

### Devoir 1. Compression d'une séquence d'entiers

Le sujet de ce travail pratique est le développement d'un outil de compression<sup>1</sup>. L'entrée de la procédure de compression est une séquence de nombres naturels  $\{x_i: i = 0, 1, 2, \dots\}$  (abstraction d'un fichier d'octets, par exemple). Notre procédure encode l'entrée à l'aide d'une transformation qui produit une séquence  $\{y_i: i = 0, 1, 2, \dots\}$  de *petits* entiers quand l'entrée est compressible. Ensuite, on encode les entiers en une séquence de chaînes de bits  $\{z_i: i = 0, 1, 2, \dots\}$ , en utilisant moins de bits pour les petits entiers. Il est important que les deux étapes sont inversibles instantanément : on peut calculer  $x_i \leftrightarrow y_i$  et  $y_i \leftrightarrow z_i$  à chaque  $i = 0, 1, 2, \dots$ . En conséquence, la procédure est une méthode d'encodage déchiffable.

#### 1.1 Move-to-front

La transformation **move-to-front** (MTF) sert à encoder une séquence d'entiers par une autre. Pour encoder la séquence  $x_0, x_1, x_2, \dots$  d'entiers non-négatifs, on maintient une permutation  $\pi$ . Initialement, on a identité :  $\pi[n] = n$  pour tout  $n$ .

Pour transmettre le prochain entier  $n$ , on identifie son unique position  $i$  dans la permutation ( $\pi[i] = n$ ), et on émet  $i$ . En même temps, on performe une rotation circulaire de  $\pi[0..i]$  qui place  $n$  en position 0. Dans la permutation résultante  $\pi'$ , on a  $\pi'[0] = \pi[i]$ ,  $\pi'[j] = \pi[j-1]$  pour  $j = 1, 2, \dots, i$  et  $\pi'[j] = \pi[j]$  pour tout  $j > i$ .

Pour décoder, on fait la même rotation sauf que l'argument est l'indice  $i$ , et on doit retourner  $\pi[i]$  (avant la rotation).

*Exemple :* la transformation MTF de la séquence  $(1, 2, 5, 2, 1, 4, 3, 0)$  donne  $(1, 2, 5, 1, 2, 5, 5, 5)$  :

$i$	$x_i$	$\text{encodeMTF}(x_i)$	$\pi$ (après rotation)
0	1	1	$(1 \ 0 \ 2 \dots \infty)$
1	2	2	$(2 \ 1 \ 0 \ 3 \dots \infty)$
2	5	5	$(5 \ 2 \ 1 \ 0 \ 3 \ 4 \ 6 \dots \infty)$
3	2	1	$(2 \ 5 \ 1 \ 0 \ 3 \ 4 \ 6 \dots \infty)$
4	1	2	$(1 \ 2 \ 5 \ 0 \ 3 \ 4 \ 6 \dots \infty)$
5	4	5	$(4 \ 1 \ 2 \ 5 \ 0 \ 3 \ 6 \dots \infty)$
6	3	5	$(3 \ 4 \ 1 \ 2 \ 5 \ 0 \ 6 \dots \infty)$
7	0	5	$(0 \ 3 \ 4 \ 1 \ 2 \ 5 \ 6 \dots \infty)$

<sup>1</sup> Voir §5.5 de Sedgewick et Wayne pour une discussion détaillée de méthodes de compression.

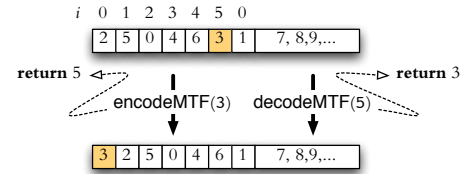


FIG. 1: La transformation MTF

## 1.2 Comment encoder les entiers ?

On veut transmettre une séquence de nombres naturels  $y_0, y_1, y_2, \dots$  sur une ligne de communication. La difficulté est qu'on ne peut utiliser que deux symboles, comme en code Morse : **0** et **1**. (Il s'agit de l'abstraction d'un fichier binaire.)

On définit donc un code  $\rho(y)$  (séquence de symboles de **0** et **1**) pour tout  $y$ , et on transmet les codes  $\rho(y_0), \rho(y_1), \dots$ , l'un après l'autre. En plus d'être *déchiffrable* (on peut récupérer la séquence  $y_0, y_1, \dots$  sans ambiguïté à partir de la chaîne concaténée),  $\rho$  devrait être *instantané* : on veut déchiffrer  $y_i$  dès qu'on reçoit tous les bits de  $\rho(y_i)$ . Et, enfin, on désire aussi que le code soit économique.



*Encodage unaire.* L'exemple le plus simple est le *code unaire*  $\alpha(n)$  :

$$\alpha(n) = \underbrace{\mathbf{1} \mathbf{1} \dots \mathbf{1}}_{n \text{ fois}} \mathbf{0} \quad (1.1)$$

Clairement, il est possible de décoder une séquence de codes concaténés  $\alpha(y_0), \alpha(y_1), \dots$ , car il suffit de compter les **1**s jusqu'au premier **0**. Mais ce code n'est pas économique de tout ...

*Encodage binaire.* L'encodage binaire standard  $\beta(n)$  est beaucoup plus efficace, parce qu'il utilise seulement  $d$  symboles où  $d$  est le plus petit nombre naturel tel que  $n < 2^d$ . Pour  $n = \sum_{i=0}^{d-1} b_i \cdot 2^i$  (avec  $b_i \in \{0, 1\}$  et  $b_{d-1} = 1$ ) on transmet la séquence  $b_{d-1}, b_{d-2}, \dots, b_0$ , en commençant par le bit de poids fort. En particulier  $\beta(0) = \lambda$  (chaîne vide),  $\beta(1) = \mathbf{1}$ ,  $\beta(2) = \mathbf{1} \mathbf{0}$ , etc. On définit le code par récurrence

$$\beta(n) = \begin{cases} \lambda & \text{vide si } n = 0 \\ \mathbf{1} & \{n = 1\} \\ \beta(\frac{n}{2}) \oplus \mathbf{0} & \{n = 2, 4, 6, 8, \dots\} \\ \beta(\frac{n-1}{2}) \oplus \mathbf{1} & \{n = 3, 5, 7, \dots\} \end{cases}$$

où  $\oplus$  dénote la concaténation. Hélas,  $\beta$  n'est pas déchiffrable :  $(1, 2, 1)$  et  $(6, 1)$  mènent à la même séquence de bits **1 1 0 1**.

*Encodage  $\gamma$ .* On peut construire un code instantané en préfixant  $\beta(n)$  par sa longueur encodé en unaire : si  $\beta(n)$  comprend  $d$  symboles, alors on écrit  $d$  fois **1**, suivi par un seul **0**, et l'encodage binaire sur  $d$  bits.

$$\gamma(n) = \alpha(|\beta(n)|) \oplus \beta(n), \quad (1.2)$$

où  $\oplus$  dénote la concaténation et  $|\dots|$  la longueur en bits. Le code  $\gamma$  est instantanément déchiffrable : compter les **1**s au début ( $= d$ ), jusqu'au premier **0**, et décoder les  $d$  bits suivants contenant  $\beta(y_i)$ .



### 1.3 Travail à faire

#### Développement algorithmique (30 points)

*Move-to-front.* ► Développer une structure de données<sup>2</sup> qui implante le type abstrait avec opérations `encodeMTF` et `decodeMTF`. La structure est initialisée avec le maximum  $M$  tel que tout entier à l'entrée sera entre 0 et  $(M - 1)$ . Montrez l'implantation des opérations avec tous les détails.

<sup>2</sup> **Indice:** Pour stocker la permutation actuelle  $\pi$ , on peut utiliser une liste chaînée ou un tableau. Initialiser avec la permutation d'identité.

*Encodage universel.* On écrit les chaînes de bits à l'aide d'une interface de sortie : la méthode `writeBit( $b$ )` écrit le bit  $b$  dans le fichier. ► Donner un algorithme *récurif* `encodeOmega( $n$ )` pour écrire<sup>3</sup> l'encodage  $\omega_1(n)$ . (Noter que l'algorithme doit aussi calculer  $|\beta(n)|$  et écrire  $\beta(n)$ .)

<sup>3</sup> **Indice:** Appelez `writeBit(0)` pour `0` et `writeBit(1)` pour `1`.

#### Implémentation Java (30 points)

► Implémenter un outil de compression en Java qui emploie la transformation MTF et l'encodage universel  $\omega$ .

- ★ L'exécutable s'appelle `MTFOmega`, et il prend 1 argument de la ligne de commande qui est le nom d'un fichier à lire
- ★ Le fichier d'entrée est ouvert par `FileReader` (comme un flot de caractères), et lu à l'aide de la méthode `read()` (qui retourne une valeur de 16 bits, entre 0 et 65535)
- ★ Chaque entier  $x_i$  lu est écrit comme `encodeOmega(encodeMTF( $x_i$ ))`
- ★ La sortie est un fichier binaire, écrit sur `System.out` (on écrit un octet par la méthode `write`). Ici, il faut encoder la séquence de bits en une séquence d'octets. Utilisez le principe (gros-boutiste) (*big-endian*), et implémentez<sup>4</sup> une structure de buffer/tampon pour empaqueter les bits.

Exemple d'application à la ligne de commande (code dans `tp1.jar`) :

```
% java -cp tp1.jar MTFOmega fichier.txt > fichier.zzz
```

► Tester l'efficacité de la méthode développée sur quelques exemples : choisir quelques fichiers (5–10) sur le site <https://introcs.cs.princeton.edu/java/data/>, et comparer la taille du fichier comprimé avec celle d'un autre logiciel de compression au choix (par exemple, `gzip`). Compiler les résultats dans un tableau : pour chaque jeu de données, donner en octets (1) la taille originale, (2) la taille comprimée selon `MTFOmega`, et (3) la taille comprimée selon l'autre outil.

<sup>4</sup> Vous avez le droit d'adapter le code de `BinaryOut.java` de Sedgewick et Wayne, avec mention de la source.