

# Special course in Computer Science: Advanced Text Algorithms

## Lecture 2: Pattern Matching Algorithms

Eugen Czeizler

Department of IT, Abo Akademi

<http://combio.abo.fi/teaching/textalg/>

(slides originally by I. Petre, E. Czeizler, V. Rogojin)

# Pattern Matching Problem

- The pattern matching problem is one of the most investigated problems in text algorithms.
- Implementations of algorithms for this problem are used daily for accessing information
  - use Google search engine
  - make a data base query
  - use the “search” or “replace” commands in text editors
- Simple example of a pattern matching problem:
  - Find occurrences of the word “advanced” in the text “Special course in Computer Science: advanced text algorithms”

# Pattern Matching Problem

- A **pattern** is a nonempty language which does not contain the empty string:
  - a single string
  - a finite or infinite set of strings
- The **pattern matching problem**: Given a text and a pattern, find occurrences of the pattern within the text.

# Pattern Matching Problem

- We will consider for now only the case when the pattern is a single string.
- The input for the pattern matching algorithm:
  - the pattern `pat` of length `m`
  - the text `text` of length `n`
- The output of the algorithm:
  - A boolean value: 1 if `pat` occurs in the text or 0 otherwise.
  - A set of positions where the pattern occurs in the given text.

# Notations and Definitions

- A **period** of a word  $w$  is an integer  $0 < p < |w|$  such that  $w[i] = w[i+p]$  for all  $i \in \{1, \dots, |w| - p\}$ .
  - **Example:**
    - $abcabc$  has period 3:  $abc\ abc$
    - $abababa$  has periods 2, 4, and 6:  $ab\ ab\ ab\ a$  ,  $abab\ aba$  ,  $\underbrace{ababab}_6\ a$
- We denote by **period(w)** the smallest period of  $w$ .
- A word  $u$  is a **border** of a word  $v$  if  $u$  is both a prefix and a suffix of  $v$ .
  - **Example:**  $abaaaab$

# Notations and Definitions

- Just as in the case of periods, there are words which have several borders.
- Example: ababab

*ab abab*

*abab ab*

*abab ab*

*ab abab*

- Denote by **Border(w)** the longest nontrivial border of a word  $w$ , i.e.,  $\text{Border}(w) \neq w$ .

# Pattern Matching Algorithms

- We will discuss two fundamental pattern matching algorithms:
  - Knuth-Morris-Pratt (KMP) algorithm
  - Boyer-Moore (BM) algorithm
- Both of them have two stages:
  - Pre-processing stage: computing some additional tables
  - Pattern-searching stage

# Brute-force algorithm

- The scheme for a brute-force algorithm for the pattern-matching problem:

for  $i := 0$  to  $n - m$  do

    check if  $pat = text[i + 1 .. i + m]$

- This algorithm is the origin for both KMP and BM algorithms.
- The actual implementation of this algorithm depends on how we implement the function “check”, e.g., scanning the text from the left, or from the right, or other method.



# Brute-force algorithm

- Algorithm brute-force1

```
i := 0
while i ≤ n − m do begin
    j := 0;
    while j < m and pat[j + 1] = text[i + j + 1] do
        j := j + 1;
    if j = m then return(true);
    i := i + 1;
end;
return(false)
```

- the index *i* scans the positions of the text
- the index *j* scans the positions of the pattern
- if we found the pattern, then we stop
- otherwise, i.e., we have a mismatch between a character in the pattern and a character in the text, we increase *i* by 1 and we start comparing again *pat* and *text*.

# Brute-force algorithm

- Algorithm brute-force1

$i := 0$

while  $i \leq n - m$  do begin

$j := 0$ ;

    while  $j < m$  and  $pat[j + 1] = text[i + j + 1]$  do

$j := j + 1$ ;

    if  $j = m$  then return(*true*);

$i := i + 1$ ;

end;

return(*false*)

Let  $text = aaaaabaa$  and  
 $pat = aaaba$

- $i = 0$

$aaab$

$aaaabaa$

- $i = 1$

$aaab$

$aaaaabaa$

- $i = 2$

$aaaba$

$aaaabaa$

# Brute-force algorithm

- Algorithm brute-force1

```
i := 0
while i ≤ n − m do begin
    j := 0;
    while j < m and pat[j + 1] = text[i + j + 1] do
        j := j + 1;
    if j = m then return(true);
    i := i + 1;
end;
return(false)
```

- This algorithm scans both the text and the pattern from left to right.
- The function return(*x*) outputs *x* and stops the whole algorithm.
- The shifts have length 1.
- If  $\text{pat} = a^n b$  and  $\text{text} = a^{2n-1} b$ , then the brute-force algorithm performs a **quadratic number of comparisons**.

# The Brute-force algorithm

- The main drawback of the brute-force algorithm is that the shifts in the search have length 1.
- Moreover, after a shift of the pattern, the brute-force algorithm has forgotten all information about previously matched symbols.
- So it is possible that it re-compares a text symbol with different pattern symbols again and again.
- This leads to its worst case complexity of  $\Theta(nm)$  ( $n$ : length of the text,  $m$ : length of the pattern).

# The Brute-force algorithm

- This complexity is too large, and makes the use of this algorithm unpractical
  - in practice we search a given pattern within huge texts, or data bases
- Most used algorithms for pattern matching (Knuth-Morris-Pratt, Boyer-Moore) have linear time complexity and are more suitable for practical use.

# The Knuth-Morris-Pratt algorithm

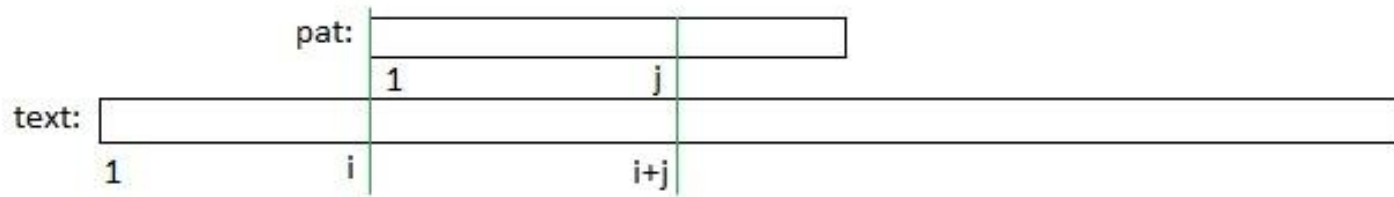
- The Knuth-Morris-Pratt algorithm uses the information gained by previous symbol comparisons.
- It never compares again a text symbol that has already matched a pattern symbol.
- As a result, the complexity of the searching phase of the Knuth-Morris-Pratt algorithm is in  $O(n)$ .

# The Knuth-Morris-Pratt algorithm

- However, a pre-processing of the pattern is necessary in order to analyze its structure.
- The pre-processing phase has a complexity of  $O(m)$ .
- Since  $m \leq n$ , the overall complexity of the Knuth-Morris-Pratt algorithm is in  $O(n)$ .

# Morris-Pratt (MP) algorithm

- Let  $\text{pat}[1..j] = \text{text}[i+1..i+j]$  for some  $j \geq 1$ .

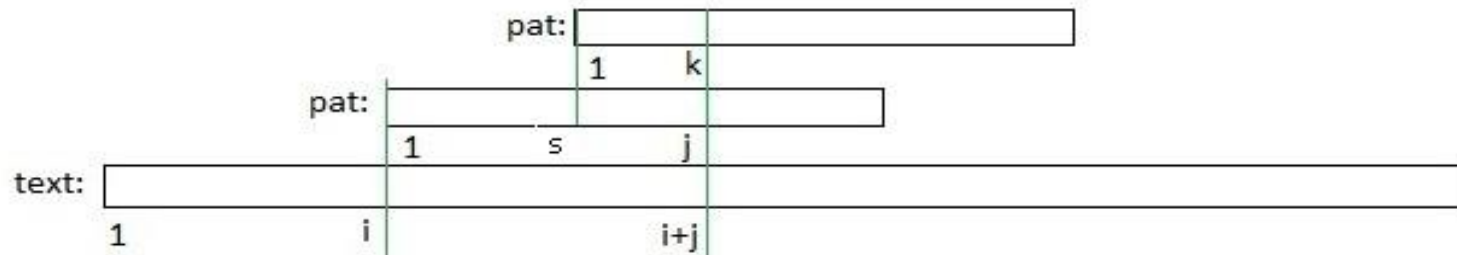


- We want to shift the search by more than 1 unit.
- Safe shift  $s$**  = a positive integer  $s$  such that there is no occurrence of the pattern at positions  $i, i+1, \dots, i+s-1$ , but there might be one at  $i+s$ .



# Safe shift

- Suppose an occurrence of the pattern starts at position  $i+s$  and let  $k=j-s$ .

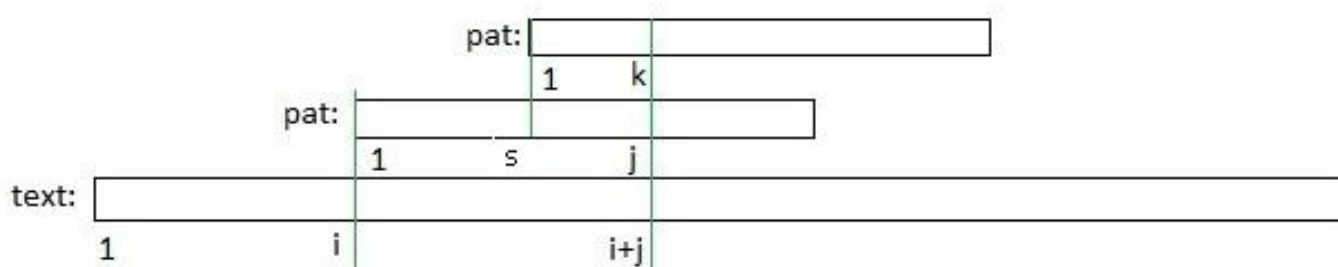


- Then  $\text{pat}[1..k]$  and  $\text{pat}[1..j]$  are suffixes of the same text:  $\text{text}[1..i+j]$
- $\text{cond}(j,k)$ :  $\text{pat}[1..k]$  is a proper suffix of  $\text{pat}[1..j]$
- So,  $s$  is a safe shift if  $k=j-s$  is the largest integer satisfying  $\text{cond}(j,k)$

# Safe shift

- $\text{cond}(j,k)$ :  $\text{pat}[1..k]$  is a proper suffix of  $\text{pat}[1..j]$
- Denote  $\text{Bord}(j) =$  the largest integer  $k$  satisfying  $\text{cond}(j,k)$
- $\text{Bord}(j) = |\text{Border}(\text{pat}[1..j])|$
- Then, the smallest safe shift is  

$$\text{MP\_Shift}(j) = j - \text{Bord}(j)$$



# Table of borders

- The table **Bord**, called the **table of borders**, is essential for this algorithm.
- It carries the information which allows us to compute the length of a safe shift ( $s=j-\text{Bord}(j)$ ) in the case of a mismatch.

# Example of a table of borders

- $\text{Bord}(j) = |\text{Border}(\text{pat}[1..j])|$
- Example: Let  $\text{pat} = \text{abababababb}$ . Then,  
 $\text{Bord}(1) = |\text{Border}(\text{pat}[1])| = |\text{Border}(a)| = 0,$   
 $\text{Bord}(2) = |\text{Border}(\text{pat}[1..2])| = |\text{Border}(ab)| = 0,$   
 $\text{Bord}(3) = |\text{Border}(\text{pat}[1..3])| = |\text{Border}(aba)| = 1,$   
 $\text{Bord}(4) = |\text{Border}(\text{pat}[1..4])| = |\text{Border}(abab)| = 2, \dots$   
 $\text{Bord} = [0, 0, 1, 2, 3, 4, 5, 6, 7, 8, 0]$
- If  $j=0$ ,  $\text{pat}[1..j] = \varepsilon$ , then the shift must be 1 and we define  $\text{Bord}(0) = -1$

# Computing the table of borders

- As a pre-processing step in the MP-algorithm we compute the table **Bord**.

```
procedure compute_Bord
   $Bord(0) := -1; t := -1;$ 
  for  $j := 1$  to  $m$  do begin
    while  $t \geq 0$  and  $pat[t+1] \neq pat[j]$  do
       $t := Bord(t);$ 
     $t := t + 1; Bord(j) := t;$ 
  end
```

# Computing the table of borders

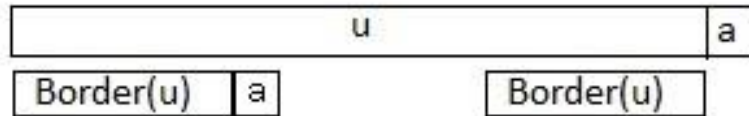
- Lemma: For a word  $u$  and a letter  $a$ , we have

$$\text{Border}(ua) = \begin{cases} \text{Border}(u)a, & \text{if } \text{Border}(u)a \text{ is a prefix of } u, \\ \text{Border}(\text{Border}(u)a), & \text{otherwise} \end{cases}$$

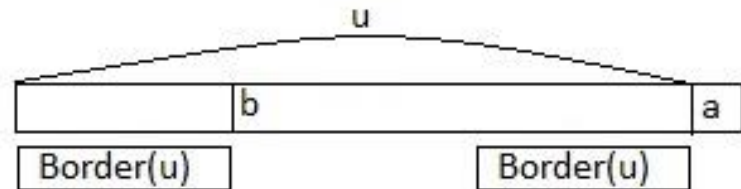
Proof: (not required for the exam)

- If  $\text{Border}(ua)$  is not empty, then  $\text{Border}(ua) = wa$ , where  $w$  is a border of  $u$
- If  $\text{Border}(u)a$  is a prefix of  $u$ , then obviously  $\text{Border}(ua) = \text{Border}(u)a$

- $\text{Border}(u)a$  is a border of  $ua$
- It is the maximal border



- Otherwise, i.e., if  $\text{Border}(u)a$  is not a prefix of  $u$ 
  - $\text{Border}(ua)$  is a prefix of  $\text{Border}(u)$  since  $|\text{Border}(ua)| \leq |\text{Border}(u)| + 1$
  - $\text{Border}(ua)$  is also a suffix of  $\text{Border}(u)a$
  - $\text{Border}(ua)$  is maximal with this property.
- So,  $\text{Border}(ua) = \text{Border}(\text{Border}(u)a)$ .



# Computing the table of borders

```
procedure compute_Bord
```

```
   $Bord(0) := -1; t := -1;$ 
```

```
  for  $j := 1$  to  $m$  do begin
```

```
    while  $t \geq 0$  and  $pat[t+1] \neq pat[j]$  do
```

```
       $t := Bord(t);$ 
```

```
     $t := t + 1; Bord(j) := t;$ 
```

```
  end
```

- Lemma: The procedure `compute_Bord` applied to a string `pat` of length `m` produces the table of borders for `pat`.

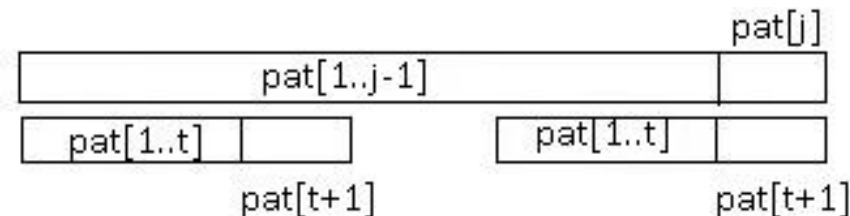
Proof:

- The table `Bord` is computed sequentially by the procedure `compute_Bord`: from the prefix of length 1 to `pat` itself.

- At the beginning of each “for cycle” `j`,  $t = \text{Bord}(j-1)$ , i.e., `pat[1..t]` is both a prefix and a suffix of `pat[1..j-1]`, and is maximal with this property.

- If  $pat[t+1] = pat[j]$ , i.e., the previous border can be extended by 1 character, i.e.,  $\text{Bord}(j) = \text{Bord}(j-1) + 1$

- Otherwise, i.e.,  $pat[t+1] \neq pat[j]$ , we compute recursively the new border in the while loop, according to the previous lemma.



# Computing the table of borders

Example:  $pat=abaababaaba$ ,  $m=11$

$Bord(0)=-1$ ;  $t=-1$

$j=1$

$t:=0$ ;  $Bord(1):=0$ ;

$j=2$

$t \geq 0$ ;  $pat[1] \neq pat[2]$

$t:=Bord(0)=-1$

$t:=0$ ;  $Bord(2)=0$ ;

$j=3$

$t \geq 0$ ;  $pat[1]=pat[3]$

$t:=1$ ;  $Bord(3):=1$ ;

$j=4$

$t \geq 0$ ;  $pat[2] \neq pat[4]$

$t:=Bord(1)=0$

$pat[1]=pat[4]$

$t:=1$ ;  $Bord(4)=1$ ;

```
procedure compute_Bord
```

```
   $Bord(0) := -1$ ;  $t := -1$ ;
```

```
  for  $j := 1$  to  $m$  do begin
```

```
    while  $t \geq 0$  and  $pat[t+1] \neq pat[j]$  do
```

```
       $t := Bord(t)$ ;
```

```
     $t := t + 1$ ;  $Bord(j) := t$ ;
```

```
  end
```



# Computing the table of borders

```
procedure compute_Bord
   $Bord(0) := -1; t := -1;$ 
  for  $j := 1$  to  $m$  do begin
    while  $t \geq 0$  and  $pat[t+1] \neq pat[j]$  do
       $t := Bord(t);$ 
     $t := t + 1; Bord(j) := t;$ 
  end
```

Example:  $pat = abaababaaba, m = 11$

$j = 5$

$t = 1 \geq 0; pat[2] = pat[5]$

$t := 2; Bord(5) = 2$

$j = 6$

$t \geq 0; pat[3] = pat[6]$

$t := 3; Bord(6) = 3$

$j = 7$

$t \geq 0; pat[4] \neq pat[7]$

$t := Bord(3) = 1$

$t \geq 0; pat[2] = pat[7]$

$t := 2; Bord(7) = 2$

# Computing the table of borders

```
procedure compute_Bord
```

```
   $Bord(0) := -1; t := -1;$ 
```

```
  for  $j := 1$  to  $m$  do begin
```

```
    while  $t \geq 0$  and  $pat[t+1] \neq pat[j]$  do
```

```
       $t := Bord(t);$ 
```

```
     $t := t + 1; Bord(j) := t;$ 
```

```
  end
```

• Lemma: The maximum number of character comparisons executed by algorithm `compute_Border` is  $2m-3$ .

Proof:

- we use a so-called “store principle” where  $t$  is seen as the number of items in store
- when  $j=1$ , we have  $t < 0$ , so no character comparison is done and  $t$  becomes 0 (i.e., the store is empty at the beginning)
- for each  $2 \leq j \leq m$  we add at most one element in the store (we execute the statement  $t := t + 1$ ). So the maximal size of the store is  $m-2$ .

# Computing the table of borders

procedure compute\_Bord

$Bord(0) := -1; t := -1;$

for  $j := 1$  to  $m$  do begin

  while  $t \geq 0$  and  $pat[t+1] \neq pat[j]$  do

$t := Bord(t);$

$t := t + 1; Bord(j) := t;$

end

- each time we execute statement  $t := Bord(t)$ ,  $t$  is decreased, i.e., we eliminate some elements from the store.

- So there are at most  $m-2$  unsuccessful character comparisons  $pat[t+1] \neq pat[j]$  (for each of them we execute the statement  $t := Bord(t)$ )

- For each  $2 \leq j \leq m$ , there is at most one successful character comparison.

- Thus the total number of character comparisons is at most  $(m-2) + (m-1) = 2m-3$

# Morris-Pratt (MP) algorithm

- Using the table of borders **Bord** and the safe shift we obtain an improved version of the algorithm brute-force1, i.e., **the Morris-Pratt (MP) algorithm**:

Algorithm MP;

$i := 0; j := 0;$

while  $i \leq n - m$  do begin

    while  $j < m$  and  $pat[j + 1] = text[i + j + 1]$  do

$j = j + 1;$

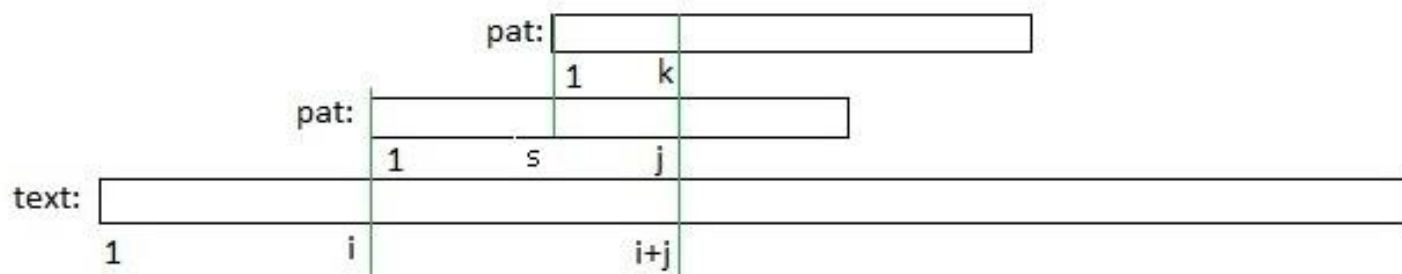
    if  $j = m$  then *return(true)*;

$i := i + MP\_Shift(j); j := \max(0, j - MP\_Shift(j));$

end;

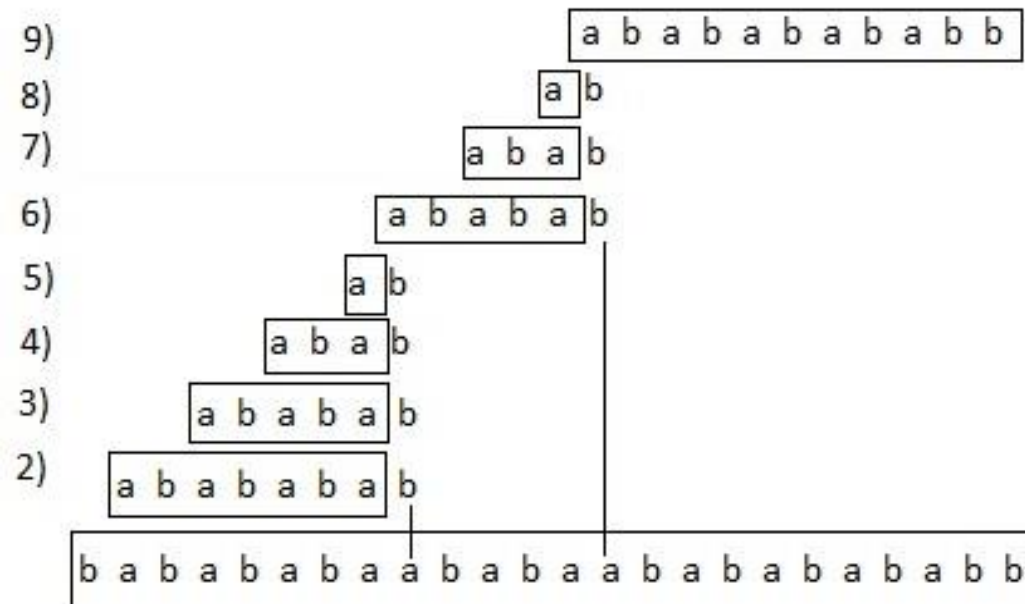
*return(false)*

# Morris-Pratt (MP) algorithm



# Morris-Pratt (MP) algorithm

- Example: Let  $\text{text} = \text{babababaabababababb}$  and  $\text{pat} = \text{ababababb}$ .



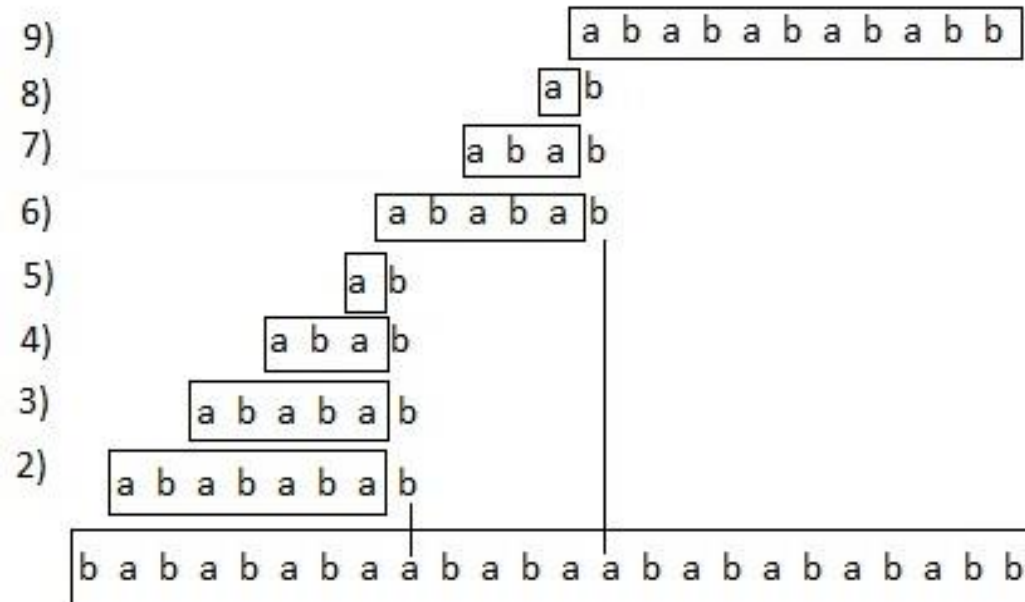
1)  $i=0, j=0 \text{ pat}[1] \neq \text{text}[1]$

2)  $i=i+\text{MP\_Shift}(j)=$   
 $=i+j-\text{Bord}(j)=0+0-(-1)=1;$   
 $j=\max(0, j-\text{MP\_Shift}(j))=0$

$\text{pat}[1]=\text{text}[2]; \text{pat}[2]=\text{text}[3]$   
 $\text{pat}[3]=\text{text}[4]; \text{pat}[4]=\text{text}[5]$   
 $\text{pat}[5]=\text{text}[6]; \text{pat}[6]=\text{text}[7]$   
 $\text{pat}[7]=\text{text}[8]; \text{pat}[8] \neq \text{text}[9]$   
 $j=7, i=1$

# Morris-Pratt (MP) algorithm

- Example: Let  
`text=babababaababaabababab`  
`abb` and `pat=ababababbb`.



$$\begin{aligned}
 3) \quad i &= i + \text{MP\_Shift}(j) = \\
 &= i + j - \text{Bord}(j) = 1 + (7 - 5) = 3 \\
 j &= \max(0, j - \text{MP\_Shift}(j)) = \\
 &= \max(0, 7 - 2) = 5 \\
 \text{pat}[6] &\neq \text{text}[9]
 \end{aligned}$$

$$\begin{aligned}
 4) \quad i &= i + \text{MP\_Shift}(j) = \\
 &= 3 + \text{MP\_Shift}(5) = 3 + (5 - 3) = 5 \\
 j &= \max(0, j - \text{MP\_Shift}(j)) = \\
 &= \max(0, 5 - 2) = 3 \\
 \text{pat}[4] &\neq \text{text}[9]
 \end{aligned}$$

# Morris-Pratt (MP) algorithm

- Lemma: The time complexity of the algorithm MP is linear in the length of the text (i.e.,  $O(n)$ ). The maximal number of character comparisons is  $2n-m$ .

## Proof:

- Let  $T(n)$  be the maximal number of comparisons  $pat[j+1]=text[i+j+1]$ .
- There are at most  $n-m+1$  unsuccessful comparisons: at most one for each  $i$  ( $0 \leq i \leq n-m$ )
- $0 \leq i+j \leq n$

Algorithm MP;

$i := 0; j := 0;$

while  $i \leq n - m$  do begin

    while  $j < m$  and

$pat[j+1] = text[i+j+1]$  do

$j = j + 1;$

    if  $j = m$  then *return(true)*;

$i := i + MP\_Shift(j);$

$j := \max(0, j - MP\_Shift(j));$

end;

*return(false)*



# Morris-Pratt (MP) algorithm

- Each time we have a successful comparison we increase  $i+j$  by 1 and this value never decreases.
- So there are at most  $n$  successful comparisons.
- If the first comparison is unsuccessful, then we do not have any successful comparison for  $i=0$ .
- So,  $T(n) \leq n + n - m = 2n - m$ .
- If  $\text{text} = a^n$  and  $\text{pat} = ab$ , then we actually have  $T(n) = 2n - m$

Algorithm MP;

$i := 0; j := 0;$

while  $i \leq n - m$  do begin

while  $j < m$  and

$\text{pat}[j + 1] = \text{text}[i + j + 1]$  do

$j = j + 1;$

if  $j = m$  then *return(true)*;

$i := i + MP\_Shift(j);$

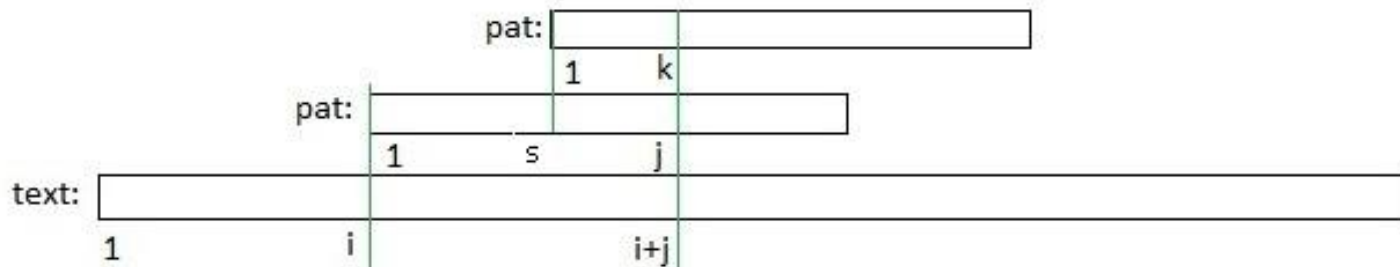
$j := \max(0, j - MP\_Shift(j));$

end;

*return(false)*

# Unnecessary comparisons in the Morris-Pratt algorithm

- If we had a mismatch in the MP-algorithm , i.e.,  $pat[j+1] \neq text[i+j+1]$ , then the next comparison is between  $text[i+j+1]$  and  $pat[k+1]$ , where  $k = Bord(j)$ .
- If  $pat[k+1] = pat[j+1]$ , then we have the same mismatch.



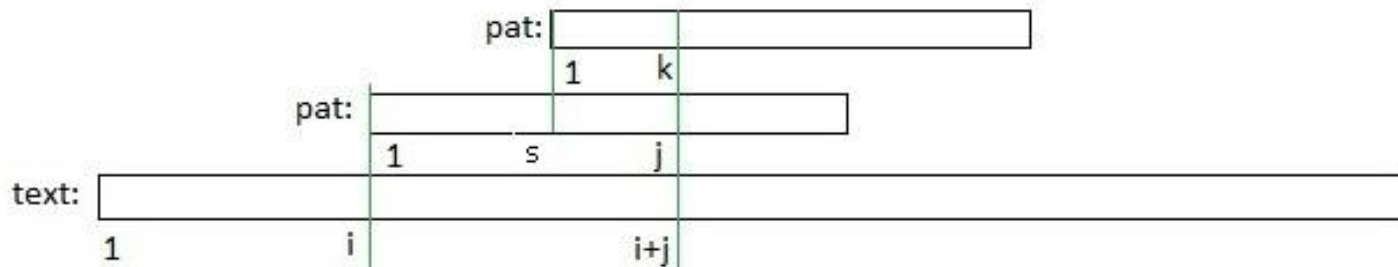
# Unnecessary comparisons in the Morris-Pratt algorithm

- So, instead of the condition used in the MP algorithm:

$\text{cond}(j,k)$ :  $\text{pat}[1..k]$  is a proper suffix of  $\text{pat}[1..j]$

we can use a stronger one for  $0 \leq j < m$ :

$\text{strong\_cond}(j,k)$ :  $\text{pat}[1..k]$  is a proper suffix of  $\text{pat}[1..j]$  and  $\text{pat}[k+1] \neq \text{pat}[j+1]$



# Knuth-Morris-Pratt (KMP) algorithm

- We define **Strong\_Bord(j)**:
  - the largest integer  $k$  satisfying **strong\_cond(j,k)**,
  - -1 otherwise,
  - **Strong\_Bord(m)=Bord(m)**.
- **Example:** For the pattern **pat=abaab** we have the following functions **Bord** and **Strong\_Bord**:

**Bord(0)=-1**

**Bord(1)=0**

**Bord(2)=0**

**Bord(3)=1**

**Bord(4)=1**

**Bord(5)=2**

**Strong\_Bord(0)=-1**

**Strong\_Bord(1)=0**

**Strong\_Bord(2)=-1**

**Strong\_Bord(3)=1**

**Strong\_Bord(4)=0**

**Strong\_Bord(5)=2**

# Knuth-Morris-Pratt (KMP) algorithm

- The function MP\_Shift from MP-algorithm is replaced by  $KMP\_Shift(j) = j - Strong\_Bord(j)$
- For the new, improved algorithm, the table **Strong\_Bord** becomes essential and is computed in the pre-processing step.

# Computation of the table Strong\_Bord

- Let  $t = \text{Bord}(j)$
- If  $\text{pat}[t+1] \neq \text{pat}[j+1]$ , then  $\text{Strong\_Bord}(j) = \text{Bord}(j)$ .
- Otherwise, i.e.,  $\text{pat}[t+1] = \text{pat}[j+1]$ , we have  
 $\text{Strong\_Bord}(j) = \text{Strong\_Bord}(t)$  ( $= \text{Strong\_Bord}(\text{Bord}(j))$ )
- $\text{Strong\_Bord}(|\text{pat}|) = \text{Bord}(|\text{pat}|)$  (as if  $\text{pat}$  is followed by an end marker)

# Computation of the table Strong\_Bord

```
procedure compute_Strong_Bord
  Strong_Bord(0) := -1; t := -1;
  for j := 1 to m do begin
    while t ≥ 0 and pat[t + 1] ≠ pat[j] do t := Strong_Bord(t);
    t := t + 1;
    if j = m or pat[t + 1] ≠ pat[j + 1] then Strong_Bord(j) := t;
    else Strong_Bord(j) := Strong_Bord(t);
  end
```

- **Example:** For the pattern  $pat = aba^{m-2}$ , the table Strong\_Bord computed by this algorithm is:  
Strong\_Bord(0) = -1, Strong\_Bord(1) = 0,  
Strong\_Bord(2) = -1, Strong\_Bord(j) = 1, for all  $3 \leq j \leq m$ .

# Computation of the table Strong\_Bord

```
procedure compute_Strong_Bord
  Strong_Bord(0) := -1; t := -1;
  for j := 1 to m do begin
    while t ≥ 0 and pat[t + 1] ≠ pat[j] do t := Strong_Bord(t);
    t := t + 1;
    if j = m or pat[t + 1] ≠ pat[j + 1] then Strong_Bord(j) := t;
    else Strong_Bord(j) := Strong_Bord(t);
  end
```

- Lemma: The maximum number of character comparisons executed by algorithm compute\_Border is  $3m-5$ .



# Knuth-Morris-Pratt (KMP) algorithm

- The algorithm KMP is the algorithm MP where **Bord** is replaced by **Strong\_Bord** and **MP\_Shift** is replaced by **KMP\_Shift**

Algorithm KMP;

$i := 0; j := 0;$

while  $i \leq n - m$  do begin

    while  $j < m$  and  $pat[j + 1] = text[i + j + 1]$  do

$j = j + 1;$

    if  $j = m$  then *return(true)*;

$i := i + KMP\_Shift(j); j := \max(0, j - KMP\_Shift(j));$

end;

*return(false)*

# Knuth-Morris-Pratt (KMP) algorithm

- Lemma: The complexity of Knuth-Morris-Pratt algorithm is in  $O(n)$ .
  - The complexity of the searching phase of the Knuth-Morris-Pratt algorithm is in  $O(n)$ .
  - The pre-processing phase has a complexity of  $O(m)$ .
  - Since  $m \leq n$ , the overall complexity of the Knuth-Morris-Pratt algorithm is in  $O(n)$ .