

Le Texte est un format omniprésent, qu'on voit partout :

- ① Comme format d'une grande partie de la culture humaine.
- ② Comme interface homme-machine (fichiers sources / UNIX)
- ③ Comme langage de la vie : chaîne d'ADN ou de protéines.

## I] Représentation informatique

### 1] Code Ascii (usage ②)

Sur 7 bit (convention 8) représente l'alphabet latin, les chiffres, la ponctuation et de symboles divers (retour à la ligne)

Avantage : un texte peut être vu comme un tableau de caractères.  
Représentation linéaire.

### 2] UTF-8 (usage ①)

Sur 8, 16 ou 32 bits représente une majorité de glyphes humains.

Avantage : compatible avec l'ASCII

Inconvénient : incompatible avec certains algorithmes

Solution : UTF-32, tableau de caractères.

Remarque : Devraient on noter  $\Sigma$  l'ensemble des lettres.

### 3] Compression

Tous les symboles n'ont pas la même fréquence d'apparition.

En français par exemple e ? s ? a ? i ? t ? n ? ...

On peut améliorer le codage en utilisant des tailles variables

Def : Un arbre binaire complet associé pour chaque nœud soit 0 soit 1 fils.

Un ABC permet de coder des lettres sans ambiguïté :



Def : Soit un codage donné par un ABC. Soit  $\lambda : \Sigma \rightarrow \mathbb{N}$  qui à chaque lettre associe la longueur de son codage. Soit  $\psi$  la fréquence des lettres.

Le coût d'un ABC pour une fréquence  $\psi$  est donné par

$$\sum_{\sigma \in \Sigma} \psi(\sigma) \lambda(\sigma)$$

NOM :

Prénom :

N° - Date - Titre de la leçon : 20 0

gloster

Algo (Huffman) : Tous les  $A \in \Sigma$  de coût  $l(A)$  minimal associé à un texte

### II) Dictionnaire

Structure de données efficace pour chercher un caractère de manière binaire de taille  $n$

Représentation binaire :  
 $a \sqrt{ab}$   
 } la mota mot de taille  $n$   
 aa ab bb } jusqu'à  $2^n$  fils pour chaque mot.

Dd Une arbre binaire de recherche sur un arbre où chaque nœud représente un mot tel que



alors  $x \leq u \leq y$   
 $v \leq x$   
 $w \leq y$

Inconvenients : Double pile des (pages)

Parvenir par les lists :  $O(n)$  avec un nombre de mots.



Dd Arbre rouge noir ou aussi un arbre où que nœuds : rouge ou noir

- Noir et noir
- chaque fils d'un nœud rouge et noir
- chaque branche à la même nombre de nœuds noirs.

Lemme de longueur des branches d'un AEN et  $O(\log n)$

Algo Insertion, suppression et recherche sur un AEN  $O(\log n)$

### III) Recherche de Motif

Soit  $T C_1 \dots C_n$  un texte

Soit  $M C_1 \dots C_m$  un mot

Peut on trouver  $i$  tel que  $\forall a \in \{C_1, \dots, C_m\} M[a] = T[C_i + a]$

#### 1) Recherche naïve

On fait une boucle sur  $i$  sur une boucle sur  $k$  :  
 complexité :  $O(m \cdot n)$

#### 2) Rabin - Karp

Principe : calcul d'un hash de  $M[C_1 \dots C_m] = h(M[C_1 \dots C_m])$

\* on pose  $h_i = h(T[C_i \dots C_i + m])$

\*  $h_i$  est calcul en  $O(n)$  à partir de  $h_{i-1}$ .

Algo On compare  $h_i$  avec  $h$  et si

$h_i = h$  : on vérifie que  $\forall k \leq m \quad M[C_k] = T[C_i + k]$

$h_{i+1} = d(h_i - T[C_i + 1]) + T[C_i + m + 1]$  [9]

Complexité :  $O(m)$

Complexité :  $O(m \cdot n)$

Complexité :  $O(m + n)$

3) Automates finis

Version naïve pour les "babas"



Comment le calculer directement :

Def : un bon de la se en mot à la fois  
Auffine se profile de la

Ex : {ε, a, abab} sont les bons de abababa

Principe Soit M = ababababa ∈

Si  $MC_{i+1} = TC_{i+1}, \dots, (i+1)$

ou  $MC_i \neq TC_i + \delta$

alors  $MC_{i+1} = TC_i + 7 - 3, (i+1)$

car 3 se la taille de plus grand bon de  $M C_{i+1}$

Coût partiellement  $O(m \cdot |Z|)$  |  $\underbrace{a \ b \ a \ a \ b \ a \ c}_{\substack{\text{a} \\ \text{b} \\ \text{c}}}$   
Coût séquentiel  $O(n)$

3) Knuth Morris Rabin

On évite l'attente de la partie précédente dans le calcul directement.

Coût partiellement  $O(n)$

Coût séquentiel  $O(n)$

IV) Compilateur et Automate (autre grande)

1) Les langages réguliers et autre grammaire par exemple se vérifie la syntaxe d'un autre.

$C = 1+2+3+\dots+n+0$  : Date valid  $\rightarrow CC/CC/CCCC$

2) Les grammaires permettent d'identifier "compilateur"

un texte de la transformation à autre :

"(1+2) x 5"  $\rightarrow$   $\sqrt[2]{1 \times 5}$

V) - Organisation de la

1) Distance d'édition

Algo calcul de nombre d'étape minimale nécessaire pour passer d'un mot  $a$ , à un mot  $b$  en

opérations suivantes :

- Suppression
- insertion
- substitution

Coût  $O(|a| + |b|)$

2) Plus longue sous séquence commune

Algo dans un tableau d'un mot  $a_1 \dots a_n$  et un mot  $b_1 \dots b_m$  tel que il existe  $\sigma : A \rightarrow B$  strict. maximale vérifier  $A[i] = B[\sigma(i)] = b_k$

Coût :  $O(|a| \cdot |b|)$  pour trouver la PLSSC de ces deux mots