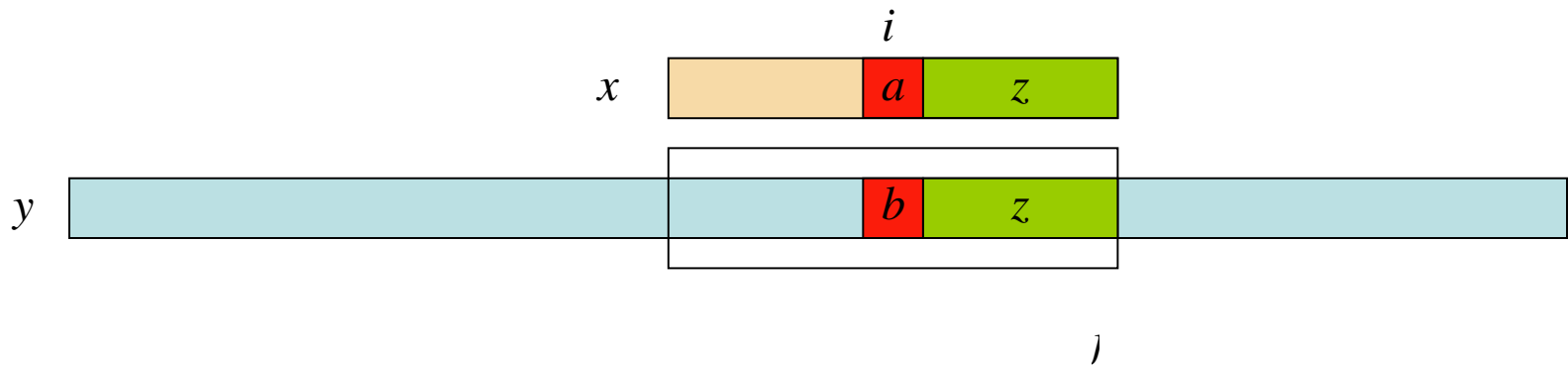


# Algorithme de Boyer-Moore

# Algorithme de Boyer-Moore

Tentative : compare les lettres du mot et les lettres de la fenêtre, de la droite vers la gauche, en commençant par la lettre la plus à droite

# Situation générale



sens des comparaisons

Le suffixe  $z = x[i+1..m-1]$  est égal au facteur  $z = y[j-m+i+2..j]$  et la lettre  $a = x[i]$  est différente de la lettre  $b = y[j-m+i+1]$

# Décalage du bon suffixe

Alors le décalage consiste à choisir le maximum entre aligner

- en face de  $y[j-m+i+2..j] = x[i+1..m-1] = z$  le facteur  $z$  de  $x[0..m-2]$  précédé par une lettre différente de  $a = x[i]$  le plus à droite ou s'il n'existe pas, à aligner le plus long préfixe de  $x$  suffixe de  $z$
- en face de  $y[j-m+i+1] = b$  l'occurrence de  $b$  la plus à droite dans  $x$

# Décalage du bon suffixe

Alors le décalage consiste à choisir le maximum entre aligner

- en face de  $y[j-m+i+2..j] = x[i+1..m-1] = z$  le facteur  $z$  de  $x[0..m-2]$  précédé par une lettre différente de  $a = x[i]$  le plus à droite ou s'il n'existe pas, à aligner le plus long préfixe de  $x$  suffixe de  $z$ 
  - ne dépend que de  $x$
- en face de  $y[j-m+i+1] = b$  l'occurrence de  $b$  la plus à droite dans  $x$ 
  - dépend de  $A$

# Décalage du bon suffixe

On définit deux conditions, pour  $0 \leq i \leq m-1$ ,  $1 \leq d \leq m$  et  $b \in A$

la condition de suffixe  $Cs$

$$\bullet \quad Cs(i,d) = \begin{cases} 0 < d \leq i+1 \text{ et } x[i-d+1..m-d-1] \leq_{suff} x \\ \text{ou} \\ i+1 < d \text{ et } x[0..m-d-1] \leq_{suff} x \end{cases}$$

# Décalage du bon suffixe

On définit la condition  $Co$ ,

pour  $0 \leq i \leq m-1$  et  $1 \leq d \leq m$  par

$$Co(i,d) = \begin{cases} 0 < d \leq i \text{ et } x[i-d] \neq x[i] \\ \text{ou} \\ i < d \end{cases}$$

# Décalage du bon suffixe

Alors la table du bon suffixe est définie comme suit, pour  $0 \leq i \leq m-1$

$$\textit{bon-suff}[i] = \min \{ d \mid Cs(i,d) \text{ et } Co(i,d) \text{ sont satisfaites} \}$$



# Décalage de la dernière occurrence

On définit la table

$$\textit{dern-occ} : A \rightarrow \{ 1, 2, \dots, m \}$$

de la façon suivante

$$\textit{dern-occ}[a] = \min \{ \{ k \mid 1 \leq k \leq m-1 \text{ et } x[m-1-k] = a \} \cup \{m\} \}$$

pour  $a \in A$

algo BM ( $x, m, y, n$ )

$j \leftarrow m-1$

tantque  $j < n$  faire

$i \leftarrow m-1$

tantque  $i \geq 0$  et  $x[i] = y[j-m+1+i]$  faire

$i \leftarrow i-1$

si  $i < 0$  alors

signaler une occurrence de  $x$

$j \leftarrow j + \text{pér}(x)$

sinon

$j \leftarrow j + \max\{\text{bon-suff}[i],$   
 $\text{dern-occ}[y[j-m+1+i]]-m+1+i\}$

# Décalage

## Remarque

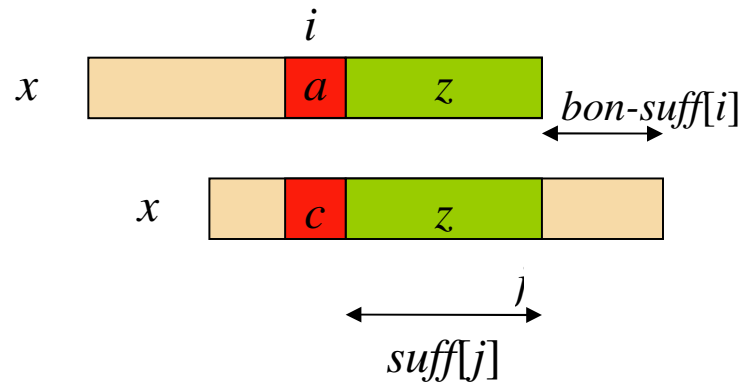
$$\textit{pér}(x) = \textit{bon-suff}[0]$$

$$\forall a \in A$$

# Phase de prétraitement : calcul de la table *dern-occ*

```
algo DERN-OCC( $x, m$ )  
  pour  $a \in A$  faire  
     $\text{dern-occ}[a] \leftarrow m$   
  pour  $i \leftarrow 0$  à  $m-2$  faire  
     $\text{dern-occ}[x[i]] \leftarrow m-1-i$   
  retourner  $\text{dern-occ}$ 
```

# Phase de prétraitement : calcul de la table *bon-suff*



# Calcul de la table *bon-suff*

On définit la table *suff* par :

$$suff[i] = |lsc(x, x[0..i])|$$

où  $lsc(x,y)$  = longueur du plus long suffixe  
commun à  $x$  et  $y$

algo SUFFIXES( $x, m$ )

$g \leftarrow m-1$

$suff[m-1] \leftarrow m$

**pour**  $i \leftarrow m-2$  **à** 0 **faire**

**si**  $i > g$  **et**  $suff[i+m-1-f] \neq i-g$  **alors**

$suff[i] \leftarrow \min\{suff[i+m-1-f], i-g\}$

**sinon**

$g \leftarrow \min\{i, g\}$

$f \leftarrow i$

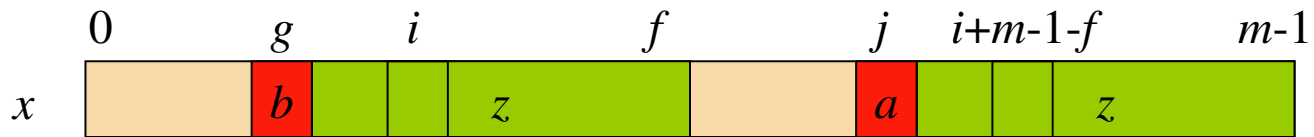
**tantque**  $g \geq 0$  **et**  $x[g] = x[g+m+1-f]$  **faire**

$g \leftarrow g-1$

$suff[i] \leftarrow f-g$

**retourner**  $suff$

# Invariants de l'algorithme SUFFIXES( $x, m$ )



$$z = lsc(x, x[0..f])$$



# Complexité de l'algorithme SUFFIXES( $x, m$ )

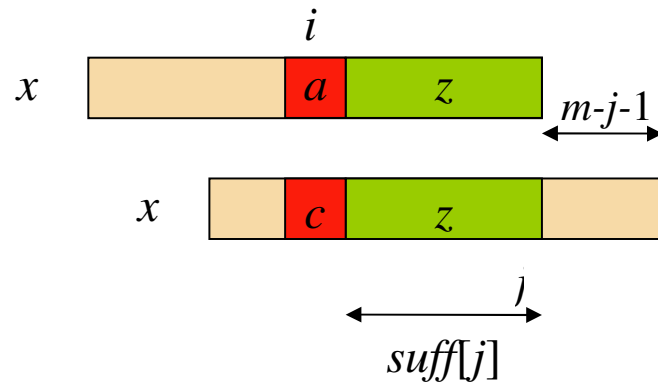
## Proposition 1

L'algorithme SUFFIXES( $x, m$ ) calcule les valeurs de la table *suff* en temps et espace  $O(m)$ .

Preuve (idée) :

- Chaque lettre comparée positivement n'est pas recomparée :  $m$  comparaisons.
- Chaque lettre comparée négativement implique une décrémentement de  $i$  :  $m-1$  comparaisons.
- Au total :  $2m-1$  comparaisons.

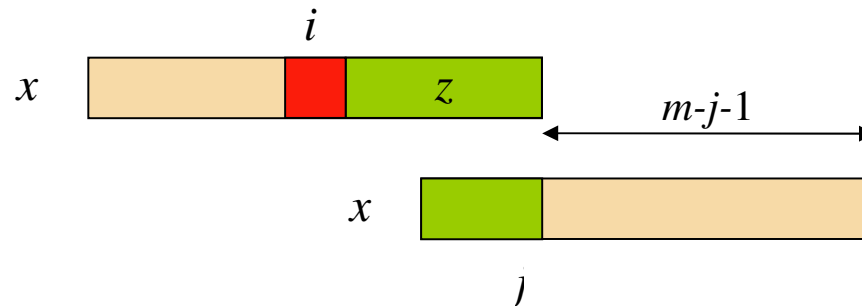
# Phase de prétraitement : calcul de la table *bon-suff*



$$\text{bon-suff}[i] = m-j-1$$

avec  $i = m-1-\text{suff}[j]$

# Phase de prétraitement : calcul de la table *bon-suff*



$$\textit{bon-suff}[i] = m-j-1$$

algo BON-SUFFIXES( $x, m, suff$ )

$i \leftarrow 0$

pour  $j \leftarrow m-2$  à  $-1$  faire

si  $j = -1$  ou  $suff[j] = j+1$  alors

tantque  $i < m-1-j$  faire

$bon-suff[i] \leftarrow m-1-j$

$i \leftarrow i+1$

pour  $j \leftarrow 0$  à  $m-2$  faire

$bon-suff[m-1-suff[j]] \leftarrow m-1-j$

retourner  $bon-suff$

# Complexité de l'algorithme BON-SUFFIXE( $x, m, suff$ )

## Proposition 2

L'algorithme BON-SUFFIXE( $x, m, suff$ ) calcule les valeurs de la table *bon-suff* en temps et espace  $O(m)$ .

Preuve (idée) :

- $j$  varie de  $m-2$  à  $-1$  dans la première boucle **pour** ;
- $i$  varie de  $0$  à  $m-1-\min\{j\}=m$  dans la première boucle **pour** ;
- $j$  varie de  $0$  à  $m-2$  dans la deuxième boucle **pour**.

# Complexité de l'algorithme de Boyer-Moore

## Théorème 3 [Cole 1994]

Lors de la localisation d'un mot  $x$  de longueur  $m$  tel que  $pér(x) > m/3$  dans un texte  $y$  de longueur  $n$ , l'algorithme  $BM(x, m, y, n)$  effectue moins de  $4n$  comparaisons entre des lettres de  $x$  et des lettres de  $y$ .

# Complexité de l'algorithme de Boyer-Moore

## Théorème 4 [Cole 1994]

Lors de la localisation d'un mot  $x$  de longueur  $m$  tel que  $pér(x) > m/2$  dans un texte  $y$  de longueur  $n$ , l'algorithme  $BM(x, m, y, n)$  effectue moins de  $3n$  comparaisons entre des lettres de  $x$  et des lettres de  $y$ .

# Améliorations

- Algorithme Turbo-BM : au pire  $2n$  comparaisons entre des lettres du mot et des lettres du texte
- Algorithme Apostolico-Giancarlo : au pire  $3n/2$  comparaisons entre des lettres du mot et des lettres du texte