



# Chapitre 3

## Codage de l'information

### 3.1. Vocabulaire

Quelle que soit la nature de l'information traitée par un ordinateur (image, son, texte, vidéo), elle l'est toujours sous la forme d'un ensemble de nombres écrits en base 2, par exemple 01001011.

Le terme **bit** (b minuscule dans les notations) signifie « binary digit », c'est-à-dire 0 ou 1 en numérotation binaire. Il s'agit de la plus petite unité d'information manipulable par une machine numérique. Il est possible de représenter physiquement cette information binaire par un signal électrique ou magnétique, qui, au-delà d'un certain seuil, correspond à la valeur 1.

L'**octet** (en anglais *byte* ou B majuscule dans les notations) est une unité d'information composée de 8 bits. Il permet par exemple de stocker un caractère comme une lettre ou un chiffre.

Une unité d'information composée de 16 bits est généralement appelée **mot** (en anglais *word*).

Une unité d'information de 32 bits de longueur est appelée **mot double** (en anglais *double word*, d'où l'appellation *dword*).

Beaucoup d'informaticiens ont appris que 1 kilooctet valait 1024 octets. Or, depuis décembre 1998, l'organisme international *IEC* a statué sur la question<sup>1</sup>.

Voici les unités standardisées :

- Un kilooctet (ko) =  $10^3$  octets
- Un mégaoctet (Mo) =  $10^6$  octets
- Un gigaoctet (Go) =  $10^9$  octets
- Un téraoctet (To) =  $10^{12}$  octets
- Un pétaoctet (Po) =  $10^{15}$  octets
- Un exaoctet (Eo) =  $10^{18}$  octets
- Un zettaoctet (Zo) =  $10^{21}$  octets
- Un yottaoctet (Yo) =  $10^{24}$  octets
- Un brontoctet =  $10^{27}$  octets (non officiel)

Avez-vous déjà acheté un disque dur et constaté, en l'utilisant pour la première fois, que sa taille réelle était sensiblement plus petite que celle annoncée par le fabricant ? Sachez que le fabricant ne vous a techniquement pas menti. Il a profité d'une confusion courante entre deux systèmes de préfixes d'unités.

En effet, lors du développement des premiers ordinateurs, les informaticiens ont décidé d'utiliser le préfixe « kilo » pour désigner 1024, puisqu'elle est raisonnablement proche de 1000. Cette tendance s'est poursuivie ensuite : un groupe de 1024 kilooctets a été appelé un mégaoctet, un groupe de 1024 mégaoctets a été appelé gigaoctets, et ainsi de suite. Alors que le passage successif entre les préfixes kilo, méga, téra, ..., correspond en principe à un facteur 1000, il correspond donc à un facteur 1024 en informatique. Un mégaoctet devrait en principe valoir  $1000 \times 1000$  octets, c'est-à-dire 1'000'000 d'octets, mais il vaut  $1024 \times 1024$  octets en informatique, c'est-à-dire 1'048'576 octets... ce qui correspond à une différence de 4.63 % !

<sup>1</sup> <http://physics.nist.gov/cuu/Units/binary.html>

## 3.2. Les bases décimale, binaire et hexadécimale

Nous utilisons le système décimal (base 10) dans nos activités quotidiennes. Ce système est basé sur dix symboles, de 0 à 9, avec une unité supérieure (dizaine, centaine, etc.) à chaque fois que dix unités sont comptabilisées. C'est un système **positionnel**, c'est-à-dire que l'endroit où se trouve le symbole définit sa valeur. Ainsi, le 2 de 523 n'a pas la même valeur que le 2 de 132. En fait, 523 est l'abréviation de  $5 \cdot 10^2 + 2 \cdot 10^1 + 3 \cdot 10^0$ . On peut selon ce principe imaginer une infinité de systèmes numériques fondés sur des bases différentes.

En informatique, outre la base 10, on utilise très fréquemment **le système binaire** (base 2) puisque l'algèbre booléenne est à la base de l'électronique numérique. Deux symboles suffisent : 0 et 1.

On utilise aussi très souvent **le système hexadécimal** (base 16) du fait de sa simplicité d'utilisation et de représentation pour les mots machines (il est bien plus simple d'utilisation que le binaire). Il faut alors six symboles supplémentaires : A (qui représente le 10), B (11), C (12), D (13), E (14) et F (15).

Le tableau ci-dessous montre la représentation des nombres de 0 à 15 dans les bases 10, 2 et 16.

Décimal	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Binaire	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111
Hexadécimal	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F

### 3.2.1. Conversion décimal - binaire

Convertissons 01001101 en décimal à l'aide du schéma ci-dessous :

$$\begin{array}{cccccccc}
 2^7 & 2^6 & 2^5 & 2^4 & 2^3 & 2^2 & 2^1 & 2^0 \\
 0 & 1 & 0 & 0 & 1 & 1 & 0 & 1
 \end{array}$$

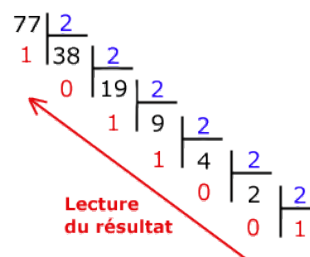
Le nombre en base 10 est  $2^6 + 2^3 + 2^2 + 2^0 = 64 + 8 + 4 + 1 = 77$ .

Allons maintenant dans l'autre sens et écrivons 77 en base 2. Il s'agit de faire une suite de divisions euclidiennes par 2. Le résultat sera la juxtaposition des restes.

Le schéma ci-contre explique la méthode mieux qu'un long discours :

77 s'écrit donc en base 2 : 1001101.

Si on l'écrit sur un octet, cela donne : 01001101.



### 3.2.2. Conversion hexadécimal - binaire

Pour convertir un nombre binaire en hexadécimal, il suffit de faire des **groupes de quatre bits** (en commençant depuis la droite). Par exemple, convertissons 1001101 :

Binaire	0100	1101
Pseudo-décimal	4	13
Hexadécimal	4	D

1001101 s'écrit donc en base 16 : 4D.

Pour convertir d'hexadécimal en binaire, il suffit de lire ce tableau de bas en haut.

#### Exercice 3.1

Donnez la méthode pour passer de la base décimale à la base hexadécimale (dans les deux sens).

**Exercice 3.2**

Complétez le tableau ci-dessous. L'indice indique la base dans laquelle le nombre est écrit.

	Bases		
	2	10	16
1001010110 <sub>2</sub>			
2002 <sub>10</sub>			
A1C <sub>16</sub>			



**Exercice 3.3**

Écrivez en Python un programme permettant de convertir un nombre d'une base de départ  $d$  vers une base d'arrivée  $a$  ( $d$  et  $a$  compris entre 2 et 16).

### 3.3. Représentation des nombres entiers

#### 3.3.1. Représentation d'un entier naturel

Un entier naturel est un nombre entier positif ou nul. Le choix à faire (c'est-à-dire le nombre de bits à utiliser) dépend de la fourchette des nombres que l'on désire utiliser. Pour coder des nombres entiers naturels compris entre 0 et 255, il nous suffira de 8 bits (un octet) car  $2^8 = 256$ . D'une manière générale un codage sur  $n$  bits pourra permettre de représenter des nombres entiers naturels compris entre 0 et  $2^n - 1$ .

**Exemples :**  $9 = 00001001_2$ ,  $128 = 10000000_2$ , etc.

#### 3.3.2. Représentation d'un entier relatif

Un entier relatif est un entier pouvant être négatif. Il faut donc coder le nombre de telle façon que l'on puisse savoir s'il s'agit d'un nombre positif ou d'un nombre négatif, et il faut de plus que les règles d'addition soient conservées. L'astuce consiste à utiliser un codage que l'on appelle **complément à deux**. Cette représentation permet d'effectuer les opérations arithmétiques usuelles naturellement.

- **Un entier relatif positif** ou nul sera représenté en binaire (base 2) comme un entier naturel, à la seule différence que le bit de poids fort (le bit situé à l'extrême gauche) représente le signe. Il faut donc s'assurer pour un entier positif ou nul qu'il est à zéro (0 correspond à un signe positif, 1 à un signe négatif). Ainsi, si on code un entier naturel sur 4 bits, le nombre le plus grand sera 0111 (c'est-à-dire 7 en base décimale).
  - Sur 8 bits (1 octet), l'intervalle de codage est  $[-128, 127]$ .
  - Sur 16 bits (2 octets), l'intervalle de codage est  $[-32768, 32767]$ .
  - Sur 32 bits (4 octets), l'intervalle de codage est  $[-2147483648, 2147483647]$ .

D'une manière générale le plus grand entier relatif positif codé sur  $n$  bits sera  $2^{n-1} - 1$ .

- **Un entier relatif négatif** sera représenté grâce au codage en complément à deux.

#### Principe du complément à deux

1. Écrire la valeur absolue du nombre en base 2. Le bit de poids fort doit être égal à 0.
2. Inverser les bits : les 0 deviennent des 1 et vice versa. On fait ce qu'on appelle le complément à un.
3. On ajoute 1 au résultat (les dépassements sont ignorés).

## Codage de l'information

Cette opération correspond au calcul de  $2^n - |x|$ , où  $n$  est la longueur de la représentation et  $|x|$  la valeur absolue du nombre à coder.

Ainsi  $-1$  s'écrit comme  $256 - 1 = 255 = 11111111_2$ , pour les nombres sur 8 bits.

### Exemple

On désire coder la valeur  $-19$  sur 8 bits. Il suffit :

1. d'écrire 19 en binaire : 00010011
2. d'écrire son complément à 1 : 11101100
3. et d'ajouter 1 : 11101101

La représentation binaire de  $-19$  sur 8 bits est donc 11101101.

On remarquera qu'en additionnant un nombre et son complément à deux on obtient 0. En effet,  $00010011 + 11101101 = 00000000$  (avec une retenue de 1 qui est éliminée).

### Astuce

Pour transformer de tête un nombre binaire en son complément à deux, on parcourt le nombre de droite à gauche en laissant inchangés les bits jusqu'au premier 1 (compris), puis on inverse tous les bits suivants. Prenons comme exemple le nombre 20 : 00010100.

1. On garde la partie à droite telle quelle : 00010**100**
2. On inverse la partie de gauche après le premier un : **1110**1100
3. Et voici  $-20$  : 11101100

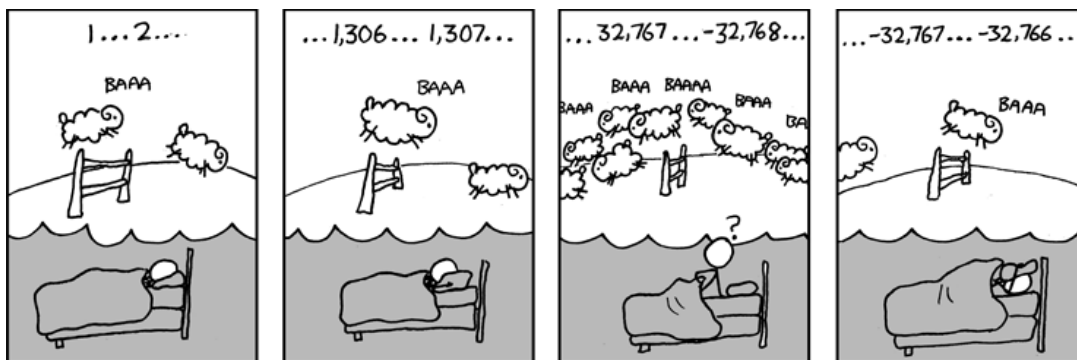
### Exercice 3.4

1. Codez les entiers relatifs suivants sur 8 bits (16 si nécessaire) : 456,  $-1$ ,  $-56$ ,  $-5642$ .
2. Que valent en base dix les trois entiers relatifs suivants :

01101100  
11101101  
10101010101010 ?

### Exercice 3.5

Expliquez ce rêve étrange (source de l'image : <http://xkcd.com/571>).



Le 4 juin 1996, une fusée Ariane 5 a explosé 40 secondes après l'allumage. La fusée et son chargement avaient coûté 500 millions de dollars. La commission d'enquête a rendu son rapport au bout de deux semaines. Il s'agissait d'une erreur de programmation dans le système inertiel de référence. À un moment donné, un nombre codé en virgule flottante sur 64 bits (qui représentait la vitesse horizontale de la fusée par rapport à la plate-forme de tir) était converti en un entier sur 16 bits. Malheureusement, le nombre en question était plus grand que 32768 (le plus grand entier que l'on peut coder sur 16 bits) et la conversion a été incorrecte.

**Exercice 3.6**

Certains logiciels utilisent la représentation POSIX du temps, dans laquelle le temps est représenté comme un nombre de secondes écoulées depuis le 1<sup>er</sup> janvier 1970 à minuit (0 heure). Sur les ordinateurs 32 bits, la plupart des systèmes d'exploitation représentent ce nombre comme un **nombre entier signé de 32 bits**.

- Quel est le nombre de secondes maximum que l'on peut représenter ?
- À quelle date cela correspond-il (jour, mois, année, heures, minutes, secondes).

*Indications :*

- afin de tenir compte des années bissextiles, comptez par cycles de 4 ans composés de  $4 \cdot 365 + 1 = 1461$  jours ;
  - l'an 2000 est une année bissextile.
- Que se passera-t-il une seconde plus tard ? Quel sera le nombre de secondes affiché (en base 10) ? À quelle date cela correspond-il ?

**3.4. Représentation des nombres réels**

En base 10, l'expression 652,375 est une manière abrégée d'écrire :

$$6 \cdot 10^2 + 5 \cdot 10^1 + 2 \cdot 10^0 + 3 \cdot 10^{-1} + 7 \cdot 10^{-2} + 5 \cdot 10^{-3}.$$

Il en va de même pour la base 2. Ainsi, l'expression 110,101 signifie :

$$1 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 + 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3}.$$

**3.4.1. Conversion de binaire en décimal**

On peut ainsi facilement convertir un nombre réel de la base 2 vers la base 10. Par exemple :

$$110,101_2 = 1 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 + 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} = 4 + 2 + 0,5 + 0,125 = 6,625_{10}.$$

**Exercice 3.7**

Transformez  $0,0101010101_2$  en base 10.

Transformez  $11100,10001_2$  en base 10.

**3.4.2. Conversion de décimal en binaire**

Le passage de base 10 en base 2 est plus subtil. Par exemple : convertissons 1234,347 en base 2.

- La partie entière se transforme comme au § 3.2.1 :  $1234_{10} = 10011010010_2$
- On transforme la partie décimale selon le schéma suivant :

$0,347 \cdot 2 = 0,694$	$0,347 = 0,0\dots$
$0,694 \cdot 2 = 1,388$	$0,347 = 0,01\dots$
$0,388 \cdot 2 = 0,766$	$0,347 = 0,010\dots$
$0,766 \cdot 2 = 1,532$	$0,347 = 0,0101\dots$
$0,532 \cdot 2 = 1,064$	$0,347 = 0,01011\dots$
$0,104 \cdot 2 = 0,208$	$0,347 = 0,010110\dots$
$0,208 \cdot 2 = 0,416$	$0,347 = 0,0101100\dots$
$0,416 \cdot 2 = 0,832$	$0,347 = 0,01011000\dots$
$0,832 \cdot 2 = 1,664$	$0,347 = 0,010110001\dots$

On continue ainsi jusqu'à la précision désirée...

**Attention ! Un nombre à développement décimal fini en base 10 ne l'est pas forcément en base 2.**

Cela peut engendrer de mauvaises surprises. Prenons par exemple ce programme Python :



```
i=0.0
while i<1:
    print(i)
    i+=0.1
```

Résultat attendu	Résultat réel
0.0	0.0
0.1	0.1
0.2	0.2
0.3	0.300000000000000004
0.4	0.4
0.5	0.5
0.6	0.6
0.7	0.7
0.8	0.7999999999999999
0.9	0.8999999999999999
	0.9999999999999999

Le problème vient du fait que l'on n'arrive pas à représenter exactement 0.3 en binaire.



Le 25 février 1991, pendant la Guerre du Golfe, une batterie américaine de missiles Patriot, à Dharan (Arabie Saoudite), a échoué dans l'interception d'un missile Scud irakien. Le Scud a frappé un baraquement de l'armée américaine et a tué 28 soldats. La commission d'enquête a conclu à un calcul incorrect du temps de parcours, dû à un problème d'arrondi. Les nombres étaient représentés en virgule fixe sur 24 bits. Le temps était compté par l'horloge interne du système en dixième de seconde. Malheureusement, 1/10 n'a pas d'écriture finie dans le système binaire :  $1/10 = 0,1$  (dans le système décimal) =  $0,0001100110011001100110011\dots$  (dans le système binaire). L'ordinateur de bord arrondissait 1/10 à 24 chiffres, d'où une petite erreur dans le décompte du temps pour chaque 1/10 de seconde. Au moment de l'attaque, la batterie de missile Patriot était allumée depuis environ 100 heures, ce qui a entraîné une accumulation des erreurs d'arrondi de 0,34 s. Pendant ce temps, un missile Scud parcourt environ 500 m, ce qui explique que le Patriot soit passé à côté de sa cible.

### Exercice 3.8

Transformez en base 2 :  $0,5625_{10}$  ;  $0,15_{10}$  ;  $12,9_{10}$  .

### 3.4.2. La norme IEEE 754

IEEE, que l'on peut prononcer « i 3 e » :  
*Institute of Electrical and Electronics Engineers*

La norme IEEE 754 définit la façon de coder un nombre réel. Cette norme se propose de coder le nombre sur 32 bits et définit trois composantes :

- le signe est représenté par un seul bit, le bit de poids fort
- l'exposant est codé sur les 8 bits consécutifs au signe
- la mantisse (les bits situés après la virgule) sur les 23 bits restants



Certaines conditions sont toutefois à respecter pour les exposants :

- l'exposant 00000000 est interdit.
- l'exposant 11111111 est interdit. On s'en sert toutefois pour signaler des erreurs, on appelle alors cette configuration du nombre NaN, ce qui signifie « Not a number ».
- il faut rajouter 127 (01111111) à l'exposant pour une conversion de décimal vers un nombre réel binaire. Les exposants peuvent ainsi aller de -254 à 255.

La formule d'expression des nombres réels est ainsi la suivante :

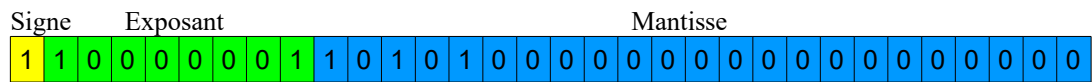
$$(-1)^S * 2^{(E - 127)} * (1 + F)$$

- $S$  est le bit de signe et l'on comprend alors pourquoi 0 est positif ( $-1^0=1$ ),
- $E$  est l'exposant auquel on doit bien ajouter 127 pour obtenir son équivalent codé,
- $F$  est la partie fractionnaire.

### Exemple

Traduisons en binaire, en utilisant la norme IEEE 754, le nombre  $-6,625$ .

- Codons d'abord la valeur absolue en binaire :  $6,625_{10} = 110,1010_2$
- Nous mettons ce nombre sous la forme : **1, partie fractionnaire**  
 $110,1010 = 1,101010 \cdot 2^2$  ( $2^2$  décale la virgule de 2 chiffres vers la droite)
- La partie fractionnaire étendue sur 23 bits est donc **101 0100 0000 0000 0000**.
- **Exposant** =  $127 + 2 = 129_{10} = 1000\ 0001_2$



En binaire : 11000000 11010100 00000000 00000000  
 En hexadécimal : C0 D4 00 00.

## 3.5. Le code ASCII

ASCII :  
*American  
 Standard  
 Code for  
 Information  
 Interchange*

Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char
0	00	Null	32	20	Space	64	40	@	96	60	`
1	01	Start of heading	33	21	!	65	41	A	97	61	a
2	02	Start of text	34	22	"	66	42	B	98	62	b
3	03	End of text	35	23	#	67	43	C	99	63	c
4	04	End of transmit	36	24	\$	68	44	D	100	64	d
5	05	Enquiry	37	25	%	69	45	E	101	65	e
6	06	Acknowledge	38	26	&	70	46	F	102	66	f
7	07	Audible bell	39	27	'	71	47	G	103	67	g
8	08	Backspace	40	28	(	72	48	H	104	68	h
9	09	Horizontal tab	41	29	)	73	49	I	105	69	i
10	0A	Line feed	42	2A	*	74	4A	J	106	6A	j
11	0B	Vertical tab	43	2B	+	75	4B	K	107	6B	k
12	0C	Form feed	44	2C	,	76	4C	L	108	6C	l
13	0D	Carriage return	45	2D	-	77	4D	M	109	6D	m
14	0E	Shift out	46	2E	.	78	4E	N	110	6E	n
15	0F	Shift in	47	2F	/	79	4F	O	111	6F	o
16	10	Data link escape	48	30	0	80	50	P	112	70	p
17	11	Device control 1	49	31	1	81	51	Q	113	71	q
18	12	Device control 2	50	32	2	82	52	R	114	72	r
19	13	Device control 3	51	33	3	83	53	S	115	73	s
20	14	Device control 4	52	34	4	84	54	T	116	74	t
21	15	Neg. acknowledge	53	35	5	85	55	U	117	75	u
22	16	Synchronous idle	54	36	6	86	56	V	118	76	v
23	17	End trans. block	55	37	7	87	57	W	119	77	w
24	18	Cancel	56	38	8	88	58	X	120	78	x
25	19	End of medium	57	39	9	89	59	Y	121	79	y
26	1A	Substitution	58	3A	:	90	5A	Z	122	7A	z
27	1B	Escape	59	3B	;	91	5B	[	123	7B	{
28	1C	File separator	60	3C	<	92	5C	\	124	7C	
29	1D	Group separator	61	3D	=	93	5D	]	125	7D	}
30	1E	Record separator	62	3E	>	94	5E	^	126	7E	~
31	1F	Unit separator	63	3F	?	95	5F	_	127	7F	□

La norme **ASCII** (on prononce généralement « aski ») établit une correspondance entre une représentation binaire des caractères de l'alphabet latin et les symboles, les signes, qui constituent cet alphabet. Par exemple, le caractère *a* est associé à 1100001 (97) et *A* à 1000001 (65).

## Codage de l'information

La norme ASCII permet ainsi à toutes sortes de machines de stocker, analyser et communiquer de l'information textuelle. En particulier, la quasi-totalité des ordinateurs personnels et des stations de travail utilisent l'encodage ASCII. Le code ASCII de base représentait les caractères sur 7 bits (c'est-à-dire 128 caractères possibles, de 0 à 127).

- Les codes 0 à 31 ne sont pas des caractères. On les appelle caractères de contrôle car ils permettent de faire des actions telles que :
  - retour à la ligne (*Carriage return*)
  - bip sonore (*Audible bell*)
- Les codes 65 à 90 représentent les majuscules
- Les codes 97 à 122 représentent les minuscules (il suffit de modifier le 6ème bit pour passer de majuscules à minuscules, c'est-à-dire ajouter 32 au code ASCII en base décimale).

Le code ASCII a été mis au point pour la langue anglaise, il ne contient donc pas de caractères accentués, ni de caractères spécifiques à une langue. Le code ASCII a donc été étendu à 8 bits pour pouvoir coder plus de caractères (on parle d'ailleurs de code **ASCII étendu**...).

Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char
128	80	Ç	160	A0	á	192	C0	Ł	224	E0	α
129	81	ù	161	A1	í	193	C1	ł	225	E1	β
130	82	é	162	A2	ó	194	C2	ŧ	226	E2	Γ
131	83	â	163	A3	ú	195	C3	ł̇	227	E3	π
132	84	ä	164	A4	ñ	196	C4	—	228	E4	Σ
133	85	à	165	A5	Ñ	197	C5	†	229	E5	σ
134	86	ã	166	A6	ª	198	C6	‡	230	E6	μ
135	87	ç	167	A7	º	199	C7	‡	231	E7	ι
136	88	ê	168	A8	¿	200	C8	Ł̇	232	E8	φ
137	89	ë	169	A9	ƒ	201	C9	Ŧ	233	E9	θ
138	8A	è	170	AA	ƒ̇	202	CA	Ł̈	234	EA	Ω
139	8B	ï	171	AB	½	203	CB	Ŧ̇	235	EB	ϑ
140	8C	î	172	AC	¼	204	CC	‡̇	236	EC	∞
141	8D	ì	173	AD	ı	205	CD	=	237	ED	∞̇
142	8E	Ë	174	AE	«	206	CE	‡̈	238	EE	ε
143	8F	Ā	175	AF	»	207	CF	±	239	EF	∩
144	90	É	176	B0	⋯	208	DO	Ł̈	240	FO	≡
145	91	æ	177	B1	⋮	209	D1	Ŧ̈	241	F1	±̇
146	92	Æ	178	B2	⋭	210	D2	Ŧ̈̇	242	F2	≥
147	93	ó	179	B3		211	D3	Ł̈̈	243	F3	≤̇
148	94	ö	180	B4	†	212	D4	Ł̈̈̇	244	F4	{
149	95	ò	181	B5	‡	213	D5	Ŧ̈̈̇	245	F5	}
150	96	û	182	B6	‡̇	214	D6	Ŧ̈̈̈	246	F6	÷
151	97	ù	183	B7	Ŧ̈̈̈̇	215	D7	‡̈̈̈̇	247	F7	≈
152	98	ÿ	184	B8	Ŧ̈̈̈̈̇	216	D8	‡̈̈̈̈̇	248	F8	°
153	99	ÿ̇	185	B9	‡̈̈̈̈̈̇	217	D9	‡̈̈̈̈̈̇̇	249	F9	•
154	9A	Û	186	BA	‡̈̈̈̈̈̈̇	218	DA	Ŧ̈̈̈̈̈̈̇	250	FA	·
155	9B	ϕ	187	BB	Ŧ̈̈̈̈̈̈̈̇	219	DB	■	251	FB	√
156	9C	£	188	BC	Ł̈̈̈̈̈̈̈̈̇	220	DC	■	252	FC	∞̈̈̈̈̈̈̈̈̇
157	9D	¥	189	BD	Ł̈̈̈̈̈̈̈̈̈̇	221	DD	■	253	FD	∞̈̈̈̈̈̈̈̈̈̇̇
158	9E	€	190	BE	Ł̈̈̈̈̈̈̈̈̈̈̇	222	DE	■	254	FE	■
159	9F	f	191	BF	Ŧ̈̈̈̈̈̈̈̈̈̈̈̇	223	DF	■	255	FF	□

Cette norme s'appelle **ISO-8859** et se décline par exemple en ISO-8859-1 lorsqu'elle étend l'ASCII avec les caractères accentués d'Europe occidentale, et qui est souvent appelée **Latin-1** ou Europe occidentale.



Il existe d'autres normes que l'ASCII, comme l'**Unicode** par exemple, qui présentent l'avantage de proposer une version unifiée des différents encodages de caractères complétant l'ASCII mais aussi de permettre l'encodage de caractères autres que ceux de l'alphabet latin. Unicode définit des dizaines de milliers de codes, mais les 128 premiers restent compatibles avec ASCII.

**UTF-8** (abréviation de l'anglais *Universal Character Set Transformation Format 1 - 8 bits*) est un codage de caractères informatiques conçu pour coder l'ensemble des caractères du « répertoire universel de caractères codés », aujourd'hui totalement compatible avec le standard Unicode, en restant compatible avec la norme ASCII limitée à l'anglais de base, mais très largement répandue depuis des décennies.

Par sa nature, UTF-8 est d'un usage de plus en plus courant sur Internet, et dans les systèmes devant échanger de l'information. Il s'agit également du codage le plus utilisé dans les systèmes GNU, Linux et compatibles pour gérer le plus simplement possible des textes et leurs traductions dans tous les systèmes d'écritures et tous les alphabets du monde.



Toutes ces normes différentes et leurs incompatibilités partielles sont la cause des problèmes que l'on rencontre parfois avec les caractères accentués. C'est pour cette raison qu'il vaut mieux, quand on écrit des courriels à l'étranger, n'utiliser que des caractères non accentués.

## 3.6. Codes détecteurs/correcteurs d'erreurs

Un code correcteur est une technique de codage basée sur la **redondance**. Elle est destinée à corriger les erreurs de transmission d'un message sur une voie de communication peu fiable.

La théorie des codes correcteurs ne se limite pas qu'aux communications classiques (radio, câble coaxial, fibre optique, etc.) mais également aux supports de stockage comme les disques compacts, la mémoire RAM et d'autres applications où l'intégrité des données est importante.

### Comment détecter et/ou corriger des erreurs ?

On peut transmettre un nombre soit en chiffres, soit en lettres :

1. On envoie « 0324614103 ». S'il y a des erreurs de transmission, par exemple si je reçois « 0323614203 », je ne peux pas les détecter.
2. On envoie « zéro trente-deux quatre cent soixante et un quarante et un zéro trois ». S'il y a des erreurs de transmission, par exemple si je reçois « zérb trente-deu quate cent soixante en un quaranhe et on zéro tros », je suis capable de corriger les erreurs et de retrouver le bon numéro.

Dans le premier cas, l'information est la plus concise possible. Dans le deuxième cas au contraire, le message contient plus d'informations que nécessaire. C'est cette redondance qui permet la détection et la correction d'erreurs.

### Pourquoi ces codes ?

- Des canaux de transmission imparfaits entraînant des erreurs lors des échanges de données.
- La probabilité d'erreur sur une ligne téléphonique est de  $10^{-7}$  (cela peut même atteindre  $10^{-4}$ ). Avec un taux d'erreur de  $10^{-6}$  et une connexion à 1 Mo/s, en moyenne 8 bits erronés sont transmis chaque seconde...

### Principe général

- Chaque suite de bits à transmettre est augmentée par une autre suite de bits dite « de redondance » ou « de contrôle ».
- Pour chaque suite de  $k$  bits transmise, on ajoute  $r$  bits. On dit alors que l'on utilise un code  $C(n, k)$  avec  $n = k + r$ .
- À la réception, les bits ajoutés permettent d'effectuer des contrôles.



Richard Hamming  
(1915-1998)

### 3.6.1. La distance de Hamming

La distance de Hamming doit son nom à Richard **Hamming**. Elle est décrite dans un article fondateur pour la théorie des codes. Elle est utilisée en télécommunication pour compter le nombre de bits altérés dans la transmission d'un message d'une longueur donnée.

**Exemple :** la distance de Hamming entre 1011101 et 1001001 est 2 car deux bits sont différents.

Il est souhaitable d'avoir une certaine distance entre les mots envoyés, afin de détecter s'il y a eu une erreur de transmission. Par exemple, si l'on a trois messages à transmettre de trois bits, il vaut mieux choisir les codages qui sont à distance 2 les uns des autres, par exemple 000, 110 et 101. En effet, si un seul bit est altéré, on recevra un message impossible. Par contre, en utilisant 000, 001 et 010, un bit altéré pourrait passer inaperçu.

### 3.6.2. Somme de contrôle

La somme de contrôle (en anglais « checksum ») est un cas particulier de contrôle par redondance. Elle permet de détecter les erreurs, mais pas forcément de les corriger.

Le principe est d'ajouter aux données des éléments dépendant de ces dernières et simples à calculer. Cette redondance accompagne les données lors d'une transmission ou bien lors du stockage sur un support quelconque. Plus tard, il est possible de réaliser la même opération sur les données et de comparer le résultat à la somme de contrôle originale, et ainsi conclure sur la corruption potentielle du message.

#### Bit de parité

Transmettons sept bits auxquels viendra s'ajouter un **bit de parité**. On peut définir le bit de parité comme étant égal à zéro si la somme des autres bits est paire et à un dans le cas contraire. On parle de parité paire. Si la somme des bits est impaire, c'est qu'il y a eu une erreur de transmission.

**Exemple :** 1010001 (7 bits) devient 11010001 (8 bits)

Cette approche permet de détecter les nombres d'erreurs impaires, mais un nombre pair d'erreurs passera inaperçu.

### 3.6.3. Le code ISBN

L'ISBN (*International Standard Book Number*) est un numéro international qui permet d'identifier, de manière unique, chaque livre publié. Il est destiné à simplifier la gestion informatique des livres dans les bibliothèques, librairies, etc.



Le numéro ISBN-10 se compose de quatre segments, trois segments de longueur variable et un segment de longueur fixe, la longueur totale de l'ISBN comprend dix chiffres (le 1<sup>er</sup> janvier 2007, la longueur a été étendue à 13 chiffres en ajoutant un groupe initial de 3 chiffres).

Si les quatre segments d'un ancien code ISBN à 10 chiffres sont notés A - B - C - D :

- A identifie un groupe de codes pour un pays, une zone géographique ou une zone de langue.
- B identifie l'éditeur de la publication.
- C correspond au numéro d'ordre de l'ouvrage chez l'éditeur.
- D est un chiffre-clé calculé à partir des chiffres précédents et qui permet de vérifier qu'il n'y a pas d'erreurs. Outre les chiffres de 0 à 9, cette clé de contrôle peut prendre la valeur X, qui représente le nombre 10.

#### Calcul du chiffre-clé d'un numéro ISBN-10

- On attribue une pondération à chaque position (de 10 à 2 en allant en sens décroissant) et on fait la somme des produits ainsi obtenus.

- On conserve le reste de la division euclidienne de ce nombre par 11. La clé s'obtient en retranchant ce nombre à 11. Si le reste de la division euclidienne est 0, la clé de contrôle n'est pas 11 ( $11 - 0 = 11$ ) mais 0.
- De même, si le reste de la division euclidienne est 1, la clé de contrôle n'est pas 10 mais la lettre X.

Le nombre 11 étant premier, une erreur portant sur un chiffre entraînera automatiquement une incohérence du code de contrôle. La vérification du code de contrôle peut se faire en effectuant le même calcul sur le code ISBN complet, en appliquant la pondération 1 au dixième chiffre de la clé de contrôle (si ce chiffre-clé est X, on lui attribue la valeur 10) : la somme pondérée doit alors être un multiple de 11.

### Exemple

Pour le numéro ISBN (à 9 chiffres) 2-940043-41, quelle est la clé de contrôle ?

Code ISBN	2	9	4	0	0	4	3	4	1
Pondération	10	9	8	7	6	5	4	3	2
Produit	20	81	32	0	0	20	12	12	2

**La question qui tue :**  
pourquoi est-il indispensable de donner une pondération à chaque position ?

La somme des produits est **179**, dont le reste de la division euclidienne par 11 est **3**. La clé de contrôle est donc  $11 - 3 = 8$ . L'ISBN au complet est : 2-940043-41-8.

La vérification de la clé complète à 10 chiffres donne la somme pondérée  $179 + 8 = 187$ , qui est bien un multiple de 11.

### Depuis le 1<sup>er</sup> janvier 2007, les ISBN sont passés à 13 chiffres

#### Pourquoi une réforme du système ISBN ?

- Pour augmenter la capacité de numérotation du système ISBN qui était devenu insuffisant, notamment en raison de l'augmentation du nombre de publications électroniques.
- Pour rendre l'ISBN complètement compatible avec l'EAN-13 qui sert à la génération des codes-barres (voir § 3.8.1) :

Le **code EAN** (*European Article Numbering*) est un code-barres utilisé par le commerce et l'industrie conformément aux spécifications d'*EAN International*.



#### Qu'est-ce qui a changé ?

- On fait précéder l'ISBN du préfixe « 978 » et on recalcule la clé de contrôle.
- L'ISBN à 13 chiffres est identique à l'EAN-13 servant à la génération du code-barres et figurant sous celui-ci.
- Le préfixe « 979 » sera introduit quand les segments ISBN existants seront épuisés.
- Les segments identifiant les éditeurs ne demeureront pas les mêmes lorsque le préfixe passera en 979.
- Depuis le 1<sup>er</sup> janvier 2007, l'ISBN à 13 chiffres est utilisé pour toutes les commandes et toutes les transactions commerciales, manuelles ou électroniques.
- Depuis le 1<sup>er</sup> janvier 2007, les codes-barres sont surmontés de l'ISBN à 13 chiffres segmenté (avec des tirets) alors que l'EAN-13 correspondant est indiqué en dessous des codes à barres (sans tiret ni espace).

#### Algorithme permettant de générer l'EAN Bookland à partir de l'ISBN

1. Garder les 9 premiers chiffres de l'ISBN.
2. Ajouter le préfixe Bookland de l'EAN.
3. Entrer/lire les facteurs de pondération constants associés à chaque position de l'EAN (1 3 1

3 1 3 1 3 1 3 1 3).

4. Multiplier chaque chiffre par le facteur de pondération qui lui est associé.
5. Diviser la somme par le nombre du module (10) pour trouver le reste.
6. Si le reste vaut 0, le numéro de contrôle est aussi 0. Sinon, le numéro de contrôle est 10 moins le reste trouvé. Ajouter ce suffixe.

Convertissons l'ISBN 0-8436-1072-7 :

ISBN				0	8	4	3	6	1	0	7	2	7
Ajout du préfixe 978	9	7	8	0	8	4	3	6	1	0	7	2	
Facteurs	1	3	1	3	1	3	1	3	1	3	1	3	
Produits	9	21	8	0	8	12	3	18	1	0	7	6	

$$9 + 21 + 8 + 0 + 8 + 12 + 3 + 18 + 1 + 7 + 6 = 93.$$

Le reste de la division par 10 est 3. Le chiffre de contrôle est donc  $10 - 3 = 7$ . Le code EAN-ISBN est donc : 978-0-8436-1072-7.

### Exercice 3.9

Quel est le numéro de contrôle du code ISBN à 9 chiffres 2-35288-041 ?

Donnez le code EAN Bookland à partir de l'ISBN ci-dessus.

### 3.6.4. Formule de Luhn

Cette formule génère un chiffre de vérification, qui est généralement annexé à un numéro d'identité partiel pour générer un identifiant complet.

Cet identifiant (numéro partiel + chiffre de vérification) est soumis à l'algorithme suivant pour vérifier sa validité :

1. On alterne les facteurs de pondération 1 et 2 sous les chiffres du code, **en les écrivant de droite à gauche.**
2. Chaque chiffre du code est multiplié par son facteur de pondération. Si un chiffre multiplié par 2 est plus grand que 9 (comme c'est le cas par exemple pour 8 qui devient 16), alors il faut le ramener à un chiffre entre 1 et 9 en soustrayant 9.
3. On effectue la somme de tous les produits obtenus.
4. Si cette somme est un multiple de 10, alors le nombre original est valide.

#### Exemple

Considérons l'identification du nombre 972-487-086. La première étape consiste à doubler un chiffre sur deux en partant de l'avant-dernier jusqu'au début, et de faire la somme de tous les chiffres, doublés ou non (si un chiffre est supérieur à 9, on retranche 9, d'où la 4<sup>ème</sup> ligne). Le tableau suivant montre cette étape :

Code	9	7	2	4	8	7	0	8	6
Facteurs	1	2	1	2	1	2	1	2	1
Produits	9	14	2	8	8	14	0	16	6
Chiffres entre 0 et 9	9	5	2	8	8	5	0	7	6
									<b>Somme=50</b>

La somme, égale à 50, est un multiple de 10, donc le nombre est valide.

Chaque véhicule ferroviaire dispose d'un numéro d'identification unique le distinguant de tout autre véhicule ferroviaire. La numérotation des « plaques d'immatriculation » des trains a été uniformisée par l'**Union internationale des chemins de fer (UIC)** : chaque locomotive, chaque automotrice et chaque voiture de voyageurs est identifiée par un numéro à douze chiffres.

Hans Peter **Luhn** (1896-1964) est un informaticien allemand qui a travaillé pour la société IBM.



Le numéro écrit sur la voiture de voyageurs ci-dessus indique (voir [6]) :

- le type de véhicule : 50 = service commercial national
- le pays : 85 = Suisse
- le type de voiture : 3- = voiture mixte 1<sup>ère</sup> – 2<sup>e</sup> classe, -9 = neuf compartiments
- la vitesse maximale : 43 = vitesse maximale comprise entre 121 et 140 km/h
- le numéro de série : 829
- le chiffre 3 situé après le tiret est une clé de contrôle pour la vérification informatique.

### Exercice 3.10

Vérifiez que ce numéro d'immatriculation satisfait la condition de la formule de Luhn.

### 3.6.5. Vérification des CCP

CCP signifie Compte Courant Postal. Pour vérifier qu'un numéro de CCP est valide, on va utiliser le tableau ci-dessous, tiré de la documentation de Postfinance :

retenue	0	1	2	3	4	5	6	7	8	9	contrôle
0	0	9	4	6	8	2	7	1	3	5	0
1	9	4	6	8	2	7	1	3	5	0	9
2	4	6	8	2	7	1	3	5	0	9	8
3	6	8	2	7	1	3	5	0	9	4	7
4	8	2	7	1	3	5	0	9	4	6	6
5	2	7	1	3	5	0	9	4	6	8	5
6	7	1	3	5	0	9	4	6	8	2	4
7	1	3	5	0	9	4	6	8	2	7	3
8	3	5	0	9	4	6	8	2	7	1	2
9	5	0	9	4	6	8	2	7	1	3	1

Vérifions le numéro de CCP 17-603303-5.

1. On commence toujours avec une retenue de 0.
2. Le premier chiffre du CCP est 1. Le tableau indique que la nouvelle retenue est 9 (intersection de la ligne « 0 » et de la colonne « 1 »).
3. Avec une retenue de 9 et un deuxième chiffre égal à 7, on trouve que la deuxième retenue est 7 (intersection de la ligne « 9 » et de la colonne « 7 »).
4. On continue ainsi jusqu'à l'avant-dernier chiffre du CCP (le dernier chiffre est la clé de

contrôle qui se lira dans la colonne « contrôle ») :

Retenues	Chiffres du CCP
0	1
9	7
7	6
6	0
7	3
0	3
6	0
7	3
0	0

La dernière retenue est 0. Cela correspond dans la colonne rose à la clé de contrôle 0.

Le dernier chiffre du CCP devrait donc être un 0 et non pas un 5, ce qui signifie que le numéro de CCP est incorrect.

### Exercice 3.11

Vérifiez le numéro de CCP 10-546323-6.

## 3.6.6. Code de Hamming

Un code de Hamming permet la **détection** et la **correction** automatique d'**une** erreur si elle ne porte que sur un bit du message. Un code de Hamming est **parfait**, ce qui signifie que pour une longueur de code donnée, il n'existe pas d'autre code plus compact ayant la même capacité de correction. En ce sens, son rendement est maximal.

### Structure d'un code de Hamming

- les  $m$  bits du message à transmettre et les  $n$  bits de contrôle de parité.
- longueur totale :  $2^n - 1$
- longueur du message :  $m = (2^n - 1) - n$
- on parle de code  $x-y$  :  $x$  est la longueur totale du code ( $n+m$ ) et  $y$  la longueur du message ( $m$ )
- les bits de contrôle de parité  $C_i$  sont en position  $2^i$  pour  $i = 0, 1, 2, \dots$
- les bits du message  $D_j$  occupent le reste du message.

7	6	5	4	3	2	1
$D_3$	$D_2$	$D_1$	$C_2$	$D_0$	$C_1$	$C_0$

Structure d'un code de Hamming 7-4

### Exemples de code de Hamming

- un mot de code 7-4 a un coefficient d'efficacité de  $4/7 = 57\%$
- un mot de code 15-11 a un coefficient d'efficacité de  $11/15 = 73\%$
- un mot de code 31-26 a un coefficient d'efficacité de  $26/31 = 83\%$

### Retrouver l'erreur dans un mot de Hamming

Si les bits de contrôle de parité  $C_2, C_1, C_0$  ont tous la bonne valeur, il n'y a pas d'erreur ; sinon la valeur des bits de contrôle indique la position de l'erreur entre 1 et 7.

Le code de Hamming présenté ici ne permet de retrouver et corriger qu'une erreur.

Pour savoir quels bits sont vérifiés par chacun des bits de contrôle de parité, on peut construire le tableau ci-après.

	7	6	5	4	3	2	1
$C_2$	1	1	1	1	0	0	0
$C_1$	1	1	0	0	1	1	0
$C_0$	1	0	1	0	1	0	1

- La première ligne indique la position des bits.
- Dans chaque colonne, on écrit le numéro de la position verticalement en binaire.
- Les 1 indiquent alors quels bits sont vérifiés. Ainsi :
  - $C_2$  vérifie les bits 4, 5, 6, 7. La somme de ces bits doit être paire.
  - $C_1$  vérifie les bits 2, 3, 6, 7. La somme de ces bits doit être paire.
  - $C_0$  vérifie les bits 1, 3, 5, 7. La somme de ces bits doit être paire.

### Exemple d'un code de Hamming 7-4

On souhaite envoyer le message 1010. Complétons le mot de Hamming correspondant :

7	6	5	4	3	2	1
1	0	1		0		

- $C_0$  vaut 0 pour pouvoir rendre pair  $1+1+0$  (les bits d'indices 7, 5, 3).
- $C_1$  vaut 1 pour pouvoir rendre pair  $1+0+0$  (les bits d'indices 7, 6, 3).
- $C_2$  vaut 0 pour pouvoir rendre pair  $1+0+1$  (les bits d'indices 7, 6, 5).

7	6	5	4	3	2	1
1	0	1	0	0	1	0

Imaginons que l'on reçoive le mot 0010010 (le bit de poids fort a été altéré).

$C_0$  a la mauvaise valeur, car  $0+1+0+0$  est impair, donc il y a une erreur en position 7, 5, 3 ou 1.

$C_1$  a la mauvaise valeur, car  $0+0+0+1$  est impair, donc il y a une erreur en position 7, 6, 3 ou 2.

$C_2$  a la mauvaise valeur, car  $0+0+1+0$  est impair, donc il y a une erreur en position 7, 6, 5 ou 4.

Écrivons le nombre binaire  $C_2C_1C_0$  où  $C_i$  vaut 0 si le bit de contrôle  $C_i$  a la bonne valeur et 1 sinon. On obtient ici 111, ce qui correspond à 7 en binaire. Le bit erroné est le numéro 7.

Que se passe-t-il si c'est un des bits de contrôle qui est altéré ? Imaginons que l'on ait reçu 1010011 (cette fois-ci, c'est le bit de poids faible qui a été altéré).

$C_0$  a la mauvaise valeur, car  $1+1+0+1$  est impair. Il y a une erreur en position 7, 5, 3 ou 1.

$C_1$  a la bonne valeur, car  $1+0+0+1$  est pair. Il n'y a pas d'erreur en position 7, 6, 3 et 2.

$C_2$  a la bonne valeur, car  $1+0+1+0$  est pair. Il n'y a pas d'erreur en position 7, 6, 5 et 4.

Ici,  $C_2C_1C_0$  vaut 001. Le bit erroné est donc le numéro 1.

### Exercice 3.12

Vous voulez envoyer le mot 1011. Quels bits de contrôle devez-vous lui adjoindre et quelle séquence transmettez-vous alors ?

### Exercice 3.13

Y a-t-il une erreur dans le mot suivant (Hamming 7-4) : 1101101 ?

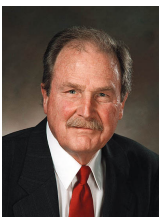
**Exercice 3.14**

Soit un mot de Hamming 15–11 suivant :

15	14	13	12	11	10	9	8	7	6	5	4	3	2	1
1	0	1	1	0	1	1	1	1	0	1	1	0	1	1

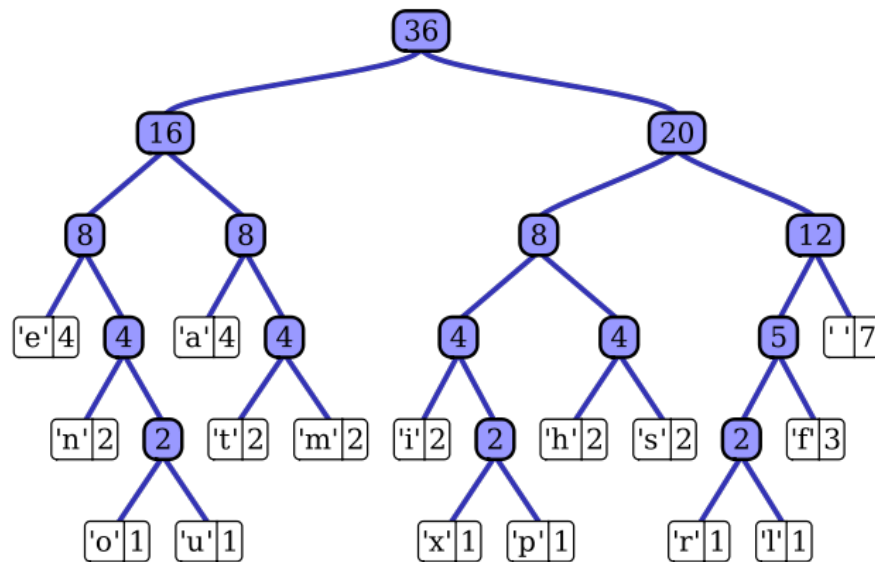
1. Quels sont les bits de contrôle de parité ?
2. Quelles positions contrôlent chacun de ces bits ?
3. Quel est le message reçu ?
4. Est-ce que le message reçu correspond au message transmis ?

### 3.7. Codage de Huffman



David A. Huffman  
(1925-1999)

Le codage de Huffman (1952) est une méthode de **compression statistique de données** qui permet de réduire la longueur du codage d'un alphabet. Cela signifie que l'on va d'abord calculer les fréquences des caractères du texte (ou les couleurs d'une image). Puis on va construire un arbre binaire (voir ci-dessous) dont les extrémités seront les lettres. Les lettres les plus fréquentes seront en haut de l'arbre, les moins fréquentes en bas. Le codage sera obtenu en parcourant l'arbre depuis le haut en allant vers la lettre : un mouvement à gauche sera codé 0 et un mouvement à droite sera codé 1. Ainsi, dans l'arbre ci-dessous, « a » sera codé « 010 », et « x » « 10010 ».



Un exemple d'arbre obtenu avec la phrase « this is an example of a huffman tree »

#### Le principe

Huffman propose de **recoder les données qui ont une occurrence très faible sur une longueur binaire supérieure à la moyenne, et recoder les données très fréquentes sur une longueur binaire très courte.**

Ainsi, pour les données rares, nous perdons quelques bits gagnés pour les données répétitives. Par exemple, dans un fichier ASCII le « w » apparaissant 10 fois aura un code très long : 0101000001000. Ici la perte est de 40 bits (10 x 4 bits), car sans compression, il serait codé sur 8 bits au lieu de 12. Par contre, le caractère le plus fréquent comme le « e » avec 200 apparitions sera codé par 1. Le gain sera de 1400 bits (7 x 200 bits). On comprend l'intérêt d'une telle méthode.

Le codage de Huffman a une **propriété de préfixe** : une séquence binaire ne peut jamais être à la fois représentative d'un élément codé et constituer le début du code d'un autre élément. Si un

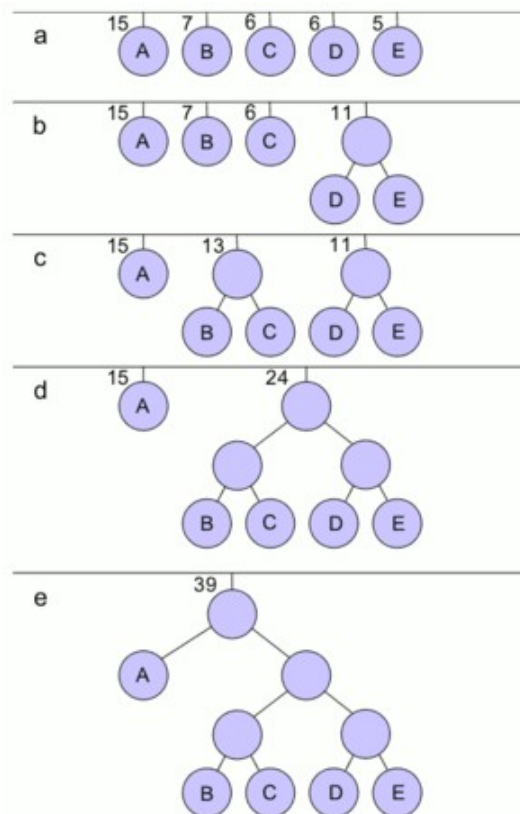


caractère est représenté par la combinaison binaire 100 alors la combinaison 10001 ne peut être le code d'aucune autre information. Dans ce cas, l'algorithme de décodage interpréterait les 5 bits comme deux mots : 100-01. Cette caractéristique du codage de Huffman permet une codification à l'aide d'une structure d'arbre binaire.

**Méthode**

Au départ, chaque lettre a une étiquette correspondant à sa fréquence (ou sa probabilité) d'apparition. Il y a autant d'arbres (à 1 sommet) que de lettres.

L'algorithme choisit à chaque étape les deux arbres d'étiquettes minimales, soit  $x$  et  $y$ , et les remplace par l'arbre formé de  $x$  et  $y$  et ayant comme étiquette la somme de l'étiquette de  $x$  et de l'étiquette de  $y$ . La figure ci-dessous représente les étapes de la construction d'un code de Huffman pour l'alphabet source  $\{A, B, C, D, E\}$ , avec les fréquences  $F(A)=15, F(B)=7, F(C)=6, F(D)=6$  et  $F(E)=5$ .



Le code d'une lettre est alors déterminé en suivant le chemin depuis la racine de l'arbre jusqu'à la feuille associée à cette lettre en concaténant successivement un 0 ou un 1 selon que la branche suivie est à gauche ou à droite. Ainsi, sur la figure ci-dessus,  $A=0, B=100, C=101, D=110, E=111$ .

Par exemple, le mot « ABBE » serait codé 0100100111. Pour décoder, on lit simplement la chaîne de bits de gauche à droite. Le seul découpage possible, grâce à la propriété du préfixe, est 0-100-100-111.

Ce principe de compression est aussi utilisé dans le codage d'image TIFF (Tagged Image Format File) spécifié par *Microsoft Corporation* et *Aldus Corporation*.

Il existe des méthodes qui ne conservent pas exactement le contenu d'une image (méthodes non conservatives) mais dont la représentation visuelle reste correcte. Entre autres, il y a la méthode JPEG (Join Photographic Experts Group) qui utilise la compression de type Huffman pour coder les informations d'une image.

Malgré son ancienneté, cette méthode est toujours remise au goût du jour, et offre des performances appréciables.

**Exercice 3.15**

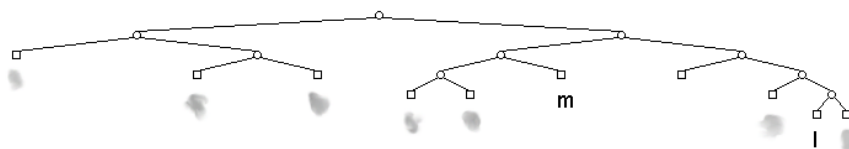
Construisez un codage de Huffman du message « CECI EST UN CODAGE DE HUFFMAN. ». Notez qu'il y a plusieurs codages de Huffman possibles, mais la longueur du message codé sera toujours la même. N'oubliez pas les espaces et le point.  
Vérifiez la propriété du préfixe.

**Exercice 3.16**

De quel roman célèbre est tirée la phrase compressée ci-dessous ?

1111101100001101010101101111010001001101011010001110

L'arbre de Huffman (malheureusement en partie effacé) est le suivant :



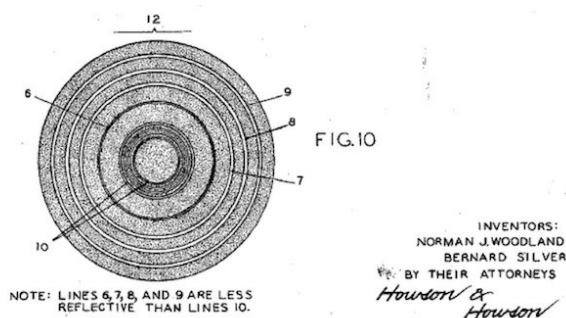
Voici la fréquence d'apparition des lettres :

a	C	e	h	l (i)	l (L)	m	s		.
2	1	2	1	1	3	2	1	2	1

### 3.8. Codes-barres

Un code-barres est la représentation d'une donnée numérique ou alphanumérique sous forme d'un symbole constitué de barres et d'espaces dont l'épaisseur varie en fonction des données ainsi codées. Ils sont destinés à une lecture automatisée par un capteur électronique.

Le premier code-barres était... circulaire et date de 1949. Voici le brevet déposé par ses inventeurs Norman J. Woodland et Bernard Silver :



Il ne sera jamais utilisé dans l'industrie. Le premier scanner permettant de lire des code-barres ne vit le jour qu'au début des années 1970 et les ordinateurs nécessaires au traitement des données n'en étaient qu'à leur balbutiements.

#### 3.8.1. Codes-barres EAN

Le **numéro EAN** (*European Article Numbering*) identifie des articles ou des unités logistiques de façon unique, codé sous forme de codes-barres. L'EAN est composé de 8 ou 13 chiffres représentés sous forme de séquences de barres noires et blanches formant un code-barres. Ce type de code-barres se trouve sur la presque totalité de produits courants (alimentation, vêtements, droguerie, papeterie,

électroménager, etc.). Le code est lu lors du passage aux caisses des commerces.

Il existe des codes EAN-8 des codes EAN-13, composés respectivement de 8 ou 13 chiffres :

- les codes EAN-8 sont réservés à l'usage sur des produits de petite taille ;
- les codes EAN-13 sont utilisés sur tous les autres produits.

Le système des codes EAN, comme tous les systèmes de codes-barres, fait appel à des notions d'arithmétique modulaire. Sa structure tient compte des contraintes physiques liées aux conditions de leur impression et de leur lecture. En effet, la reconnaissance d'un code-barres nécessite :

- de pouvoir séparer et mesurer les largeurs des barres, à des distances de lecture variable, sur différents types de capteurs et en l'absence de toute horloge ou mesure de référence,
- de pouvoir le faire en outre dans n'importe quel sens de lecture et quelles que soient les couleurs effectivement utilisées.

### Structure du code EAN-13

- Le caractère de **Début** codé 101 (1 = noir, 0 = blanc)
- Le second caractère du **Préfixe** (6). Le premier caractère du préfixe n'est pas codé.
- Les cinq caractères du **Numéro de Participant** (12345)
- Le **Séparateur Central** est codé 01010
- Les cinq caractères du **Numéro d'Article** (67890)
- Le **Check Digit** (0)
- Le caractère de **Fin** codé 101.



Les caractères situés à gauche du Séparateur Central du symbole EAN 13 utilisent deux jeux de codification nommés élément A et élément B. Ceux situés à droite utilisent le jeu de codification nommé élément C. Le premier caractère du préfixe détermine l'alternance des éléments A et B à utiliser pour le codage des 6 caractères situés à gauche du séparateur central.

### Table de codage des caractères EAN-13

Chacun des chiffres composant l'EAN-13 peut, selon sa position dans le code, avoir **trois représentations distinctes** nommées. La table ci-dessous indique comment coder chaque caractère d'un EAN 13 selon qu'il se trouve à gauche ou à droite du Séparateur Central. Les caractères de gauche utilisant le Set A ou B en fonction de la valeur du premier caractère du Préfixe.

Chaque élément peut être représenté en binaire par une suite de 7 bits :

- un X ou 1 correspondant à une barre élémentaire noire,
- un \_ ou 0 correspondant à une barre élémentaire blanche.

Voici les représentations des 10 chiffres comme éléments A, B ou C :

	élément A	élément B	élément C		élément A	élément B	élément C
0	[ _ XX X ]	[ _ X XXX ]	[ XXX _ X ]	soit, graphiquement			
1	[ XX _ X ]	[ XX XX ]	[ XX XX ]				
2	[ _ X XX ]	[ XX XX ]	[ XX XX ]				
3	[ XXXX _ X ]	[ X _ X ]	[ X _ X ]				
4	[ X _ XX ]	[ XXXX _ X ]	[ X XXX ]				
5	[ XX _ X ]	[ XXXX _ X ]	[ X XXX ]				
6	[ X XXXX ]	[ _ X X ]	[ X X ]				
7	[ XXXX XX ]	[ X _ X ]	[ X X ]				
8	[ XX XXX ]	[ _ X X ]	[ X X ]				
9	[ _ X XX ]	[ X XXX ]	[ XXX _ X ]				

Les éléments ci-dessus peuvent aussi être représentés sous la forme équivalente suivante où seules les largeurs des quatre barres successives (décrites de gauche à droite) sont indiquées (les éléments A et C se distinguent uniquement par la couleur de la première barre) :

chiffre	élt A ou C	élt B
0	3211	1123
1	2221	1222
2	2122	2212
3	1411	1141
4	1132	2311
5	1231	1321
6	1114	4111
7	1312	2131
8	1213	3121
9	3112	2113



2 : 2122 (A), 3 : 1411 (A), 4 : 2311 (B), 5 : 1231 (A),  
 6 : 4111 (B), 7 : 2131 (B), 8 : 1213 (C), 9 : 3112 (C),  
 1 : 2221 (C), 2 : 2122 (C), 4 : 1132 (C), 8 : 1213 (C).

Dans le code-barres ci-dessus, le préfixe commence par 1, ce qui donne l'alternance AABABB. Les autres alternances en fonction du premier chiffre du préfixe sont données ci-après.

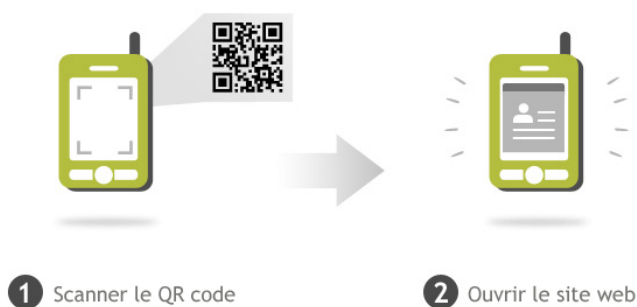
1	AABABB	2	AABBAB	3	AABBBA
4	ABAABB	5	ABBAAB	6	ABBBAA
7	ABABAB	8	ABABBA	9	ABBABA
0	AAAAAA (UPC-A)				

### Remarques

- Les codages B et C d'un même chiffre (donc situés ci-dessus sur la même ligne du premier tableau) sont toujours symétriques l'un de l'autre.
  - Il en résulte que si le code est lu à l'envers, un élément de type B apparaît comme un élément de type C et un élément de type C apparaît de type B.
  - Par contre, un élément de type A lu à l'envers est incorrect. Il y a donc toujours au moins un élément de type A dans le code afin de déterminer le sens de lecture : pour le code EAN-8 ce sont tous les éléments de gauche, et pour le code EAN-13, le deuxième élément est toujours de type A.
- Le codage d'un même chiffre comme élément C est toujours le complément de son codage comme élément A (les 1 et les 0 sont permutées).
  - Il en résulte que si un élément qui devrait être de type A est en fait de type C, c'est que les couleurs sont inversées : blanc  $\rightarrow$  1 et noir  $\rightarrow$  0
- Il découle des remarques précédentes que le codage EAN permet de lire le code :
  - quel que soit le sens de lecture (c'est pourquoi on peut, à une caisse de supermarché, présenter un produit dans n'importe quelle position) ;
  - quelle que soit la couleur d'impression du code-barres (blanc sur fond noir, ou noir sur fond blanc par exemple).
- Le codage EAN peut permettre de différencier 10 caractères distincts (ici, ce sont les dix chiffres de 0 à 9) ; en revanche, il ne pourrait pas servir à coder des caractères supplémentaires (donc pas de lettre).

### 3.8.2. Codes QR

Le code QR (ou *QR code* en anglais) est un code-barres en deux dimensions (ou code à matrice) constitué de modules noirs disposés dans un carré à fond blanc. Ce cours contient un QR code à chaque début de chapitre. Le nom *QR* est l'acronyme de l'anglais *Quick Response*, car son contenu de données peut être décodé rapidement.

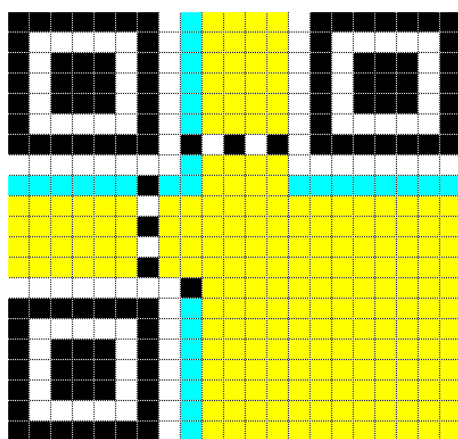


Le code QR a été créé par l'entreprise japonaise Denso-Wave en 1994 pour le suivi des pièces de voiture dans les usines de Toyota.

Les codes QR peuvent mémoriser des adresses web, du texte, des numéros de téléphone, des SMS ou autres types de données lisibles par les smartphones et les téléphones mobiles équipés d'une application de lecture (lecteur de code QR ou *QR reader* en anglais).

Ils peuvent stocker jusqu'à 7089 caractères numériques, 4296 caractères alphanumériques ou 2953 octets. Par rapport au code-barres traditionnel qui ne peut stocker que de 10 à 13 caractères, ils ont l'avantage de pouvoir stocker beaucoup d'informations tout en étant petits et rapides à scanner.

Le code QR est défini comme un standard ISO (IEC18004).



Les cases noires et blanches servent à l'orientation du code QR. Grâce à ces trois carrés, l'image peut être lue dans tous les sens.

Dans la partie cyan se trouvent des informations sur le format.

Les informations sont codées dans la partie jaune.

L'avantage du code QR est sa facilité et rapidité d'utilisation et de création. Pour lire un code QR, il suffit de lancer l'application de lecture et viser le code dans le mobile. En ce qui concerne l'écriture, il y a plusieurs sites web qui permettent de générer librement les codes QR.

Les codes QR utilisent le système Reed-Solomon pour la correction d'erreur : le code contient jusqu'à 30 % de redondance. Ainsi, on peut modifier un code QR pour le rendre plus artistique et il restera lisible :



## Sources

- [1] Dumas, Roch, Tannier, Varrette, *Théorie des codes : Compression, cryptage, correction*, Dunod, 2006
- [2] Martin B., *Codage, cryptologie et applications*, Presses Polytechniques et Universitaires Romandes (PPUR), 2004
- [3] Comment ça marche, « Représentation des nombres entiers et réels », <<http://www.commentcamarche.net/contents/base/representation.php3>>
- [4] Wikipédia, « International Book Standard Number », <<http://fr.wikipedia.org/wiki/ISBN>>
- [5] Wikipédia, « Code-barres EAN », <[https://fr.wikipedia.org/wiki/Code-barres\\_EAN](https://fr.wikipedia.org/wiki/Code-barres_EAN)>
- [6] Wikipédia, « Classification UIC des voitures de chemin de fer » <[https://fr.wikipedia.org/wiki/Classification\\_UIC\\_des\\_voitures\\_de\\_chemin\\_de\\_fer](https://fr.wikipedia.org/wiki/Classification_UIC_des_voitures_de_chemin_de_fer)>
- [7] Duvallet Claude, « Les codes correcteurs et les codes détecteurs d'erreurs », <<http://litis.univ-lehavre.fr/~duvallet/enseignements/Cours/LPRODA21/UF9/LPRODA21-TD2-UF9.pdf>>
- [8] Wikipédia, « Codage de Huffman », <[http://fr.wikipedia.org/wiki/Codage\\_de\\_Huffman](http://fr.wikipedia.org/wiki/Codage_de_Huffman)>
- [9] Wikipédia, « Code QR », <[https://fr.wikipedia.org/wiki/Code\\_QR](https://fr.wikipedia.org/wiki/Code_QR)>

