

**FILIERE MP : ENS (PARIS) – ENS LYON – ENS CACHAN**

**PAGE DE GARDE DU RAPPORT DE TIPE 2012**

NOM : DELPEUCH

Prénoms : Antonin, Louis

Lycée : Louis le Grand

Classe : MP\*3

Ville : Paris

**Concours auxquels vous êtes admissible dans la banque inter-ENS :**

(Mettre une croix très visible dans la ou les case(s) vous concernant)

ENS Cachan	MP - option MP	<input type="checkbox"/>	MP - option MPI	<input type="checkbox"/>
	Informatique	<input checked="" type="checkbox"/>		
ENS Lyon	MP - option MP	<input type="checkbox"/>	MP - option MPI	<input type="checkbox"/>
	Informatique - option M	<input checked="" type="checkbox"/>	Informatique - option P	<input type="checkbox"/>
ENS (Paris)	MP - option MP	<input type="checkbox"/>	MP - option MPI	<input type="checkbox"/>
	Informatique	<input checked="" type="checkbox"/>		

**Matière dominante du TIPE :** Informatique  
(mathématiques, informatique ou physique)

**Titre du TIPE :** Complexité de Kolmogorov et compression de Huffman

**Nombre de pages** (à porter dans les cases ci-dessous) :

Texte 

T	5
---	---

 Illustrations 

I	10
---	----

 Bibliographie 

B	1
---	---

**Résumé imprimé (6 lignes) :**

On définit la complexité de Kolmogorov, qui mesure l'information contenue dans un mot, et on en étudie quelques propriétés. Cette complexité n'étant pas calculable, on cherche à la majorer. On présente pour cela l'algorithme de compression de Huffman. Pour justifier les objectifs de cet algorithme, on introduit l'inégalité de Kraft, dont les applications à la théorie algorithmique de l'information sont par ailleurs fructueuses. Enfin, on illustre ces concepts par l'étude de la compression de certains préfixes du mot de Thue-Morse.

A. PARIS....., le 05/06/2012  
Signature du (de la) candidat(e)

Signature du professeur responsable de  
la classe préparatoire dans la discipline

J. CAZOR

Cachet de  
l'établissement



# Complexité de Kolmogorov et compression de Huffman

Antonin Delpeuch

14 juin 2012

## Introduction

Un paradoxe, dit de Berry, est de définir un entier par « le plus petit entier naturel non descriptible en français par une expression de dix-huit mots ou moins. »

Peut-on définir la longueur de la plus courte description d'un objet? Quelles en sont les propriétés? Ces idées doivent bien sûr être formalisées. Pour cela, utilise la notion de fonction *calculable* (ou *semi-calculable* pour une fonction partielle), de  $\Sigma^*$  dans lui-même, où  $\Sigma = \{0, 1\}$ . On note  $Dom(f)$  le domaine de  $f$ . Si  $x \in Dom(f)$ , on note  $f(x) \downarrow$ , et  $f(x) \uparrow$  sinon.

Dans tout ce qui suit, les mots (et donc les programmes) sont écrits sur l'alphabet binaire, sauf mention du contraire.

## 1 Complexité de Kolmogorov

### 1.1 Définition

Il existe une fonction calculable  $\mathcal{U}$  telle que pour toute fonction semi-calculable  $f$  il existe  $\rho_f \in \Sigma^*$  tel que pour tout  $x \in Dom(f)$ ,  $\mathcal{U}(\rho_f \cdot x) = f(x)$ .  $\mathcal{U}$  agit en fait comme un interpréteur universel, et  $\rho_f$  représente un code source de  $f$  pour cet interpréteur.

**Définition 1.** La complexité de Kolmogorov d'un mot  $\omega$  de  $\Sigma^*$ , notée  $C(\omega)$ , est la taille du plus petit programme qui génère  $\omega$ .

$$\begin{aligned} C : \Sigma^* &\longrightarrow \mathbb{N} \\ \omega &\longmapsto \min\{|s|, s \in \Sigma^*, \mathcal{U}(s) = \omega\} \end{aligned}$$

On appelle *programme minimal* d'un mot  $\omega$  un programme de taille  $C(\omega)$  qui génère  $\omega$ . On note  $\hat{\omega}$  l'un d'entre eux.

$C$  est bien définie car pour tout  $x \in \Sigma^*$ , il existe un programme qui génère  $x$  : l'identité de  $\Sigma^*$  est calculable, donc  $\rho_{id} \cdot x$  convient.

### 1.2 Rareté des mots à faible complexité

**Propriété 1.** Il y a au plus  $2^{n-p+1} - 1$  mots  $w$  de taille  $n$  tels que  $C(w) \leq n - p$ .

*Démonstration.* Soit  $n \in \mathbb{N}$ . Considérons l'ensemble  $F$  des mots de  $\Sigma^n$  dont un programme minimal est de taille inférieure à  $n - p$ .

$F$  s'injecte dans l'ensemble des programmes de taille inférieure ou égale à  $n - p$ , qui est de cardinal  $2^{n-p+1} - 1$ . En effet,  $\mathcal{U}$  associe à tout programme un unique mot. Donc  $\text{Card}(F) \leq 2^{n-p+1} - 1$ .  $\square$

Les mots qui peuvent être compressés de manière significative sont donc rares. Cela renforce l'idée qu'un mot *pris au hasard* est un mot qui a une complexité de Kolmogorov forte. On peut d'ailleurs utiliser le même argument (avec  $p = 1$ ) pour montrer la proposition suivante.

**Propriété 2.** Pour tout  $n$ , il existe  $\omega \in \Sigma^*$  tel que  $C(\omega) \geq n$ .

### 1.3 Des inégalités vérifiées par $C$

**Propriété 3.**  $C(\omega) \leq |\omega| + O(1)$

On a obtenu l'inégalité précédente en justifiant l'existence d'un programme générateur dans la définition de la complexité.

**Propriété 4.**  $C(\omega\omega) \leq C(\omega) + O(1)$

*Démonstration.* La fonction  $\omega \mapsto \omega\omega$  est calculable. On la compose avec une fonction générant  $\omega$ , le codage en résultant ne diffère que d'une constante indépendante de  $\omega$ .  $\square$

**Propriété 5.**  $C(\widehat{\omega}) = C(\omega, C(\omega)) + O(1)$

*Démonstration.* Si on connaît un  $\widehat{\omega}$ , on obtient  $C(\omega) = |\widehat{\omega}|$  et  $\omega = \mathcal{U}(\widehat{\omega})$ . Donc  $C(\omega, C(\omega)) \leq C(\widehat{\omega}) + O(1)$ . Réciproquement, si on connaît  $\omega$  et  $C(\omega)$ , on obtient un  $\widehat{\omega}$  en exécutant en parallèle tous les programmes de longueur  $C(\omega)$  (il y en a un nombre fini) et en choisissant le premier qui s'arrête en renvoyant  $\omega$ . D'où la deuxième inégalité.  $\square$

Enfin, l'écriture binaire des entiers naturels donne le résultat suivant :

**Propriété 6.**  $C(n) \leq \log n + O(1)$

## 2 Inégalité de Kraft

### 2.1 Énoncé

Une idée naturelle pour décrire un mot de manière concise est de remplacer les facteurs les plus fréquents par des codes plus courts. Fixons donc  $n$  mots  $m_i$ , qu'on veut représenter chacun par un code binaire de longueur  $c_i$ . L'inégalité de Kraft donne une condition nécessaire et suffisante sur les  $c_i$  pour qu'il existe des  $x_i \in \Sigma^*$  tels que  $|x_i| = c_i$  et tels qu'il existe au plus une manière de décoder tout mot, c'est à dire que la fonction

$$\begin{aligned} \phi : \{m_1, \dots, m_n\}^* &\longrightarrow \Sigma^* \\ m_{i_1} \dots m_{i_p} &\longmapsto x_{i_1} \dots x_{i_p} \end{aligned}$$

est injective.

**Théorème 1.** (*Inégalité de Kraft, ou de McMillan*) Il existe un code  $(x_i)_{i \geq 0}$  uniquement décodable et vérifiant  $|x_i| = c_i$  pour tout  $i$  si et seulement si

$$\sum_i 2^{-c_i} \leq 1$$

Dans le cas général d'un alphabet de cardinal  $D \geq 1$ , le théorème reste vrai en remplaçant 2 par  $D$ .

*Démonstration.* On montre que l'inégalité est nécessaire. Pour l'autre sens, la preuve est celle du théorème 2.

On suppose pour commencer que la famille  $(c_i)$  est finie :  $1 \leq i \leq N$ . Soit  $u_i = |\{1 \leq j \leq N, c_j = i\}|$  le nombre de mots de taille  $i$ . On pose  $f(x) = \sum_{i \geq 0} u_i x^i$ .

$$f(x)^n = \sum_{k \geq 0} \left( \sum_{c_{i_1} + \dots + c_{i_n} = k} u_{i_1} \dots u_{i_n} \right) x^k$$

Or le coefficient devant  $x^k$  est le nombre de découpages possibles d'un mot de  $k$  lettres en  $n$  mots. Plus précisément, posons pour  $p \in \mathbb{N}$

$$\begin{aligned} \phi_p : \{m_1, \dots, m_n\}^p &\longrightarrow \Sigma^* \\ m_{i_1}, \dots, m_{i_p} &\longmapsto x_{i_1} \dots x_{i_p} \end{aligned}$$

Le coefficient devant  $x^k$  est  $|\phi_n^{-1}(\Sigma^k)| \leq 2^k$  car  $\phi$ , donc  $\phi_n$ , est injective. De plus, la famille  $(c_i)$  étant finie, les  $u_i$  sont nuls à partir d'un certain rang  $P_n$ .

$$\left(f\left(\frac{1}{2}\right)\right)^n \leq \sum_{k=0}^{P_n} 2^k \left(\frac{1}{2}\right)^k = P_n + 1$$

Cette égalité étant vraie pour tout  $n \geq 0$ , et ayant  $P_n = O(n)$ , par croissances comparées,  $f(\frac{1}{2}) \leq 1$ . Donc pour tout  $N \geq 0$ ,

$$\sum_{i=0}^N u_i 2^{-i} \leq 1$$

D'où la convergence de la série et la majoration cherchée.  $\square$

## 2.2 Complexité de Kolmogorov *préfixe*

On reprend la définition de la première partie en imposant que le domaine des fonctions calculables considérées soit *préfixe* (*prefix-free*). On obtient alors des propriétés analogues à celle du paragraphe 1.3, mais pas identiques.

On note encore  $\mathcal{U}$  une machine universelle et la complexité de Kolmogorov *préfixe* est notée  $K$ . Cette nouvelle complexité vérifie une sorte de sous-additivité, puisque la concaténation de deux programmes n'est plus ambiguë.

**Propriété 7.** *Pour tous mots  $\omega$  et  $\tau$ ,  $K(\omega\tau) \leq K(\omega) + K(\tau) + O(1)$ .*

## 2.3 Théorème KC

Le théorème KC (inégalité de Kraft, version Calculable) donne un moyen de calculer une suite  $(x_i)$  dont l'inégalité de Kraft affirme l'existence.

**Théorème 2.** *Soit  $(m_i)$  suite de mots et  $(c_i)$  suite d'entiers vérifiant  $\sum_i 2^{-c_i} \leq 1$ . Il existe une suite calculable  $(x_i)$  de mots de tailles  $(c_i)$  et une fonction partielle calculable telle que  $f(x_i) = m_i$  pour tout  $i$ .*

La preuve ci-dessous, proposée en [4], est assez lourde : on l'illustre par la figure 3.

*Démonstration.* Il suffit de donner un moyen effectif de construire des  $x_i$  vérifiant  $|x_i| = c_i$ . Pour tout  $n$ , on note  $\sigma^n = \sigma_1^n \dots \sigma_m^n$  le mot binaire tel que  $0, \sigma_1^n \dots \sigma_m^n = 1 - \sum_{j \leq n} 2^{-c_j}$ . On construit  $x_n$  par récurrence, en garantissant à chaque étape  $n$  la propriété suivante : pour tout  $m$  tel que  $x_m^n = 1$  il existe un mot  $\mu_m^n$  de longueur  $m$ , de sorte que  $S_n = \{x_i : i \leq n\} \cup \{\mu_m^n : \sigma_m^n = 1\}$  est *préfixe*.

On pose  $x_0 = 0^{d_0}$ . On a  $\sigma^0 = 11 \dots 1 = 1^{d_0}$ . Poser  $\mu_m^0 = 0^{m-1}1$  définit bien un ensemble  $S_0$  *préfixe*.

Supposons  $S_n$  défini.

Si  $\sigma_{c_{n+1}}^n = 1$ , alors  $\sigma^{n+1}$  est identique à  $\sigma^n$  hormis  $\sigma_{d_{n+1}}^{n+1}$  qui vaut 0. On peut donc poser  $x_{n+1} = \mu_{c_{n+1}}^n$  et  $\mu_m^{n+1} = \mu_m^n$  pour tous les  $m \neq c_{n+1}$ . L'hypothèse de récurrence est vérifiée au rang  $n+1$ .

Sinon, soit  $j < d_{n+1}$  maximal tel que  $\sigma_j^n = 1$ . Un tel  $j$  existe car sinon  $1 - \sum_{j \leq n} 2^{-c_j} < 2^{-c_{n+1}}$  et donc  $\sum_{j \leq n+1} 2^{-c_j} > 1$ .  $\sigma^{n+1}$  ne diffère de  $\sigma^n$  qu'aux positions  $j, j+1, \dots, c_{n+1}$  où  $\sigma^{n+1}$  vaut respectivement  $0, 1, \dots, 1$ . On pose alors  $x_{n+1} = \mu_j^n 0^{c_{n+1}-j}$ . Pour  $m < j$  ou  $m > c_{n+1}$ , on pose  $\mu_m^{n+1} = \mu_m^n$ . Pour  $j < m \leq c_{n+1}$ , on pose  $\mu_m^{n+1} = \mu_j^n 0^{m-j-1}$ . Vérifions que ces choix engendrent un  $S_{n+1}$  valide : on n'a fait que substituer à  $\mu_j^n$  un ensemble de mots commençant par  $\mu_j^n$  et n'étant pas *préfixes* les uns des autres, donc  $S_{n+1}$  est encore *préfixe*.  $\square$

Une première application de ce théorème est de montrer l'existence de codages des entiers asymptotiquement économes.

**Propriété 8.** *Pour tout  $p \in \mathbb{N}$  et tout  $\epsilon > 0$ ,  $K(n) \leq \log n + \log \log n + \dots + (1 + \epsilon) \log^p n + O(1)$ .*

*Démonstration.* Soit  $L > 0$ .

$$\sum_{n \geq q} 2^{-\log n - \dots - (1+\epsilon) \log^p n - L} = 2^{-L} \sum_{n \geq q} \frac{1}{n \log n \dots (\log^{p-1} n)^{1+\epsilon}}$$

Cette dernière série, dite de Bertrand, converge pour tout  $\epsilon > 0$ . Il ne reste donc qu'à choisir  $L$  assez grand pour que sa somme (complétée par des premiers termes pour les  $n < q$ ) soit inférieure à 1. Le théorème KC achève la preuve.  $\square$

Quelles sont les fonctions calculables majorant la complexité de Kolmogorov ? La caractérisation qui suit est une deuxième application du théorème KC, qui demande un peu plus de travail.

**Propriété 9.** Soit  $f$  une fonction calculable de  $\mathbb{N}$  dans  $\mathbb{N}$ . S'équivalent :

- i.  $K(n) \leq f(n) + O(1)$
- ii.  $\sum_n 2^{-f(n)} < \infty$

### 3 Compression de Huffman

#### 3.1 Principe de la compression

Soit  $t$  un mot sur l'alphabet  $\{0,1\}^n$ . On cherche un morphisme  $\phi$  par lequel le mot  $t$  sera encodé en un minimum de bits. On cherche donc à minimiser

$$\sum_{c \in \{0,1\}^n} |t|_c \cdot |\phi(c)|$$

où  $|t|_c$  désigne le nombre d'occurrences de  $c$  dans  $t$ .

On veut d'autre part que  $\phi$  soit injectif : les  $(|\phi(c)|)_c$  vérifient alors l'inégalité de Kraft, ce qui a pour conséquence qu'on peut imposer sans augmenter le coût ci-dessus que le codage soit préfixe, d'après la construction utilisée dans la preuve du théorème 2.

#### 3.2 Algorithme de Huffman

L'algorithme de Huffman est un algorithme glouton qui crée un arbre définissant un codage optimal pour le mot  $t$ . À l'initialisation, on crée une feuille par caractère présent dans le texte et on associe à chaque feuille le nombre d'occurrences du caractère dans  $t$ . On obtient ainsi une forêt de feuilles.

Tant que la forêt n'est pas réduite à un seul arbre, on sélectionne deux arbres de fréquences minimales, et on les remplace par un seul arbre dont la racine a deux fils : les racines de ces deux arbres. On attribue à ce nouvel arbre un nombre d'occurrences égal à la somme des nombre d'occurrences de ses deux fils. Ce procédé est représenté sur la figure 1.

À la fin de l'algorithme, on obtient un arbre définissant un codage qu'on peut prouver optimal. [2]

#### 3.3 Borne supérieure pour la complexité de Kolmogorov

On dispose alors d'un algorithme qui, pour tout mot dont on veut majorer la complexité de Kolmogorov, génère un programme affichant ce mot, d'une taille potentiellement plus courte que la solution qui consiste à écrire intégralement le mot en question dans le code. Plus précisément, on compresse le mot avec l'algorithme de Huffman, et on assemble dans un même programme l'arbre qui représente le codage et le texte compressé. Le programme restitue alors le mot original.

Ce procédé se révèle en général inutile, mais c'est sans surprise puisque la plupart des mots ont une complexité de Kolmogorov très proche de leur longueur. Et le programme généré est même plus long que la solution triviale, bien souvent : même si le codage est optimal, il faut ajouter l'arbre de Huffman et une certaine quantité de code pour opérer la décompression, ce qui est coûteux.



## Conclusion

L'intérêt des méthodes de compression statistique dans le cadre de l'étude de la complexité de Kolmogorov semble donc limité. Il faudrait disposer de compresseurs exploitant des redondances beaucoup plus complexes que de simples observations statistiques. Pourtant, des recherches récentes ont montré que l'utilisation de compresseurs classiques pour approximer la complexité de Kolmogorov permet, par exemple, de mesurer la *similarité* entre deux mots. Ces techniques trouvent des applications dans des algorithmes de classification non supervisée.

## Références

- [1] Olivier Carton. *Langages formels, calculabilité et complexité*. Vuibert, 2008.
- [2] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction à l'algorithmique*. Dunod, 2002.
- [3] Jean-Paul Delahaye. *Information, complexité et hasard*. Hermès, 1999.
- [4] Rod Downey and Denis Hirschfeldt. *Algorithmic Randomness and Complexity*. Springer Verlag, 2010.

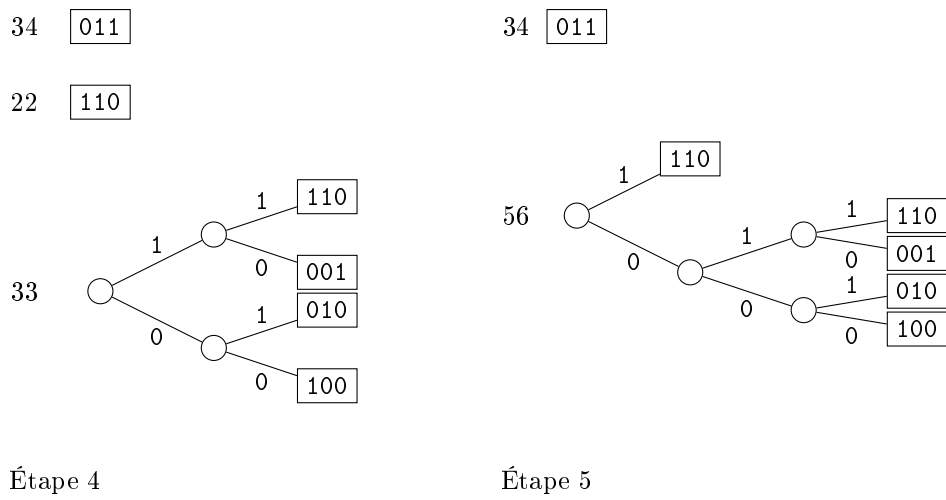
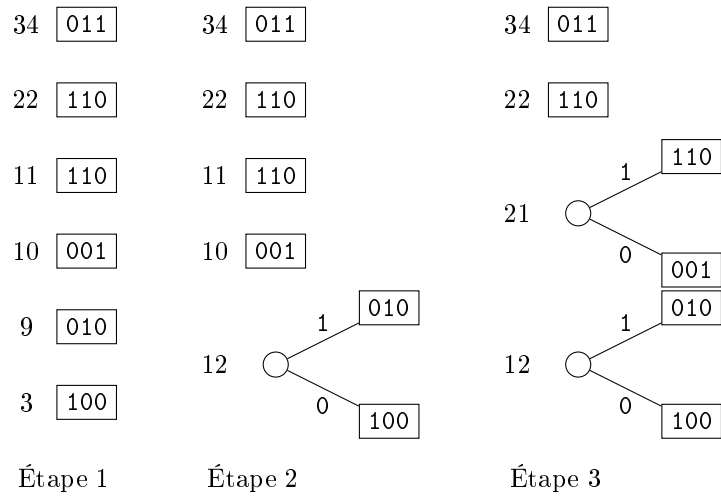


FIGURE 1 – Le déroulement de l’algorithme de Huffman

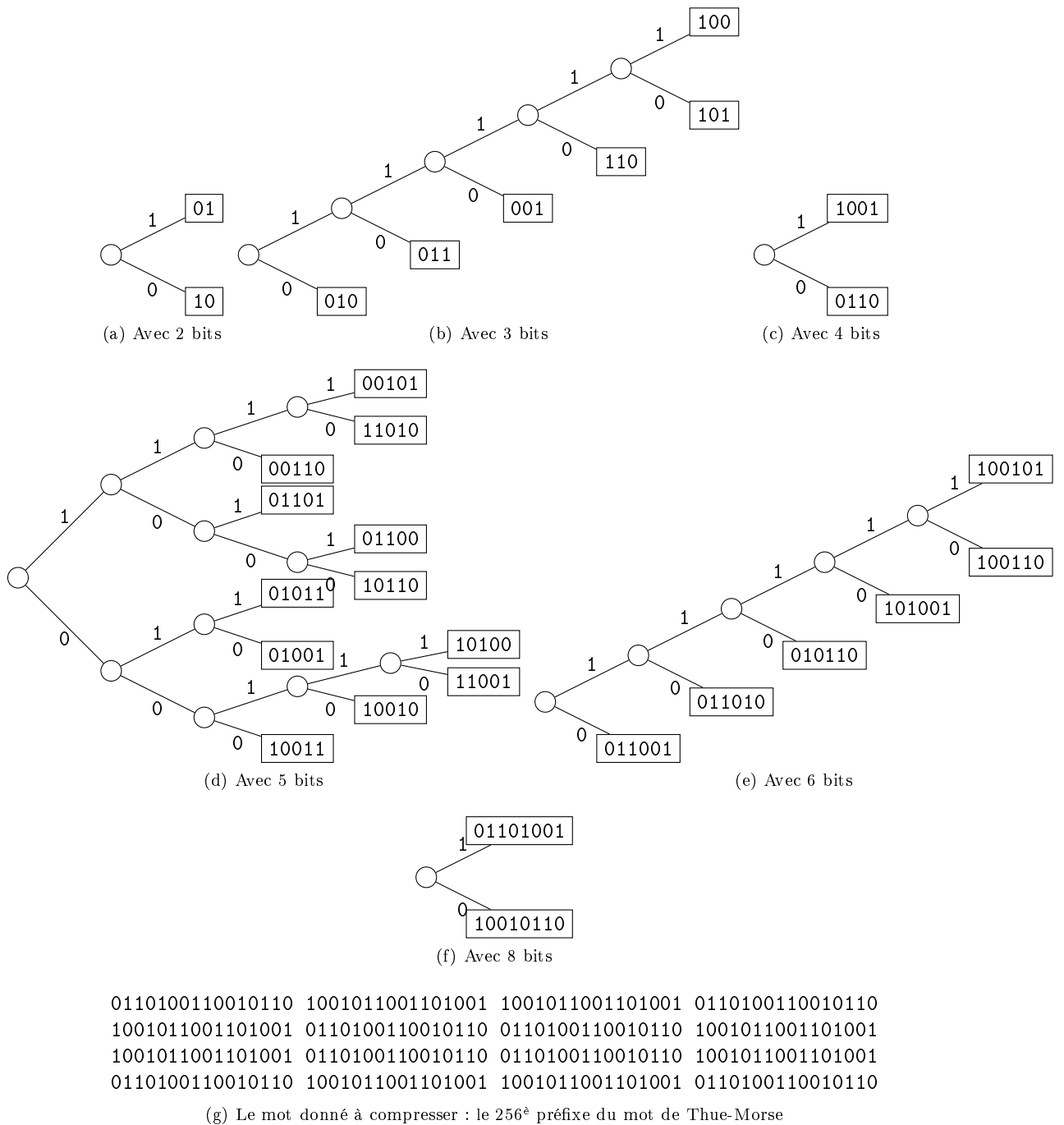
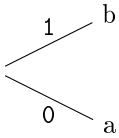


FIGURE 2 – Les arbres de Huffman obtenus en choisissant différentes tailles de caractère

Dans cette figure, on omet les étiquettes : 

On détaille les trois premières étapes de l'algorithme pour  $c_0 = 3, c_1 = 1, c_2 = 5$  qui correspondent aux trois cas décrits dans la preuve.

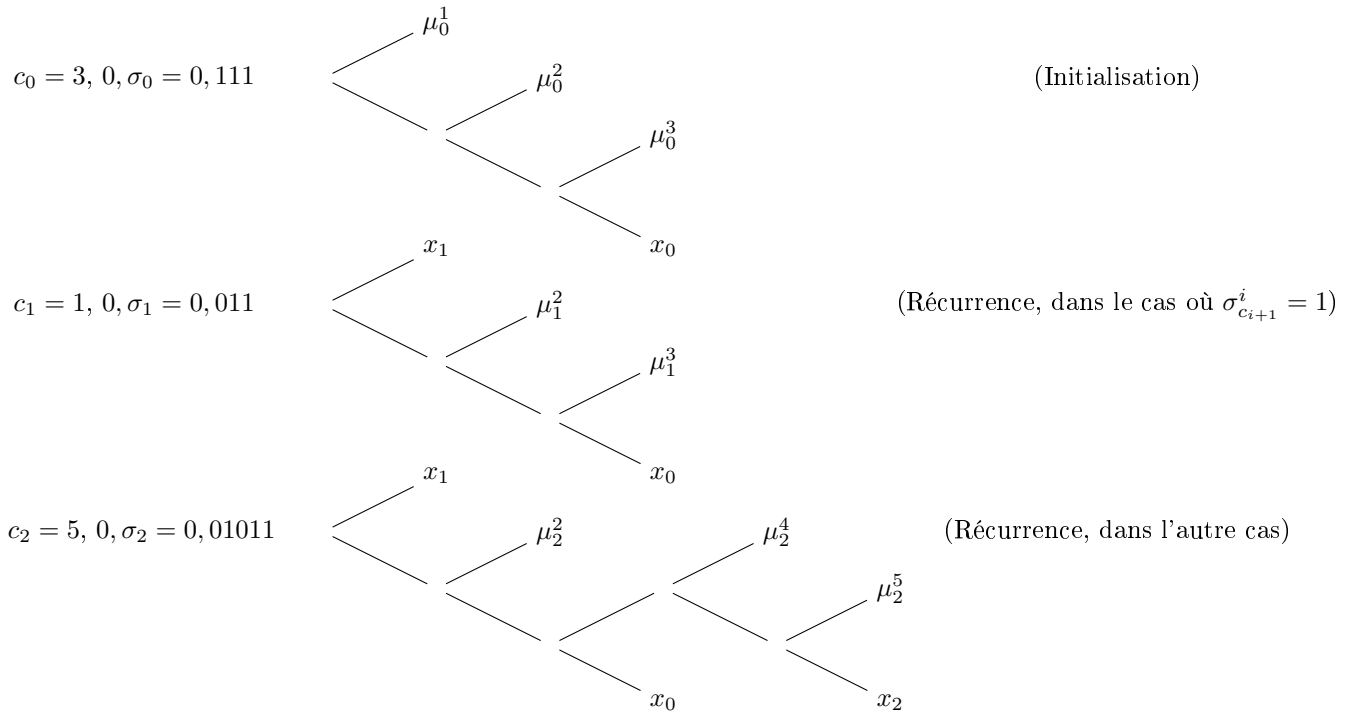


FIGURE 3 – Construction par étapes de l'arbre des  $x_i$

## A Code source

### A.1 Module Compression

```
1  (*****  
2  * OUTILS GÉNÉRIQUES POUR LA COMPRESSION AVEC UN CODE PRÉFIXE *  
3  *****)  
4  
5  open Utilitaires  
6  
7  (* Un mot est une liste de nombres, a priori. En pratique, on  
8  * travaille volontiers avec un alphabet binaire. *)  
9  type mot = int list;;  
10  
11  (* Type de l'arbre de Huffman *)  
12  type codage =  
13  | Code of mot  
14  | Noeud of codage * codage;;  
15  
16  (* Ordre sur les arbres pondérés : nécessaire pour la gestion du tas. *)  
17  let ordre_tas a b = (fst a) > (fst b);;  
18  
19  (* Taille initiale de la table de hachage en fonction de la taille des caractères :  
20  * dans le pire des cas, 2p caractères.  
21  * En pratique, moins sur des textes particuliers. *)  
22  let taille_table_de_hachage p = p ;;  
23  
24  (* Décoder un texte compressé en utilisant l'arbre de Huffman passé en paramètre  
25  * *)  
26  let rec decoder arbre pos lst = match lst,pos with  
27  | ([],Noeud(_)) -> []  
28  | (_,Code(m)) -> concat_listes (decoder arbre arbre lst) m  
29  | (0::t,Noeud(p,q)) -> decoder arbre p t  
30  | (1::t,Noeud(p,q)) -> decoder arbre q t  
31  | (_::t,n) -> decoder arbre n t (* Pour que le filtrage soit complet *);;  
32  
33  (* Coupe le texte en mots de la taille donnée *)  
34  let rec decouper_mots taille_mot restant accu texte = match texte,restant with  
35  | (_,0) -> accu::(decouper_mots taille_mot taille_mot [] texte)  
36  | ([],n) -> []  
37  | (h::t,n) -> decouper_mots taille_mot (n-1) (h::accu) t ;;  
38  
39  (* Construit un tableau indexé par les mots donnant le codage compressé de ces  
40  * mots. *)  
41  let construire_traduction arbre taille_mot =  
42  (* La table de compression, notée trad (qui associe à chaque caractère son codage)  
43  * est une table de hachage *)  
44  let trad = Hashtbl.create (taille_table_de_hachage taille_mot) in  
45  
46  (* On réalise un parcours en profondeur *)  
47  let rec parcourir_arbre trad accu = fonction  
48  | Code(m) -> Hashtbl.add trad m (rev accu)  
49  | Noeud(p,q) -> (parcourir_arbre trad (0::accu) p);  
50  (parcourir_arbre trad (1::accu) q) in  
51  (parcourir_arbre trad [] arbre); trad ;;  
52
```

```

53  (* Compression d'un texte à partir d'un arbre *)
54  let compresser arbre taille_mot texte =
55    let trad = (construire_traduction arbre taille_mot) in
56
57    (* Prend une liste de caractères et renvoie la liste de leurs
58     * équivalents compressés *)
59    let rec compresser_mots trad = map (fun h -> Hashtbl.find trad h) in
60
61    (* Concatène les mots pour obtenir le texte compressé *)
62    let concatene_mots = it_list (fun a b -> a @ b) [] in
63
64    (* On assemble le tout... *)
65    concatene_mots (compresser_mots trad (decouper_mots taille_mot taille_mot [] texte));;

```

## A.2 Module Tas

```

1  open Compression
2
3  (* Structure retenue pour le tas :
4   * un tableau dont la première case est le nombre d'éléments stockés dans le tas
5   * et les éléments qui suivent après, selon la hiérarchie suivante : *)
6
7  let pere i = i/2;;
8  let gauche i = 2*i;;
9  let droite i = 2*i + 1;;
10
11  let taille_tas t = (fst t.(0));;
12  let incr_taille_tas t =
13    t.(0) <- ((taille_tas t)+1,(snd t.(0)));;
14  let decr_taille_tas t =
15    t.(0) <- ((taille_tas t)-1,(snd t.(0)));;
16
17  (* L'ordre avec lequel le tas est trié est donné par la fonction
18   * ordre_tas qui doit être définie avant l'inclusion de ce fichier
19   * Signification : ordre_tas a b renvoie true si a < b *)
20
21  (* Échange deux éléments dans un tableau *)
22  let echange tab i j =
23    let temp = tab.(i) in
24    tab.(i) <- tab.(j);
25    tab.(j) <- temp;;
26
27  (* Entasse un élément dont les sous-arbres sont valides *)
28  let rec entasser_haut tab = fonction
29    | 1 -> ()
30    | n -> if (ordre_tas tab.(pere n) tab.(n)) then
31      begin
32        echange tab n (pere n);
33        entasser_haut tab (pere n)
34      end;;
35
36  (* Entasse un élément dont le père est valide *)
37  let rec entasser_bas tab elem =
38    let max = ref elem in
39    if ((gauche elem) <= (taille_tas tab) && (ordre_tas tab.(!max) tab.(gauche elem))) then
40      max := (gauche elem);

```

```

41   if ((droite elem) <= (taille_tas tab) && (ordre_tas tab.(!max) tab.(droite elem))) then
42     max := (droite elem);
43   if !max <> elem then
44     begin
45       echange tab !max elem;
46       entasser_bas tab !max
47     end;;
48
49
50   ( Insère un élément dans un tas *)
51   let insere_dans_tas tab elem =
52     if (Array.length tab) <= (taille_tas tab)+1 then
53       failwith "Tas plein";
54     incr_taille_tas tab;
55     tab.(taille_tas tab) <- elem;;
56
57   ( Supprime le minimum du tas *)
58   let pop_minimum tab =
59     if (taille_tas tab) = 0 then
60       failwith "Tas vide";
61     let mini = tab.(1) in
62     tab.(1) <- tab.(taille_tas tab);
63     decr_taille_tas tab;
64     entasser_bas tab 1;
65     mini;;

```

### A.3 Module Huffman

```

1   ( Définit des fonctions classiques d'arithmétique ou de manipulation de listes *)
2   open Utilitaires
3   ( Définit le codage et le décodage d'un texte avec un arbre de codage préfixe donné *)
4   open Compression
5
6   (*****
7    * GÉNÉRATION DE L'ARBRE DE HUFFMAN *
8    *****)
9
10  ( Affichage du tableau du nombre d'occurrences de chaque caractère *)
11  let afficher_occurences compte =
12    let afficher_paire w occur =
13      print_list "" " " : " w;
14      print_int occur;
15      print_newline () in
16
17    Hashtbl.iter afficher_paire compte ;;
18
19  ( Incrémente la case d'une table de hachage correspondant à la clé w *)
20  let incr_occur tbl w =
21    let oldval =
22      try
23        Hashtbl.find tbl w
24      with Not_found -> 0 in
25
26    if oldval <> 0 then
27      Hashtbl.remove tbl w;
28      Hashtbl.add tbl w (oldval + 1);;

```

```

29
30 (* Compte les occurrences de chaque mot de taille donnée dans le texte
31 * (toujours avec notre définition particulière d'une occurrence d'un
32 * mot : il faut que le mot commence à une position multiple de
33 * taille_mot) *)
34 let compter_occurrences taille_mot texte =
35   let compte = Hashtbl.create (taille_table_de_hachage taille_mot) in
36   let rec compte_entiers = function
37     | [] -> ()
38     | h::t -> incr_occur compte (h);
39               compte_entiers t in
40     compte_entiers (decouper_mots taille_mot taille_mot [] texte);
41
42     afficher_occurrences compte;
43     compte;;
44
45
46 (* Définit les fonctions de manipulation d'un tas *)
47 open Tas
48
49 (* Construit le tas initial des feuilles *)
50 let construire_tas taille_mot compte =
51   let tab = Array.make ((Hashtbl.length compte)+1) (0,Code([])) in
52   tab.(0) <- (0,Code([]));
53
54   (* Copie des occurrences non nulles dans le tas *)
55   let ajouter_mot_dans_tas mot occur =
56     incr_taille_tas tab;
57     tab.(taille_tas tab) <- (occur,Code(mot)) in
58
59   Hashtbl.iter ajouter_mot_dans_tas compte ;
60
61   (* Création d'un tas valide *)
62   for i = 1 to (taille_tas tab) do
63     entasser_haut tab i;
64   done;
65
66   tab;;
67
68 (* Fusionne les noeuds tant qu'il y en a plusieurs *)
69 let rec fusionner_noeuds tab =
70   if (taille_tas tab) > 1 then
71     begin
72       let (p1,n1) = (pop_minimum tab) in
73       let (p2,n2) = (pop_minimum tab) in
74       insere_dans_tas tab (p1+p2, Noeud(n1,n2));
75       fusionner_noeuds tab
76     end
77   else snd tab.(1);;
78
79
80 (* Mise en commun : génération de l'arbre de Huffman d'un texte, pour une
81 * taille de mot particulière *)
82 let arbre_huffman taille_mot texte =
83   let etape1 = compter_occurrences taille_mot texte in
84   let etape2 = construire_tas taille_mot etape1 in

```

```

85   let etape3 = fusionner_noeuds etape2 in
86   etape3;;
87
88   (*****
89   * GÉNÉRATION DU DÉCODEUR EN bf *
90   * *****)
91
92   let bf_affichage_mot m =
93     let rec afficher_mot etat_init restant = match etat_init,restant with
94       | (0,[]) -> ""
95       | (1,[]) -> "-"
96       | (_,[]) -> "[-]" (* Pour que le filtrage soit complet *)
97       | (0,0:h) -> strconcat "." (afficher_mot 0 h)
98       | (1,1:h) -> strconcat "." (afficher_mot 1 h)
99       | (0,1:h) -> strconcat "+." (afficher_mot 1 h)
100      | (1,0:h) -> strconcat "-." (afficher_mot 0 h)
101      | (_,c:h) -> afficher_mot c h (* Pour que le filtrage soit complet *)
102    in
103    afficher_mot 0 m;;
104
105   let bf_condition est_racine a b =
106     let debut = (if est_racine then "" else ",") in
107     String.concat "" [debut; ">+<[>>"; b; "<-<-]>[->"; a; "<<"];;
108
109   let decompresseur_bf arbre =
110     let rec traduire_arbre est_racine = function
111       | Code(m) -> bf_affichage_mot (rev m)
112       | Noeud(a,b) -> bf_condition est_racine (traduire_arbre false a) (traduire_arbre false b) in
113     String.concat "" ["+[-"; (traduire_arbre true arbre); "+"]];;

```

#### A.4 Calcul du mot de Thue-Morse

```

1   (*****
2   * Première méthode de génération du mot de Thue-Morse :
3   *
4   * Définir le morphisme  $\sigma : 0 \mapsto 01 ; 1 \mapsto 10$  calculant
5   *  $\sigma(w)$  à partir de  $w$  représenté sous la forme d'une liste.
6   *
7   * Il suffit alors d'itérer le morphisme jusqu'à avoir le
8   * nombre désiré de termes
9   * *****)
10
11  let rec print_list = function
12    | [] -> print_newline ()
13    | t::q -> print_int t; print_list q ;;
14
15  let rec morphisme = function
16    | [] -> []
17    | 0::q -> 0::(1::(morphisme q))
18    | _::q -> 1::(0::(morphisme q)) ;;
19
20  let rec iterer x = function
21    | 0 -> x
22    | n -> morphisme (iterer x (n-1)) ;;
23
24  (*****

```

```

25  * Autre méthode, avec l'évaluation paresseuse
26  *
27  * Implémentation de l'évaluation paresseuse
28  * largement inspirée de l'article de Pierre Weiss
29  * dans la Lettre de Caml no. 2
30  * *****)
31
32  (* Le type glaçon est l'outil de base de l'évaluation paresseuse :
33  * Un objet peut être soit connu, soit calculable plus tard en
34  * exécutant la fonction ne prenant aucun argument. *)
35
36  type 'a glaçon =
37  | Connu of 'a
38  | Gele of (unit -> 'a) ;;
39
40  (* Une liste paresseuse est constituée, comme les listes classiques,
41  * de [] et de hd::tl. La différence ici est que tl n'est pas forcément
42  * connu, grace au type glaçon. L'avantage est de pouvoir gérer des
43  * listes « infinies » : on ne les calcule que jusqu'à un certain rang. *)
44
45  type 'a liste_par =
46  | Nil
47  | Cons of 'a cellule
48  and 'a cellule = {hd : 'a; mutable tl : 'a liste_par glaçon} ;;
49
50  (* La fonction force permet de connaître un élément de plus de la liste
51  * donnée. Elle renvoie la queue de cette liste. *)
52
53  let force cel =
54  match cel.tl with
55  | Connu v -> v
56  | Gele f -> let resultat_f = f () in
57  cel.tl <- Connu (resultat_f);
58  resultat_f ;;
59
60  let print_bool = function
61  | false -> print_string "0"
62  | true -> print_string "1" ;;
63
64  (* Un do_list pour liste paresseuse : il faut s'arrêter au bout de
65  * n étapes. *)
66
67  let rec do_list_par f n = function
68  | Nil -> print_newline ()
69  | Cons cel ->
70  if n = 0 then print_newline ()
71  else
72  begin
73  f cel.hd;
74  do_list_par f (n-1) (force cel);
75  end ;;
76
77  (* Fonction map pour les listes paresseuses *)
78
79  let rec map_par f = function
80  | Nil -> Nil

```

```

81   | Cons cel -> Cons {hd = f cel.hd; tl = Gele (fun () -> map_par f (force cel))} ;;
82
83   (* Application : calcul du mot infini de Thue-Morse
84   * On utilise la relation de récurrence suivante :
85   * - Le premier terme est 0.
86   * - Connaissant les 2n premiers termes, les 2n suivants
87   *   sont les négations des 2n premiers, dans le même ordre. *)
88
89   let rec recopier_en_niant = fonction
90   | Nil -> Nil
91   | Cons cel -> Cons {hd = not cel.hd; tl = (match cel.tl with
92     | Connu v -> Connu (recopier_en_niant v)
93     | Gele _ -> Gele (fun () -> recopier_en_niant tm))}
94   and
95   tm = Cons {hd = false; tl = Gele (fun () -> recopier_en_niant tm) } ;;
96
97   do_list_par print_bool 36 tm ;;
98
99
100  (* Troisième méthode
101  *
102  * Généralisable pour obtenir un point fixe de
103  * n'importe quel morphisme de la forme
104  *  $\sigma : a \mapsto a\omega$ 
105  *  $b \mapsto \sigma(b)$ 
106  *
107  *  $\vdots$ 
108  *  $h \mapsto \sigma(h)$ 
109  * dont le point fixe commençant par a est
110  *  $a \omega \sigma(\omega) \sigma^2(\omega) \dots$ 
111  *)
112  (* Équivalent de tl pour les listes paresseuses *)
113
114  let tail_par = fonction
115  | Nil -> Nil
116  | Cons cel -> (force cel);;
117
118  (* Morphisme du mot de Thue-Morse *)
119  let rec phi = fonction
120  | Nil -> Nil
121  | Cons cel -> let reste = Gele (fun () -> phi (force cel)) in
122    match cel.hd with
123    | 0 -> Cons {hd = 0; tl = Connu
124      (Cons {hd = 1; tl = reste})}
125    | _ -> Cons {hd = 1; tl = Connu
126      (Cons {hd = 0; tl = reste})} ;;
127
128  let rec point_fixe = Cons { hd = 0; tl = Connu
129    (Cons { hd = 1; tl = Gele
130      (fun () -> (phi (tail_par point_fixe))) })
131    };;
132
133  do_list_par print_int 36 point_fixe ;;

```