## **Sous-mots**

mot = suite sur l'alphabet *A* sous-mot = sous-suite



$$u = u_0 u_1 \dots u_{k-1}$$
 sous-mot de  $x = x_0 x_1 \dots x_{m-1}$   
si  $x = v_0 u_0 v_1 u_1 \dots v_k u_{k-1} v_k$  avec  $v_0, v_1, \dots, v_k \hat{I}$   $A^*$ 

« sous-mot de » : relation d'ordre sur A\*

 $SC(x, y) = \{u \text{ sous-mot de } x \text{ et de } y\}$ 

### Problème 1:

Calculer PLSC(x, y) = max { | u | / u  $\hat{I}$  SC(x, y) }

## Problème 2

Déterminer un mot  $u \hat{I}$  SC(x, y) tel que | u | = PLSC(x, y)

#### **Distance**

UMLV ©

Comparaison de mots par distance

$$d(x, y) = |x| + |y| - 2.PLCS(x, y)$$

#### d est une distance

$$-d(x, y)^3 0$$

$$-d(x, y) = 0 \operatorname{ssi} x = y$$

$$-d(x, y) = d(y, x)$$

- inégalité triangulaire ?

sous-suite x subsequence z ubuesque y

$$u = un plcs de x et z$$

$$v = \text{un plcs de } z \text{ et } y$$

k = nombre d'occurrences de lettres de z communes à u et v

$$d(x, z) + d(z, y) = (|x| + |z| - 2.|u|) + (|z| + |y| - 2.|v|)$$

$$= |x| + |y| + (|z| - |u|) + (|z| - |v|) - |u| - |v|$$

$$^{3} |x| + |y| + (|v| - k) + (|u| - k) - |u| - |v|$$

$$= |x| + |y| - 2k$$

$$^{3} d(x, y)$$

## diff

#### **Comparaison de fichiers**

lettre = ligne

> cat A
Belle Marquise,
vos beaux yeux

me font mourir d'amour

> cat B

D'amour mourir me font, Belle Marquise, vos beaux yeux

> diff A B

0 a 1

> D'amour mourir me font,

3 d 3

< me font mourir d'amour

### **Applications**

gestion de versions compression par stockage de différences restauration avec "ed", "sed", ...

# **Alignement**

### Comparaisons de séquences moléculaires lettre = nucléotide



Transformation de x en y par : insertions, suppressions, remplacements Chaque opération possède un coût

#### Problème:

calculer le coût minimal de la transformation (distance d'alignement) un alignement correspondant

#### Distance d'édition

# Calcul de PLSC

# Par énumération de sous-mots temps exponentiel

# Par "programmation dynamique" temps $O(|x| \times |y|)$

Par automate même temps

# Propriété

$$X = X_0 X_1 ... X_{m-1}$$
  $y = y_0 y_1 ... y_{n-1}$   
 $u = u_0 u_1 ... u_{k-1}$  un plcs de  $x$  et  $y$ 

## **Proposition**

$$x_{m-1} = y_{n-1} P \quad u_{k-1} = x_{m-1} = y_{n-1}$$
  
et  $u_0 u_1 ... u_{k-2}$  plsc de  $x_0 x_1 ... x_{m-2}$  et  $y_0 y_1 ... y_{n-2}$   
 $x_{m-1} \quad y_{n-1}$  et  $u_{k-1} \quad x_{m-1} P \quad u$  plsc de  $x_0 x_1 ... x_{m-2}$  et  $y_0 y_1 ... y_{n-2}$   
 $x_{m-1} \quad y_{n-1}$  et  $u_{k-1} \quad y_{n-1} P \quad u$  plsc de  $x$  et  $y_0 y_1 ... y_{n-2}$ 







abac

# **Algorithme**

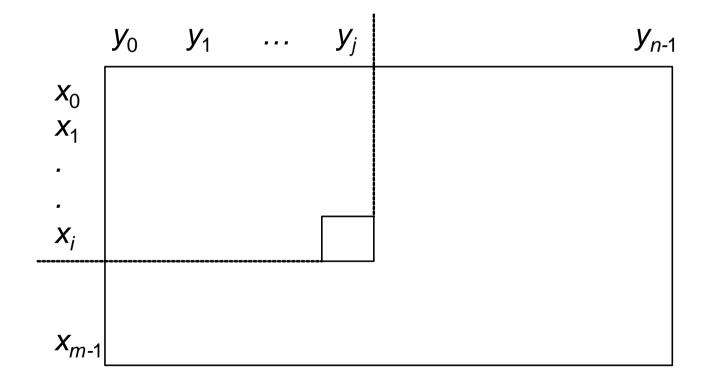
```
entier PLSC(mot x, longueur m, mot y, longueur n){
    si (m = 0 ou n = 0) alors
        retour 0;
    sinon si (x[m-1] = y[n-1]) alors
        retour PLSC(x,m-1,y,n-1)+1
    sinon
    retour max{PLSC(x,m,y,n-1),PLSC(x,m-1,y,n)}
}
```

**Stupidement exponentiel!** 

# **Programmation dynamique**

Problème décomposable en nombre fini de sous-problèmes mais chevauchement des sous-problèmes

Méthode : mémoriser les résultats intermédiaires



## Récurrence

$$L(i, j) = PLSC(x_0x_1...x_i, y_0y_1...y_j)$$
  
 $L(i, j) = 0$  si  $i = -1$  ou  $j = -1$   
 $L(i-1, j-1) + 1$  sinon si  $x_i = y_j$   
 $max(L(i, j-1), L(i-1, j))$  sinon

# Exemple

```
L -1 0 1 2 3 4 5 6 7 8
c b a c b a a b a

-1 0 0 0 0 0 0 0 0 0 0 0

0 a 0 0 1 1 1 1 1 1 1 1

1 b 0 0 1 1 1 2 2 2 2 2 2

2 c 0 1 1 1 2 2 2 2 2 2

3 d 0 1 1 1 2 2 2 2 2 2

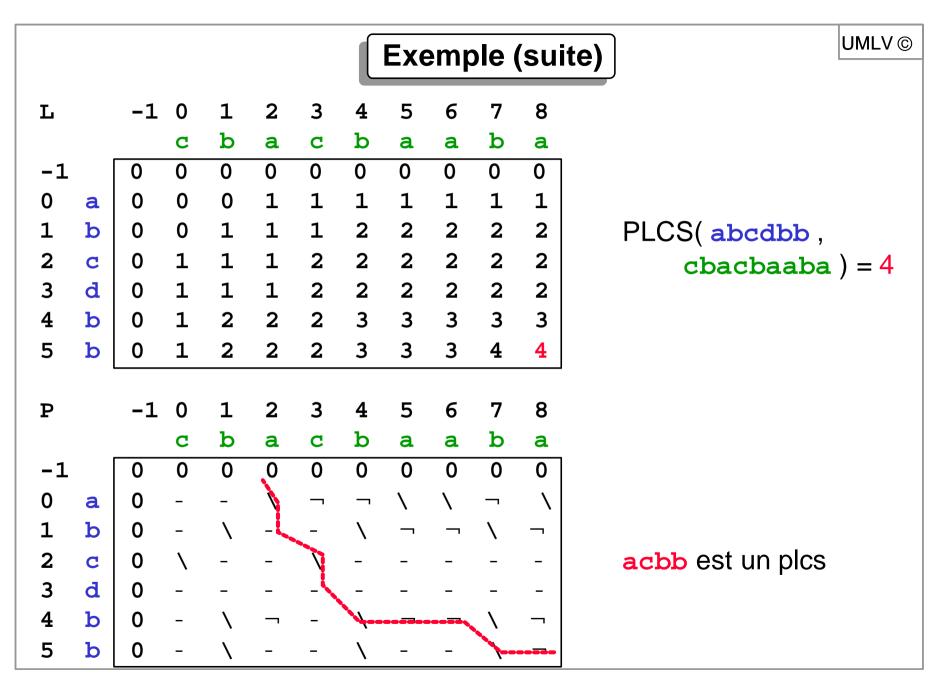
4 b 0 1 2 2 2 3 3 3 3 3

5 b 0 1 2 2 2 3 3 3 3 4 4
```

PLCS(abcdbb,cbacbaaba) = 4

# | Algorithme

```
entier PLSC(mot x, longueur m, mot y, longueur n) {
     pour i \neg -1 à m-1 faire L[i,0] \neg 0;
     pour j \neg -1 à n-1 faire L[0,j] \neg 0;
    pour i - 0 à m-1 faire
          pour j - 0 à n-1 faire
               si x[i] = y[j] alors{
                    L[i,j] - L[i-1,j-1]+1;
                    P[i,j] ¬ '\';
               }sinon si L[i-1,j] 3 L[i,j-1] alors{
                    L[i,j] - L[i-1,j] ;
                   P[i,j] ¬ '-';
               }sinon{
                    L[i,i] \neg L[i,i-1];
                   P[i,j] ¬ '¬';
     retour L[m-1,n-1];
```



## **Variantes**

#### Problème ouvert

calcul de PLSC en temps  $< O(n^2 / \log n)$ 

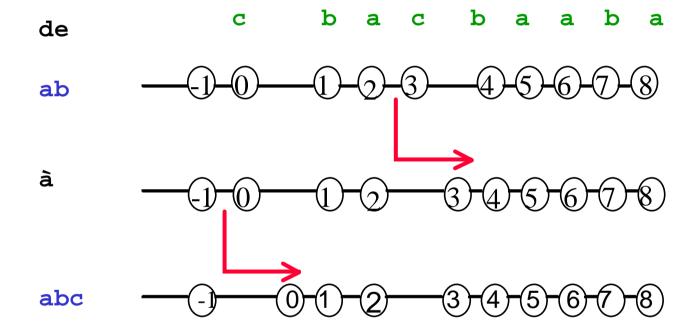
#### **Variantes**

- calcul de PLSC en espace linéaire
   [une ligne de L suffit]
- calcul d'un plsc en espace linéaire
   [algorithme de Hirschberg]
- utilisation des « dominants »
- algorithme rapide en pratique bien que temps maximal = O(n² log n) dans certaines implantations de diff [algorithme de Hunt et Szymanski]

# Algorithme de Hunt & Szymanski

# Calcul de la matrice L par ligne

# Comme un boulier



## **Algorithme**

```
X = X_0 X_1 \dots X_{m-1} Y = Y_0 Y_1 \dots Y_{n-1}
      Pos[a] = {j | y_i = a} pour a \hat{I} A
      Pour i fixé : J[k] = \{ j \mid PLSC(x_0x_1...x_i, y_0y_1...y_i) = k \}
entier PLSC(mot x, longueur m, mot y, longueur n){
      J[-1] - \{-1,0,1,...,n-1\};
      pour k - 0 à n-1 faire J[k] - E;
      pour i - 0 à m-1 faire
            pour chaque p \hat{I} Pos[x[i]] en décroissant faire{
                  k - CLASSE(p);
                  si k = CLASSE(p-1) alors{
                        (J[k],X) \neg PARTAGE(J[k],p);
                        J[k+1] \neg UNION(X,J[k+1]);
      retour CLASSE (n-1);
```

## Type abstrait "PARTAGE/FUSION"

#### Ensemble de base

```
suites (E_0, E_1, ..., E_n) d'ensembles disjoints telles que E_0 \cup E_1 \cup ... \cup E_n = \{0, 1, ..., n\}
```

## **Opérations**

Initialisation

**CLASSE**: élément ® ensemble

CLASSE(p) = k tel que  $p \hat{I} E_k$ 

PARTAGE : élément x ensemble ® ensemble x ensemble

PARTAGE(E, p) = (S, T)

 $S = \{ q\hat{I} \mid E \mid q$ 

**UNION**: ensemble x ensemble ® ensemble

## Implémentations possibles

listes, arbres, ...

Si les  $E_k$  sont des intervalles : B-arbres, 2-3-arbres, ...

## Temps de calcul de PLSC

```
2-3-arbres pour les J_k

initialisation O(n)

CLASSE(p) O(\log n)

PARTAGE(J, p) O(\log n)

UNION(X, J) O(\log n)

Temps: si Pos listes pré-calculées sur y

O(\log n \cdot \sum (|Pos(x_i)| i = 0, 1, ..., n-1)

= O(\log n \cdot \text{card}\{(i, j) / x_i = y_j\}|)

= O(n \cdot m \cdot \log n)

En pratique, sur fichiers: O(n \cdot \log n)
```

#### Pré-calcul de Pos

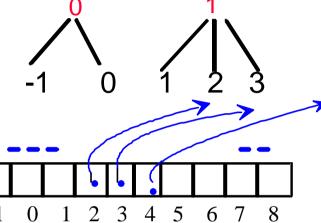
```
|y| = n \otimes O(n)
si fichiers de n lignes, par hachage O(n) en moyenne
```

# **Utilisation des 2-3 arbres**

UMLV ©

Numéro de classe à la racine

Pointeurs sur feuilles



CLASSE : par liens « père »

PARTAGE: en descendant, duplication, et restructuration

UNION : rattachement au bon niveau, et restructuration

