

MAXIME CROCHEMORE

King's College London

Maxime.Crochemore@kcl.ac.uk
<http://www.dcs.kcl.ac.uk/staff/mac/>

Examples

- ★ Naive search with backward scan (1)

a a c a a a a a b a b a a b a b . .
a b a a a a
 a b a a a a
 a b a a a a
 a b a a a a
 a b a a a a
 a b a a a a
 a b a a a a
 a b a a a a
 . . .

- ★ Naive search with backward scan (2)

a b a b a a a b a b b b a a b a b . .
a b a a a b a b
 a b a a a b a b
 a b a a a b a b
 a b a a a b a b
 a b a a a b a b
 a b a a a b a b
 a b a a a b a b
 a b a a a b a b
 . . .

Right-to-left scan — shift

```

text ... c g c t c g c g c t a t c g ...
pattern      c g c t a g c
                  c g c t a g c
                  c g c t a g c

```

- ★ Match shift: good-suffix rule (function d)
 - ★ Occurrence shift heuristics: bad-character rule
 - ★ Extra rules if memorization
-

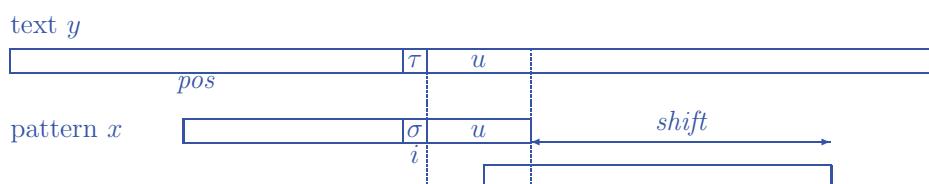
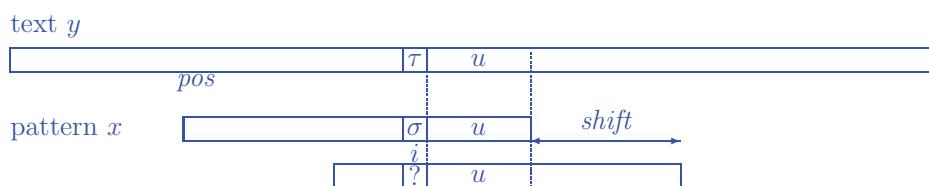
while window on text **do**

```

 $u \leftarrow$  longest common suffix of window and pattern
if  $u = \text{pattern}$  then report a match
shift window  $d(u)$  places to the right

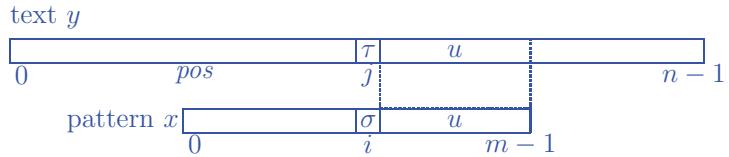
```

Match shift



- ★ **Precomputation**
of rightmost occurrences of u 's: $O(m)$
- ★ Second shift length = a period of x
- ★ Table D implements the good-suffix rule:
 $\text{shift} = d(u) = D[i]$

BM algorithm



- ★ No memorization of previous matches
-

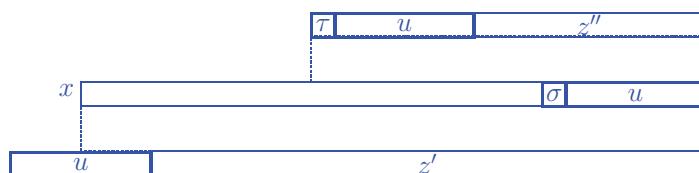
```

BM(string  $x, y$ ; integer  $m, n$ )
   $pos \leftarrow 0$ 
  while  $pos \leq n - m$  do
     $i \leftarrow m - 1$ 
    while  $i \geq 0$  and  $x[i] = y[pos + i]$  do  $i \leftarrow i - 1$ 
    if  $i = -1$  then
      output(' $x$  occurs in  $y$  at position ',  $pos$ )
       $pos \leftarrow pos + period(x)$ 
    else
       $pos \leftarrow pos + D[i]$ 

```

Suffix displacement

- ★ Displacement function d :
- $$d(u) = \min\{|z| > 0 \mid (x \text{ suffix of } uz) \text{ or } (\tau u z \text{ suffix of } x \text{ and } \tau u \text{ not suffix of } x, \text{ for } \tau \in \Sigma)\}$$
- ★ Displacement table D :
- $$D[i] = d(x[i + 1 \dots m - 1]), \text{ for } i = 0, \dots, m - 1$$

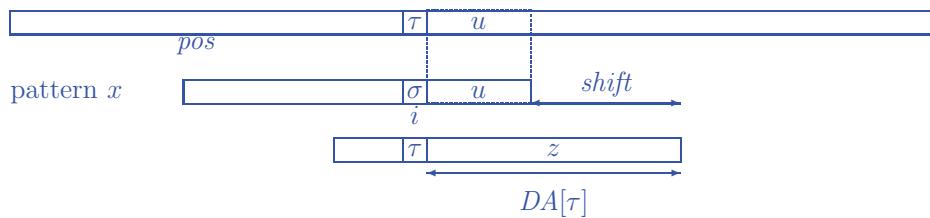


- ★ Note 1: u is a (strict) border of uz''
- ★ Note 2: $|z'|$ is a period of uz' (thus, $|z'|$ is a period of x)

Lemma 1 Table D can be computed in linear time. [see Page 22]

Occurrence shift

text y



- ★ Table DA implements the bad-character rule:

$$DA[\tau] = \min\{|z| > 0 \mid \tau z \text{ suffix of } x\} \cup \{|x|\}$$
 - ★ $shift = DA[\tau] - |u| = DA[\tau] - m + i + 1$
-

COMPUTE_DA(string x ; integer m)

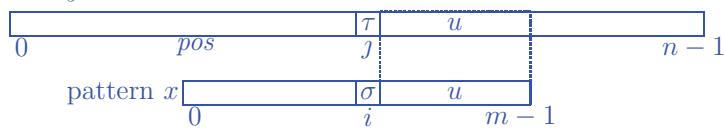
```

for all  $\sigma$  in  $\Sigma$  do
     $DA[\sigma] = m$ 
for  $i \leftarrow 0$  to  $m - 2$  do
     $DA[x[i]] = m - i - 1$ 
return  $DA$ 

```

BM with occurrence shift

text y



- ★ Use of DA in addition to D
-

BM(string x, y ; integer m, n);

```

 $pos \leftarrow 0$ 
while  $pos \leq n - m$  do
     $i \leftarrow m - 1$ 
    while  $i \geq 0$  and  $x[i] = y[pos + i]$  do
         $i \leftarrow i - 1$ 
    if  $i = -1$  then
        output('x occurs in y at position ',  $pos$ )
         $pos \leftarrow pos + period(x)$ 
    else
         $pos \leftarrow pos + \max\{D[i], DA[y[pos + i]] - m + i + 1\}$ 

```

Complexity of BM

★ Preprocessing phase

match shift	$O(m)$
occurrence shift	$O(m + \text{card } \Sigma)$

★ Search phase (finding all occurrences)

running time	$O(n \times m)$
minimum number of comparisons	n/m
maximum number of comparisons	$n \times m$

★ Extra space

for shift functions	$O(m + \text{card } \Sigma)$
can be reduced to	$O(m)$

Symbol comparisons in variants of BM

★ For finding the first occurrence

[Knuth, Morris, Pratt, 1977] $\leq 7 \times n$

[Guibas, Odlysko, 1980] $\leq 4 \times n$

[Cole, 1990] $\leq 3 \times n$

Theorem 1 If $\text{period}(x) > m/2$, BM searching algorithm performs at most $3n - n/m$ symbol comparisons. The bound is tight.

Proof difficult

Symbol comparisons in variants of BM (followed)

- ★ For finding all occurrences

[**Galil, 1979**]

$O(n)$

→ prefix memorization

$O(1)$ extra space

[**Crochemore et alii, 1991**]

$\leq 2 \times n$

→ last-suffix memorization

$O(1)$ extra space

→ **Turbo-BM**

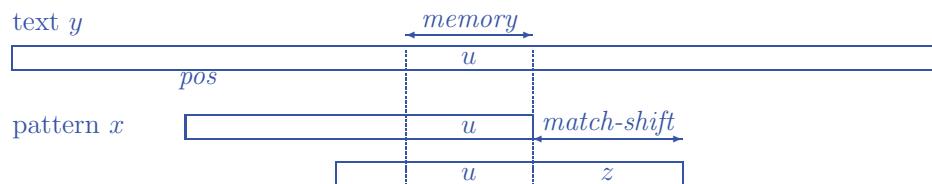
[**Apostolico, Giancarlo, 1986**]

$\leq 1.5 \times n$

→ all-suffix memorization

$O(m)$ extra space

Turbo-BM method



- ★ **Features**

- Stores the last match in case of match shift (*memory*)
- Jumps on *memory*
- Uses **turbo-shifts**

- ★ **Preprocessing**

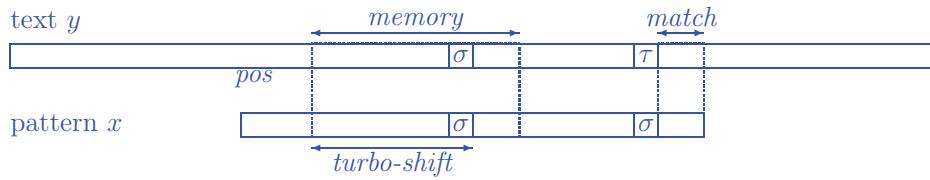
same as BM algorithm

- ★ **Search**

$O(1)$ extra space to store *memory*: (length, right position)

★ Note: *match-shift* is a period of *uz* since *u* is a border of it

Turbo-shift



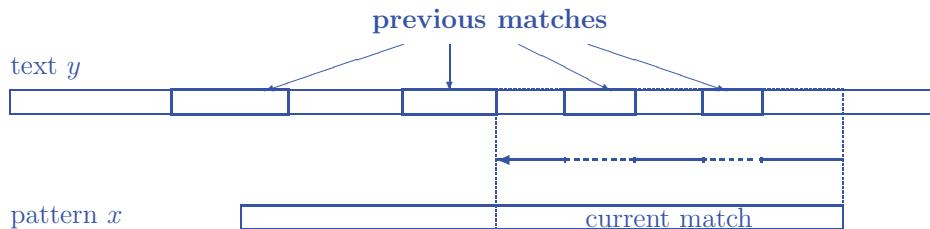
- ★ $turbo\text{-}shift = |memory| - |match|$
- ★ **Use in Turbo-BM:**
 $shift \leftarrow \max\{match\text{-}shift, occ\text{-}shift, turbo\text{-}shift\};$
if ($shift = match\text{-}shift$) **then**
 set memory;
else
 $\{ shift \leftarrow \max\{shift, match+1\}; \text{no memory;} \}$

Tight number of comparisons for Turbo-BM

Theorem 2 *The Turbo-BM searching algorithm runs in time $O(n)$. It makes no more than $2n$ symbol comparisons.*

Proof difficult

AG method



★ **Features**

- **Stores** all previous matches (suffixes of pattern)

- Uses the table Suf

$Suf[i] = \text{longest suffix of } x \text{ ending at } i \text{ in } x$



★ **Extra preprocessing**

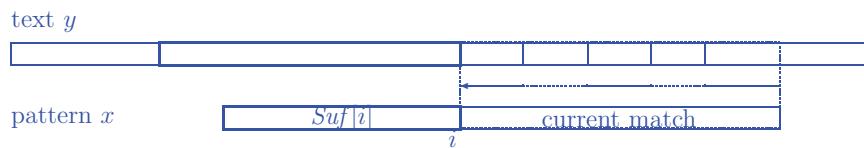
computing Suf on the pattern: $O(m)$ time and space

★ **Search**

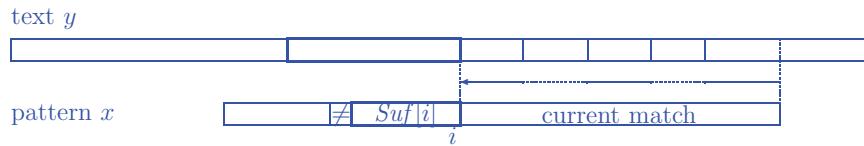
$O(m)$ extra space to store the matches

Rules for shifts in AG

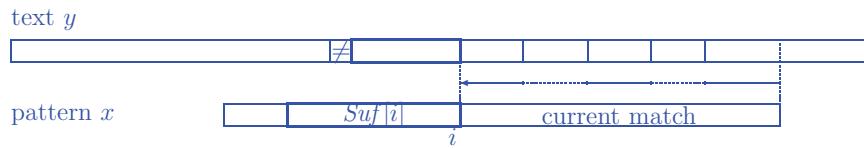
• **match**



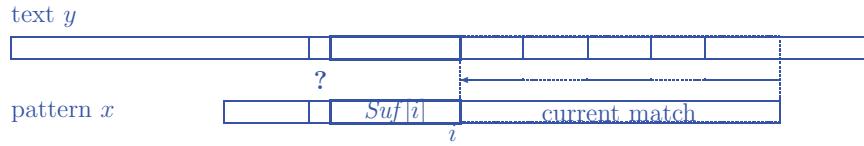
• **mismatch**



• **mismatch**



• **jump**



Tight number of comparisons for AG

Lemma 2 *The AG algorithm makes at most $\frac{n}{2}$ comparisons on text characters previously compared.*

Theorem 3 *The AG searching algorithm runs in time $O(n)$. It makes no more than $1.5n$ symbol comparisons.*

Proof rather difficult

Worst-case example

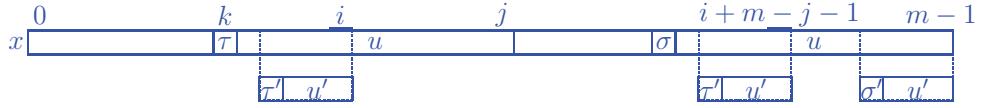
```
a a b a a a b a a b a a a b a a a b ...
a a b a a a b
a a b a a a b
a a b a a a b
a a b a a a b
a a b a a a b
a a b a a a b
a a b a a a b
a a b a a a b
a a b a a a b
```

$$\text{pattern} = a^{m-1}ba^mb, \text{ text} = (a^{m-1}ba^mb)^e$$

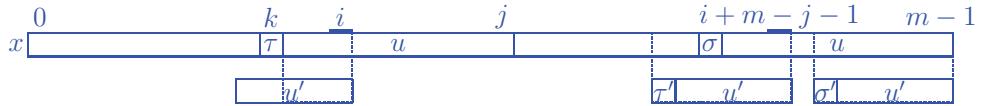
$$\text{number of comparisons} = 2m+1 + (3m+1)e = \frac{3m+1}{2m+1}n - m \approx 1.5n$$

Computing the table of suffixes

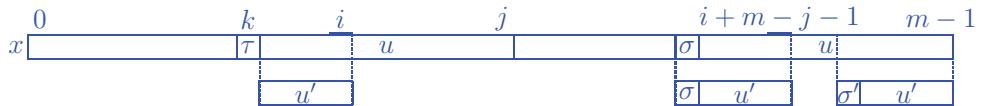
- ★ i given: $Suf[j]$ known for $i < j \leq m - 1$
 $k = \min\{j - Suf[j] \mid i < j < m\}$ (leftmost considered position)
 $\sigma \neq \tau ; \sigma' \neq \tau' ; Suf[i + m - j - 1] = |u'|$



- ★ If $Suf[i + m - j - 1] < i - k$: $Suf[i] = Suf[i + m - j - 1]$



- ★ If $Suf[i + m - j - 1] > i - k$: $Suf[i] = i - k$



- ★ If $Suf[i + m - j - 1] = i - k$: find $Suf[i]$ by scanning **from position k**
Yields a new k and a new j ($= i$)

Extra preprocessing for AG

- ★ Linear-time computation of table Suf

```

COMPUTE_SUF(string x; integer m)
  Suf[m - 1] ← m; k ← m - 1
  for i ← m - 2 downto 0 do
    if i > k and Suf[i + m - j - 1] ≠ i - k then
      Suf[i] ← min{Suf[i + m - j - 1], i - k}
    else
      j ← i; k ← min{i, k};
      while k ≥ 0 and x[k] = x[k + m - j - 1] do
        k ← k - 1;
      Suf[i] ← j - k
  return Suf

```

From Suf to D

i	0	1	2	3	4	5	6	7	8	9	10
$x[i]$	a	b	a	a	a	b	a	b	a	b	a
$Suf[i]$	1	0	3	1	1	0	3	0	5	0	11

- ★ Initializing D using the periods of x only

$Suf[2] = 3 \Rightarrow$ period 8 ; $Suf[0] = 1 \Rightarrow$ period 10 ; period 11

8	8	8	8	8	8	8	8	8	10	10	11
---	---	---	---	---	---	---	---	---	----	----	----

- ★ Accounting for occurrences inside x of its suffixes

$Suf[3] = 1 \Rightarrow D[9] \leq 7$; $Suf[8] = 5 \Rightarrow D[5] \leq 2$

8	8	8	8	8	8	8	8	10	10	11	
				2		8			10		
						8			10		
							4		7		
								6	3		
									1		
$D[i]$	8	8	8	8	8	2	8	4	10	6	1

Computing the displacement table D with the table Suf

-
- ★ Linear-time computation of table D
-

```

COMPUTE_D(string  $x$ ; integer  $m$ ; table  $D$ )
   $j \leftarrow 0;$ 
  for  $i \leftarrow m - 2$  downto  $-1$  do
    if  $i = -1$  or  $Suf[i] = i + 1$  then
      while  $j < m - i - 1$  do
         $D[j] \leftarrow m - i - 1;$ 
         $j \leftarrow j + 1;$ 
    for  $i \leftarrow 0$  to  $m - 2$  do
       $D[m - Suf[i] - 1] \leftarrow m - i - 1;$ 
  return  $D$ 

```
