



L'encodage de texte

1. Qu'est-ce qu'un encodage ?

Le jeu de caractères

Le codage des caractères est une convention qui permet, à travers un codage connu de tous, de transmettre de l'information textuelle. On définit alors un jeu de caractères (charset en anglais) appelé aussi répertoire.

A chaque caractère est assigné un numéro unique appelé code. On obtient alors un jeu de caractères codés et on peut alors créer un tableau de correspondance. Par exemple¹:

Α	В	•••	Z	0	1	 9		,	٤	?	!	_	espac e
0	1		25	26	27	35	36	37	38	39	40	41	42

Le codage du message "LA NSI, C'EST BIEN!" peut s'effectuer par l'utilisation de la table de correspondance :

L	А	es p	N	S	I	es p	С	6	Е	S	Т	es p	В	I	E	N	es p	!
11	0	42	13	18	8	42	2	38	4	18	19	42	1	8	4	13	42	40

Cependant attention, l'écriture du message ne peut pas être effectuée comme cela : 1104213188422384181942184134240

Il y a en effet ambiguité sur la longueur des caractères. La première lettre est-elle un B (1) ou un L (11) ? Et si la première lettre est un B (1) la deuxième est elle encore un B (1) ou un K (10)? etc...

Il est donc nécessaire de mettre en place un encodage. Ici un encodage simple est de coder l'information sur deux digits pour éviter toute confusion. Ainsi le message devient :

11004213180842023804181942010804134240

De cette même façon, vous êtes donc en mesure de décoder le message précédent :

01141309142017421104184204110421041836 -> BONJOUR LES ELEVES.

Ce simple exemple montre qu'une stratégie d'encodage est nécessaire.

Repères historiques

On n'a pas pas attendu l'informatique pour transmettre des messages autrement que par écrit. Le code Morse en est un parfait exemple (figure 1). Le braille est également un autre exemple (figure 2).

¹ https://zestedesavoir.com/tutoriels/1114/comprendre-les-encodages/1-theorie/



3- Encodage d'un texte





Code morse international

- 1. Un tiret est égal à trois points.
- 2. L'espacement entre deux éléments d'une même lettre est égal à un point.
- 3. L'espacement entre deux lettres est égal à trois points.
- 4. L'espacement entre deux mots est égal à sept points.

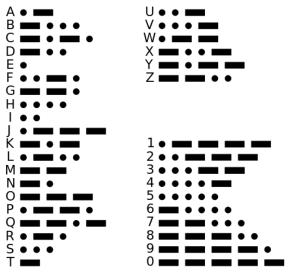


Figure 1: Table d'encodage du morse

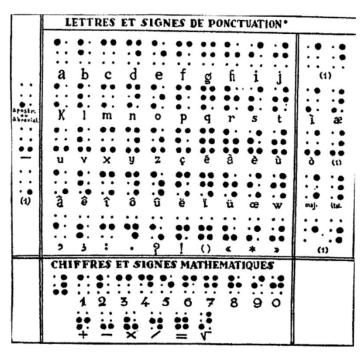


Figure 2: Table d'encodage du braille

Question 1) **Calculer** le nombre de caractères codable par les 6 points de touche d'une cellule braille

Il y a 6 caractères prenant 2 états possible (base 2). Il y a donc 2⁶ caractères codables sur une cellule braille

2.<u>Les encodages utilisés en informatique</u>

Le code ASCII

des encodages premiers historiques est l'ASCII (American Standard Code for Information Interchange, en français, le code américain normalisé pour l'échange d'informations). Cet encodage s'est imposé au début des années 1960. Le code ASCII (figure 3) contient 128 caractères qui suffisent à écrire des textes en anglais. Cependant ce code ne convient pas pour les autres langues.

Dec	Hex	Name	Char	Ctrl-char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char
0	0	Null	NUL	CTRL-@	32	20	Space	64	40	0	96	60	
1	1	Start of heading	SOH	CTRL-A	33	21	1	65	41	A	97	61	a
2	2	Start of text	STX	CTRL-B	34	22		66	42	В	98	62	b
3	3	End of text	ETX	CTRL-C	35	23	#	67	43	C	99	63	c
4	4	End of xmit	EOT	CTRL-D	36	24	\$	68	44	D	100	64	d
5	5	Enquiry	ENQ	CTRL-E	37	25	%	69	45	E	101	65	е
6	6	Acknowledge	ACK	CTRL-F	38	26	8.	70	46	F	102	66	f
7	7	Bell	BEL	CTRL-G	39	27		71	47	G	103	67	g
8	8	Backspace	BS	CTRL-H	40	28	(72	48	Н	104	68	h
9	9	Horizontal tab	HT	CTRL-I	41	29)	73	49	I	105	69	i
10	OA.	Line feed	LF	CTRL-J	42	2A		74	44	J	106	6A	j
11	OB	Vertical tab	VT	CTRL-K	43	2B	+	75	4B	K	107	6B	k
12	OC.	Form feed	FF	CTRL-L	44	2C	F	76	4C	L	108	6C	1
13	OD	Carriage feed	CR	CTRL-M	45	2D	-	77	4D	M	109	6D	m
14	Œ	Shift out	so	CTRL-N	46	2E		78	4E	N	110	6E	n
15	0F	Shift in	SI	CTRL-O	47	2F	1	79	4F	0	111	6F	0
16	10	Data line escape	DLE	CTRL-P	48	30	0	80	50	P	112	70	р
17	11	Device control 1	DC1	CTRL-Q	49	31	1	81	51	Q	113	71	q
18	12	Device control 2	DC2	CTRL-R	50	32	2	82	52	R	114	72	r
19	13	Device control 3	DC3	CTRL-S	51	33	3	83	53	S	115	73	s
20	14	Device control 4	DC4	CTRL-T	52	34	4	84	54	T	116	74	t
21	15	Neg acknowledge	NAK	CTRL-U	53	35	5	85	55	U	117	75	u
22	16	Synchronous idle	SYN	CTRL-V	54	36	6	86	56	V	118	76	٧
23	17	End of xmit block	ETB	CTRL-W	55	37	7	87	57	W	119	77	w
24	18	Cancel	CAN	CTRL-X	56	38	8	88	58	×	120	78	×
25	19	End of medium	EM	CTRL-Y	57	39	9	89	59	Y	121	79	y
26	1A	Substitute	SUB	CTRL-Z	58	ЗА	:	90	5A	Z	122	7A	z
27	18	Escape	ESC	CTRL-[59	38	;	91	5B	[123	7B	{
28	1C	File separator	FS	CTRL-\	60	3C	<	92	5C	1	124	7C	1
29	1D	Group separator	GS	CTRL-]	61	3D	-	93	5D	j	125	7D	}
30	1E	Record separator	RS	CTRL-^	62	3E	>	94	5E	^	126	7E	~
31	1F	Unit separator	US	CTRL-	63	3F	?	95	5F		127	7F	DEL

Figure 3: Table ASCII





Généralement, l'encodage est présenté en hexadécimal. Par exemple le code 0x41, correspond à la lettre « A ».

Question 2) Déterminer la signification du code suivant codé en ASCII

4C 69 66 65 20 69 6D 69 74 61 74 65 73 20 61 72 74

Life imitates art

Question 3) **Déterminer** la taille (en octets) de la phrase «Roger is an alien !». La phrase comporte 19 caractères. Chaque caractère est codé sur 1 octet. La phrase a donc une taille de 19 octets

Question 4) **Indiquer** si l'encodage ASCII permet de coder la phrase «Le cinquième élément est un bon film»

LA phrase en peut pas être encodée en ASCII puisqu'elle comporte des caractères non reconnus (é)

Le codage ASCII n'est maintenant que très rarement utilisé et ce pour des programmes spécifiques et relativement anciens.

Le codage ISO-8859-1

Le codage ASCII précédemment présenté utilise 7 bits pour l'encodage des caractères. Avec l'avancée techniques et des méthodes de contrôles, le 8ème bit a pu être utilisé pour coder plus de caractères.

Utiliser le 8ème bit permet de passer à 256 caractères possibles, cela permet notamment d'obtenir beaucoup des

							- 1	SO-885	9-1							
	x0	х1	x2	х3	x4	х5	х6	х7	x8	х9	хA	хB	хC	хD	хE	хF
0x	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	[FF]	CR	SO	SI
1x	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2x	SP	!	17	#	\$	%	&	1	()	*	+	,	-		1
Зх	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4x	@	Α	В	С	D	E	F	G	Н	- 1	J	K	L	М	N	0
5x	Р	Q	R	S	Т	U	V	W	X	Υ	Z	[١]	٨	_
6x	`	a	b	С	d	е	f	g	h	i	j	k	- 1	m	n	0
7x	р	q	r	s	t	u	V	w	x	у	z	{	- 1	}	~	DEL
8x	PAD	HOP	BPH	NBH	IND	NEL	SSA	ESA	HTS	HTJ	VTS	PLD	PLU	RI	SS2	SS3
9x	DCS	PU1	PU2	STS	CCH	MW	SPA	EPA	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
Ax	NBSP	i	¢	£	п	¥	1	§	-	©	a	"	7	-	®	-
Вх	0	±	2	3		μ	¶			1	0	>>	1/4	1/2	3/4	ż
Сх	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	ì	ĺ	î	Ï
Dx	Đ	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	В
Ex	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
Fx	ð	ñ	ò	ó	ô	õ	Ö	÷	Ø	ù	ú	û	ü	ý	þ	ÿ

Figure 4: Codage ISO-8859-1 (latin 1)

caractères de la langue française. L'encodage principalement utilisé est l'iso-8859-1 ou latin-1 (il manque cependant par exemple le caractère œ et €). La table d'encodage est représentée en figure 4, on constate qu'elle englobe la table ASCII.

Les tables iso-8859 sont aux nombres de 16. Une table également exploitable en français est l'iso-8859-15 (ou latin-9). Cette table contient notamment les deux caractères

Différences ISO 8859-15 - ISO 8859-1

Position	0xA4	0xA6	0xA8	0xB4	0xB8	0xBC	0xBD	0xBE
8859-1	п	- 1	-	•	3	1/4	1/2	3/4
8859-15	€	Š	š	Ž	ž	Œ	œ	Ÿ

Figure 5: Caractères différents sur la table latin-9







manquants décrits précédemment. Les caractères différents entre les deux tables sont indiqués sur la figure 5

L'encodage Windows-1252

Windows s'est basé sur latin-1 pour mettre au point son jeu de caractères et son encodage appelé Windows-1252 (qu'on appelle aussi CP1252 ou ANSI). Cet encodage a l'avantage de permettre l'utilisation des caractères œ et €.

Cependant la multiplication des encodages allant souvent de paire avec une modification des jeux de caractères pose un vrai problème dans la transmission de données d'un continent à l'autre. C'est à ce moment là qu'intervient l'Unicode.

L'Unicode

L'Unicode est un jeu de caractères universel qui contient l'ensemble des alphabets et des caractères spéciaux. Le développement de cette norme s'est fait au début des années 1990 et en est actuellement à sa septième version. Elle recense environ 110 000 caractères. On y trouve donc les alphabets latin, cyrillique, grec, arabe etc. Pour coder ce jeu de caractères, il faut alors 4 octets.

Question 5) Calculer le nombre de caractères différents que peut posséder la table de caractère Unicode

La table de caractère Unicode code les caractères sur 4 octets soit 32 bits. Elle possède donc au maximum 2³² caractères soit 4 294 967 296 caractères différents

UTF-8

Jusqu'à présent, on pouvait confondre jeu de caractères et encodage, mais c'est à ce moment là qu'intervient l'UTF-8. En effet, coder une lettre sur 4 octets n'est vraiment pas utile et alourdit inutilement le fichier. Prenons en exemple la lettre « A », il s'agit du 66ème, son code en hexadécimal est 0x41. Son codage en Unicode est alors 00000000 00000000 00000000 01000001. Cela fait quand même beaucoup de 0 pour rien du tout. L'UTF-8 s'appuie sur le jeu de caractères de l'Unicode mais l'encodage s'effectue à longueur variable.

L'intérêt de l'UTF-8 est d'utiliser un nombre d'octets minimal d'encoder les caractères. Au final, on peut déterminer le nombre d'octets nécessaires en UTF-8 pour faire l'encodage du caractère suivant ce tableau:

Definition du nombre à octet	3 duliaca
Représentation binaire UTF-8	Signification
0xxxxxx	1 octet codant 1 à 7 bits
110xxxxx 10xxxxxx	2 octets codant 8 à 11 bits
1110xxxx 10xxxxxx 10xxxxxx	3 octets codant 12 à 16 bits
11110xxx 10xxxxxx 10xxxxxx 10xxxxxx	4 octets codant 17 à 21 bits

Définition du nombre d'actets utilisés

Comparaison de l'Unicode et de l'UTF-8

Le développement de l'encodage UTF-8 est incontestable au cours du temps. En effet, son universalité rend son utilisation très pratique et accessible à tous. En 2010, 50% environ des sites Web utilisaient un encodage UTF-8, de nos jours c'est plus de 90%, comme le montre ces deux graphiques².

2 Wikipédia: https://fr.wikipedia.org/wiki/UTF-8



3- Encodage d'un texte





