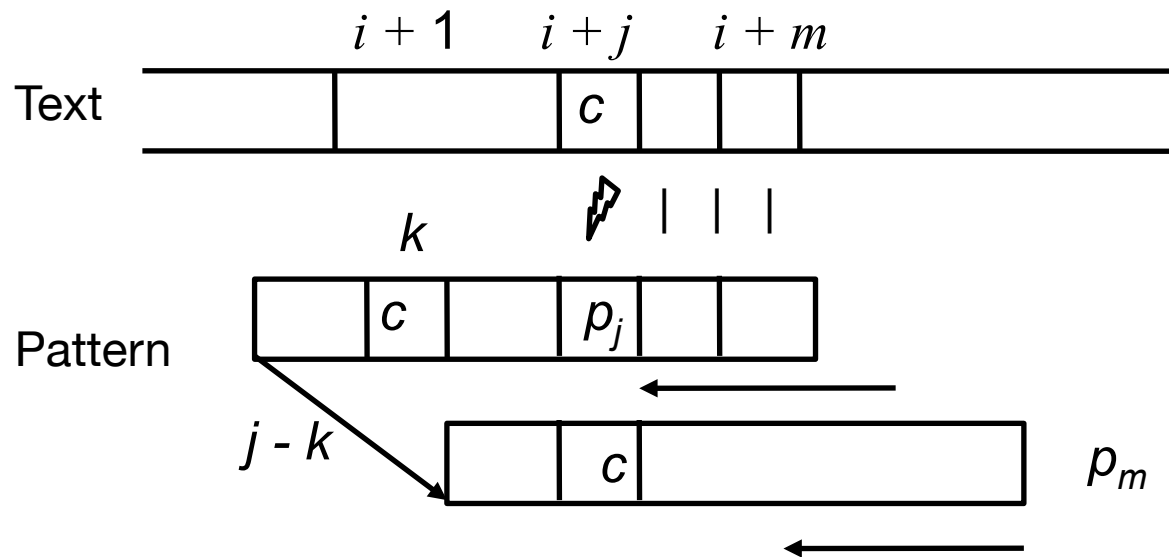


$$\Delta(i) = j$$

# BM – Last-Occurrence Function

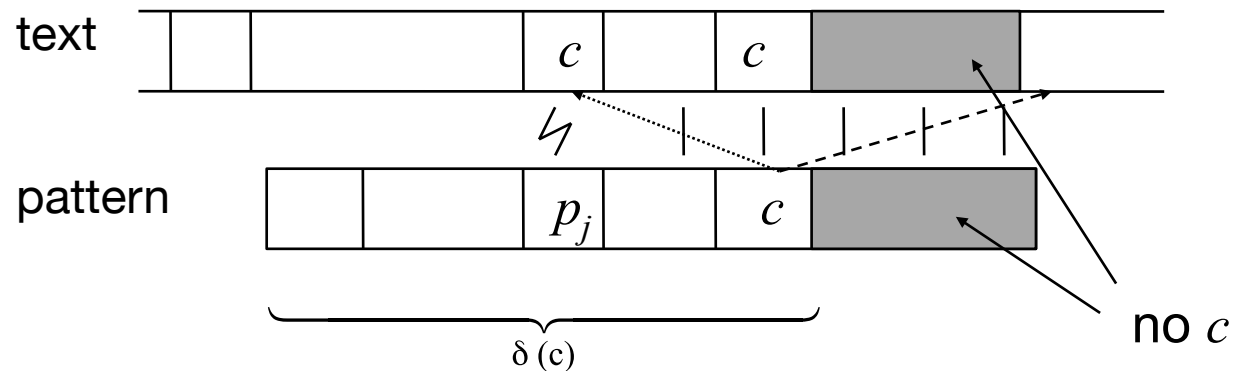
**Case 2**  $c$  occurs in the pattern. ( $\delta(c) \neq 0$ )

Shift the pattern to the right until the rightmost  $c$  in the pattern is aligned with a potential  $c$  in the text.



# BM – Last-Occurrence Function

Case 2 a:  $\delta(c) > j$

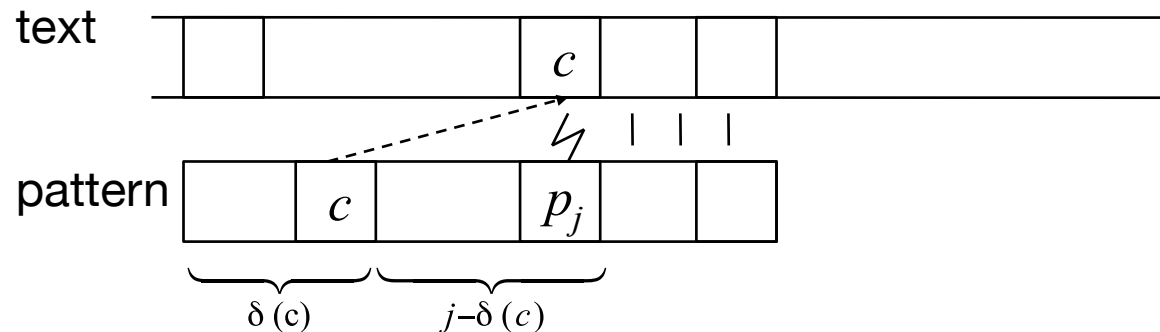


Shift the rightmost  $c$  in the pattern to a potential  $c$  in the text

shift by  $\Delta(i) = m - \delta(c) + 1$

# BM – Last-Occurrence Function

Case 2 b:  $\delta(c) < j$



Shift the rightmost c in the pattern to c in the text

shift by  $\Delta(i) = j - \delta(c)$

# Boyer-Moore Algorithm Version 1

## Algorithm *BM-search1*

**Input:** text  $T$  and pattern  $P$

**Output:** all positions of  $P$  in  $T$

```
1  $n := \text{length}(T)$ ;  $m := \text{length}(P)$ 
2 compute  $\delta$ 
3  $i := 0$ 
4 while  $i \leq n - m$  do
5    $j := m$ 
6   while  $j > 0$  and  $P[j] = T[i + j]$  do
7      $j := j - 1$ 
8   end while;
```

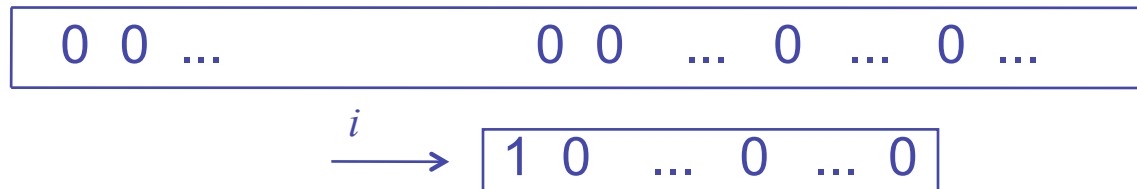
```
8   if  $j = 0$ 
9     then output position  $i$ 
10     $i := i + 1$ 
11   else if  $\delta(T[i + j]) > j$ 
12     then  $i := i + m + 1 - \delta[T[i + j]]$ 
13     else  $i := i + j - \delta[T[i + j]]$ 
14   end while;
```

# Boyer-Moore Algorithm Version 1

## Analysis:

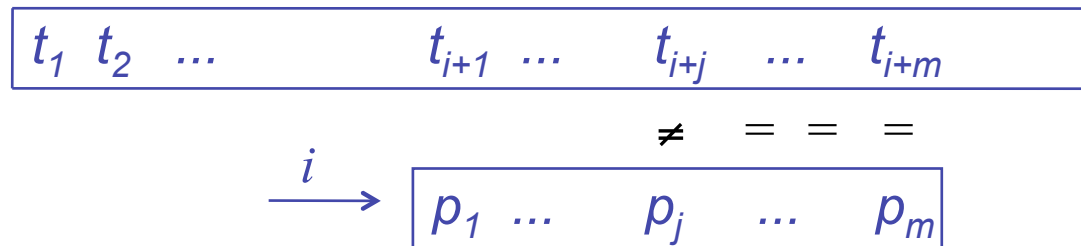
Desired running time :  $O(m + n/m)$

Worst-case running time:  $\Omega(n m)$



# Match Heuristic

Use the information collected before a mismatch  $p_j \neq t_{i+j}$  occurs.



$gsf[j]$  = position of the end of the next occurrence of the suffix  $P_{j+1} \dots m$  from the right that is not preceded by character  $P_j$   
 (good suffix function)

Possible shift:  $\gamma[j] = m - gsf[j]$

# Example for Computing $gsf$

$gsf[j]$  = position of the end of the closest occurrence of the suffix  $P_{j+1 \dots m}$  from the right that is not preceded by character  $P_j$

Pattern: banana

$gsf[j]$	inspected suffix	forbidden character	Other matches	Position
$gsf[5]$	a	n	banana <u>  </u>	2
$gsf[4]$	na	a	*** bana na <u>  </u>	0
$gsf[3]$	ana	n	banana <u>  </u>	4
$gsf[2]$	nana	a	banana <u>  </u>	0
$gsf[1]$	anana	b	banana <u>  </u>	0
$gsf[0]$	banana	$\epsilon$	banana <u>  </u>	0



# Example for Computing $gsf$

$$\Rightarrow gsf(\text{banana}) = [0,0,0,4,0,2]$$

a b a a b a b a n a n a n a n a  
≠ = = =  
b a n a n a  
b a n a n a

# Boyer-Moore Algorithm Version 2

## Algorithm *BM-search2*

**Input:** text  $T$  and pattern  $P$

**Output:** shift for all occurrences of  $P$  in  $T$

1  $n := \text{length}(T)$ ;  $m := \text{length}(P)$

2 compute  $\delta$  and  $\gamma$

3  $i := 0$

4 **while**  $i \leq n - m$  **do**

5      $j := m$

6     **while**  $j > 0$  **and**  $P[j] = T[i + j]$  **do**

7          $j := j - 1$

**end while;**

8     **if**  $j = 0$

9         **then** output position  $i$

10              $i := i + \gamma[0]$

11         **else**  $i := i + \max(\gamma[j], j - \delta[T[i + j]])$

12     **end while;**



ALBERT-LUDWIGS-  
UNIVERSITÄT FREIBURG

# Algorithm Theory

**13 Text Search - Knuth, Morris, Pratt, Boyer,  
Moore**

**Christian Schindelhauer**

Albert-Ludwigs-Universität Freiburg  
Institut für Informatik  
Rechnernetze und Telematik  
Wintersemester 2007/08

