

1 TD : Codage des caractères

Tout pour comprendre et Éviter les erreurs d'encodage

2 Au commencement était l'ASCII

pour *American Standard Code for Information Interchange*, créé en 1960 aux États-Unis.

Dans les années 50, il existait un nombre important d'encodages de caractères dans un ordinateur, les imprimantes ou les lecteurs de carte. Tous ces encodages étaient incompatibles les uns avec les autres, ce qui rendait les échanges particulièrement difficiles car il fallait utiliser des programmes pour convertir les caractères d'un encodage à un autre. Pour tenter de mettre un peu d'ordre dans tout ça, en 1960, l'American Standards Association (ASA, aujourd'hui ANSI) décide de mettre un peu d'ordre dans ce bazar en créant la norme ASCII (American Standard Code for Information Interchange).

À chaque caractère est associé un nombre binaire sur 8 bits (1 octet).

Char	Dec	Oct	Hex	Char	Dec	Oct	Hex	Char	Dec	Oct	Hex
(sp)	32	0040	0x20	@	64	0100	0x40	`	96	0140	0x60
!	33	0041	0x21	A	65	0101	0x41	a	97	0141	0x61
"	34	0042	0x22	B	66	0102	0x42	b	98	0142	0x62
#	35	0043	0x23	C	67	0103	0x43	c	99	0143	0x63
\$	36	0044	0x24	D	68	0104	0x44	d	100	0144	0x64
%	37	0045	0x25	E	69	0105	0x45	e	101	0145	0x65
&	38	0046	0x26	F	70	0106	0x46	f	102	0146	0x66
'	39	0047	0x27	G	71	0107	0x47	g	103	0147	0x67
(40	0050	0x28	H	72	0110	0x48	h	104	0150	0x68
)	41	0051	0x29	I	73	0111	0x49	i	105	0151	0x69
*	42	0052	0x2a	J	74	0112	0x4a	j	106	0152	0x6a
+	43	0053	0x2b	K	75	0113	0x4b	k	107	0153	0x6b
,	44	0054	0x2c	L	76	0114	0x4c	l	108	0154	0x6c
-	45	0055	0x2d	M	77	0115	0x4d	m	109	0155	0x6d
.	46	0056	0x2e	N	78	0116	0x4e	n	110	0156	0x6e
/	47	0057	0x2f	O	79	0117	0x4f	o	111	0157	0x6f
0	48	0060	0x30	P	80	0120	0x50	p	112	0160	0x70
1	49	0061	0x31	Q	81	0121	0x51	q	113	0161	0x71
2	50	0062	0x32	R	82	0122	0x52	r	114	0162	0x72
3	51	0063	0x33	S	83	0123	0x53	s	115	0163	0x73
4	52	0064	0x34	T	84	0124	0x54	t	116	0164	0x74
5	53	0065	0x35	U	85	0125	0x55	u	117	0165	0x75
6	54	0066	0x36	V	86	0126	0x56	v	118	0166	0x76
7	55	0067	0x37	W	87	0127	0x57	w	119	0167	0x77
8	56	0070	0x38	X	88	0130	0x58	x	120	0170	0x78
9	57	0071	0x39	Y	89	0131	0x59	y	121	0171	0x79
:	58	0072	0x3a	Z	90	0132	0x5a	z	122	0172	0x7a
;	59	0073	0x3b	[91	0133	0x5b	{	123	0173	0x7b
<	60	0074	0x3c	\	92	0134	0x5c		124	0174	0x7c
=	61	0075	0x3d]	93	0135	0x5d	}	125	0175	0x7d
>	62	0076	0x3e	^	94	0136	0x5e	~	126	0176	0x7e
?	63	0077	0x3f	_	95	0137	0x5f				

En faite, seuls 7 bits sont utilisés pour coder un caractère, le 8e bit n'est pas utilisé pour le codage des caractères. Avec 7 bits il est possible de coder jusqu'à 128 caractères ce qui est largement suffisant pour un texte écrit en langue anglaise (pas d'accents et autres lettres particulières).

- Les 32 premiers codes, de 0 à 31, ne sont pas des caractères imprimables mais des caractères "de contrôle". Par exemple le code 13 représente un retour à la ligne, et le code 7 fait produire un bip à certains ordinateurs, ce qui s'avérait utile sur les premiers IBM PC pour signaler une erreur, par exemple.
- À partir du code 32, suivent des signes de ponctuation et quelques symboles mathématiques comme ! ou + ou /, puis les chiffres arabes de 0 à 9, ainsi que les 26 lettres de l'alphabet latin, en capitales puis en minuscules.

Exercice Décoder l'expression suivante, écrite en ASCII :

```
1101100 1100101 1110011 1000000 1001110 1010011 1001001 1000000 1100011 100111 1100101 1110011 1110100
1000000 1101100 1100101 1110011 1000000 1101101 1100101 1101001 1101100 1101100 1100101 1110101 1110010
1110011
```

Vérification avec un script Python :

Aide :

- la fonction `split(" ")` permet de décomposer une chaîne de caractères en une liste, en se servant de l'espace " " comme caractère séparateur.
- `int("1101100",2)` permet de récupérer facilement la valeur en base 10 du nombre binaire `1101100`.
- la fonction `chr` renvoie le caractère correspondant à un entier.

```
>>> chr(78)
N
```

- La fonction `ord` de Python renvoie le code ASCII correspondant à un caractère. L'entier renvoyé est en base 10 (que l'on peut convertir en hexadécimal avec la fonction `hex`).

```
>>> ord('a')
97
>>> hex(ord('a'))
'0x61'
```

```
msg = "1101100 1100101 1110011
1000000 1001110 1010011
1001001 1000000 1100011 100111
1100101 1110011 1110100 1000000
1101100 1100101 1110011 1000000
1101101 1100101 1101001 1101100
1101100 1100101 1110101 1110010 1110011"
msg = msg.split(' ')
s = ""
for k in msg :
    s += chr(int(k,2))
print(s)
```

3 Et le reste du monde ?

Lorsque d'autres personnes que des américains ou des anglais ont voulu s'échanger des données faisant intervenir du texte, certains caractères (é, è, à, ñ, Ø, Ö, ß, 漢...) étaient manquants. Les 127 caractères de l'ASCII étaient largement insuffisants.

Il a donc été décidé de passer à... 256 caractères ! Il suffisait pour cela de coder les caractères non plus sur 7 bits mais sur 8 bits.

Ainsi naquit, a norme **ISO-8859-1**, une extension de l'ASCII qui utilise les huit bits de chaque octet pour représenter les caractères.

Cette norme va être principalement utilisée dans les pays européens puisqu'elle permet d'encoder les caractères utilisés dans les principales langues européennes (la norme ISO-8859-1 est aussi appelée "latin1" car elle permet d'encoder les caractères de l'alphabet dit "latin").

Pour ajouter à la complexité, la norme ISO-8859 définit pas moins de 15 versions différentes, pour satisfaire à tous les besoins mondiaux.

Donc après de nombreuses modifications successives (la dernière en date rajoutant par exemple le symbole €), la célèbre table **ISO 8859-15**, dite aussi **Latin-9** :

ISO/CEI 8859-15																
	x0	x1	x2	x3	x4	x5	x6	x7	x8	x9	xA	xB	xC	xD	xE	xF
0x	non utilisé															
1x																
2x		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3x	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4x	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5x	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6x	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7x	p	q	r	s	t	u	v	w	x	y	z	{		}	~	
8x	non utilisé															
9x																
Ax		ı	ċ	£	€	¥	Š	š	Š	©	ª	«	¬		®	ˆ
Bx	º	±	²	³	Ž	μ	¶	·	ž	¹	º	»	Œ	œ	ÿ	ı
Cx	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
Dx	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
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

3.1 Utilisation :

Les codes sont donnés en hexadécimal :

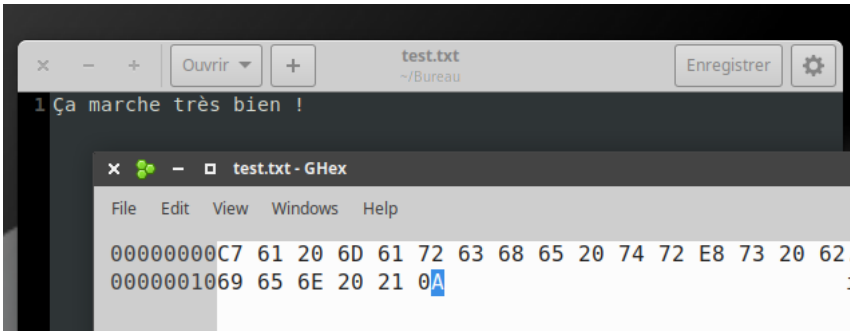
- le caractère € correspond au code hexadécimal A4, donc au nombre décimal 164.
- le caractère A correspond au code hexadécimal 41, donc au nombre décimal 65.

65... comme en ASCII ! Oui, la (seule) bonne idée aura été d'inclure les caractères ASCII avec leur même code, ce qui rendait cette nouvelle norme rétro-compatible.

Exemple :

A l'aide de notepad écrire un texte (Ça marche très bien !). Enregistrer le avec l'encodage Latin-9.

Ce fichier est ensuite ouvert avec un éditeur hexadécimal, qui permet d'observer la valeur des octets qui composent le fichier. (Comme le fichier est un .txt, le fichier ne contient que les données et rien d'autre.)



Parfait, mais comment font les Grecs pour écrire leur alphabet ?
Pas de problème, il leur suffit d'utiliser... une autre table, appelée ISO-8859-7 :

ISO/CEI 8859-7:2003																
	x0	x1	x2	x3	x4	x5	x6	x7	x8	x9	xA	xB	xC	xD	xE	xF
0x	Inutilisé															
1x																
2x	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3x	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4x	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5x	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6x	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7x	p	q	r	s	t	u	v	w	x	y	z	{		}	~	
8x	Inutilisé															
9x																
Ax	NBSP	'	'	£	€	Ɔ	ı	§	"	©	,	«	¬	SHY		—
Bx	°	±	²	³	´	µ	À	Á	Â	Ã	Ä	Å	»	Ö	½	¾
Cx	ı	Α	Β	Γ	Δ	Ε	Ζ	Η	Θ	Ι	Κ	Λ	Μ	Ν	Ξ	Ο
Dx	Π	Ρ		Σ	Τ	Υ	Φ	Χ	Ψ	Ω	Ϊ	Ϋ	ά	έ	ή	ί
Ex	ύ	α	β	γ	δ	ε	ζ	η	θ	ι	κ	λ	μ	ν	ξ	ο
Fx	π	ρ	ς	σ	τ	υ	φ	χ	ψ	ω	ϊ	ϋ	ό	ύ	ώ	

On retrouve les caractères universels hérités de l'ASCII, puis des caractères spécifiques à la langue grecque... oui mais les Thaïlandais alors ?

Pas de problème, ils ont la ISO-8859-11 :

ISO/IEC 8859-11:2001																
	x0	x1	x2	x3	x4	x5	x6	x7	x8	x9	xA	xB	xC	xD	xE	xF
0x	Inutilisé															
1x																
2x	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3x	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4x	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5x	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6x	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7x	p	q	r	s	t	u	v	w	x	y	z	{		}	~	
8x	Inutilisé															
9x																
Ax	NBSP	ก	ข	ฃ	ค	ฅ	ฆ	ง	จ	ฉ	ช	ฌ	ญ	ฎ	ฏ	
Bx	ฐ	ฑ	ฒ	ณ	ด	ต	ถ	ท	ธ	น	บ	ป	ผ	ฝ	พ	ฟ
Cx	ภ	ม	ย	ร	ฤ	ล	ฦ	ว	ศ	ห	ส	ฬ	ฬ	อ	ฮ	า
Dx	ะ	ั	า	ำ	เ	อ	ะ	็	็	็	็	็	็	็	็	B
Ex	เ	แ	โ	ใ	ใ	ใ	ใ	ใ	ใ	ใ	ใ	ใ	ใ	ใ	ใ	๓
Fx	๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๐	๑	๒	๓	๔	

Évidemment, quand tous ces gens veulent discuter entre eux, les problèmes d'encodage surviennent immédiatement : certains caractères sont remplacés par d'autres.

3.1.1 Que fait un logiciel à l'ouverture d'un fichier texte ?

Il essaie de deviner l'encodage utilisé... Parfois cela marche, parfois non.

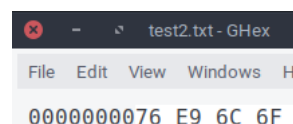


Normalement, pour un navigateur, une page web correctement codée doit contenir dans une balise `meta` le `charset` utilisé.

```
1 <!DOCTYPE html>
2 <html dir="ltr" lang="en">
3 <meta charset="utf-8" />
4 <meta name="viewport" content="width=device-width,initial-scale=1,maximu
5 <link rel="preconnect" href="//abs.twimg.com" />
6 <link rel="preconnect" href="//api.twitter.com" />
7 <link rel="preconnect" href="//pbs.twimg.com" />
8 <link rel="preconnect" href="//t.co" />
9 <link rel="preconnect" href="//video.twimg.com" />
10 <link rel="dns-prefetch" href="//abs.twimg.com" />
11 <link rel="dns-prefetch" href="//api.twitter.com" />
```

Mais parfois, il n'y a pas d'autre choix pour le logiciel d'essayer de deviner l'encodage qui semble être utilisé.

Exercice :



Le mot représenté par les octets ci-dessous est-il encodé en ASCII ou en Latin-9 ?

4 Enfin une normalisation : l'arrivée de l'UTF

En 1996, le [Consortium Unicode](#) décide de normaliser tout cela et de créer un système unique qui contiendra l'intégralité des caractères dont les êtres humains ont besoin pour communiquer entre eux.

Unicode provides a unique number for every character,
no matter what the platform,
no matter what the program,
no matter what the language.

Ils créent l'Universal character set Transformation Format : l'UTF.
Ou plutôt ils en créent... plusieurs :

- l'UTF-8 : les caractères sont codés sur 1, 2, 3 ou 4 octets.
- l'UTF-16 : les caractères sont codés sur 2 ou 4 octets.
- l'UTF-32 : les caractères sont codés sur 4 octets.

Pourquoi est-ce encore si compliqué ? En UTF-32, 32 bits sont disponibles, soit $2^{32} = 4294967296$ caractères différents encodables.

C'est largement suffisant, mais c'est surtout très très lourd !

D'autres encodages plus légers, mais plus complexes, sont donc proposés :

Arrêtons-nous sur l'UTF-8 :

Définition du nombre d'octets utilisés dans le codage (uniquement les séquences valides)			
Caractères codés	Représentation binaire UTF-8	Premier octet valide (hexadécimal)	Signi
U+0000 à U+007F	0xxxxxxx	00 à 7F	1 octet, cc
U+0080 à U+07FF	110xxxxx 10xxxxxx	C2 à DF	2 octets, c
U+0800 à U+FFFF	11100000 101xxxxx 10xxxxxx	E0 (le 2 ^e octet est restreint de A0 à BF)	3 octets, c
U+1000 à U+1FFF	11100001 10xxxxxx 10xxxxxx	E1	
U+2000 à U+3FFF	1110001x 10xxxxxx 10xxxxxx	E2 à E3	
U+4000 à U+7FFF	111001xx 10xxxxxx 10xxxxxx	E4 à E7	
U+8000 à U+BFFF	111010xx 10xxxxxx 10xxxxxx	E8 à EB	
U+C000 à U+CFFF	11101100 10xxxxxx 10xxxxxx	EC	
U+D000 à U+D7FF	11101101 100xxxxx 10xxxxxx	ED (le 2 ^e octet est restreint de 80 à 9F)	
U+E000 à U+FFFF	1110111x 10xxxxxx 10xxxxxx	EE à EF	
U+10000 à U+1FFFF	11110000 1001xxxx 10xxxxxx 10xxxxxx	F0 (le 2 ^e octet est restreint de 90 à BF)	4 octets, c
U+20000 à U+3FFFF	11110000 101xxxxx 10xxxxxx 10xxxxxx		
U+40000 à U+7FFFF	11110001 10xxxxxx 10xxxxxx 10xxxxxx	F1	
U+80000 à U+FFFFF	1111001x 10xxxxxx 10xxxxxx 10xxxxxx	F2 à F3	
U+100000 à U+10FFFF	11110100 1000xxxx 10xxxxxx 10xxxxxx	F4 (le 2 ^e octet est restreint de 80 à 8F)	

Le principe fondateur de l'UTF-8 est qu'il est **adaptatif** : les caractères les plus fréquents sont codés sur un octet, qui est la taille minimale (et qui donne le 8 de "UTF-8"). Les autres caractères peuvent être codés sur 2, 3 ou 4 octets au maximum.

UTF-8 n'utilise que l'espace dont il a besoin pour un caractère. Cela signifie donc que certains caractères n'utilisent qu'un seul octet, et d'autres deux, trois et même quatre.

UTF-8 utilise les mêmes codes qu'**ASCII** pour les 127 premiers caractères, et se sert d'octets additionnels pour représenter des caractères spéciaux comme 'é'.

Par exemple, le caractère Z serait représenté de la même façon qu'en ASCII : 01011010

Toutefois, le caractère ç devra être représenté en deux octets, car il ne fait pas partie des 127 caractères originaux.

L'encodage d'un caractère multi-octet se fait comme suit :

Les premiers bits identifient le nombre d'octets à utiliser.

- 0xxxxxx : signifie que le caractère fait 1 octet de long
- 110xxxxx signifie que le caractère fait 3 octets de long,
- 1110xxxx signifie 3 octets,
- 11110xxx 4 octets,

On note U+XXXX un caractère encodé en UTF8. Les bits restants sont utilisés pour représenter le numéro du caractère.

Latin étendu B

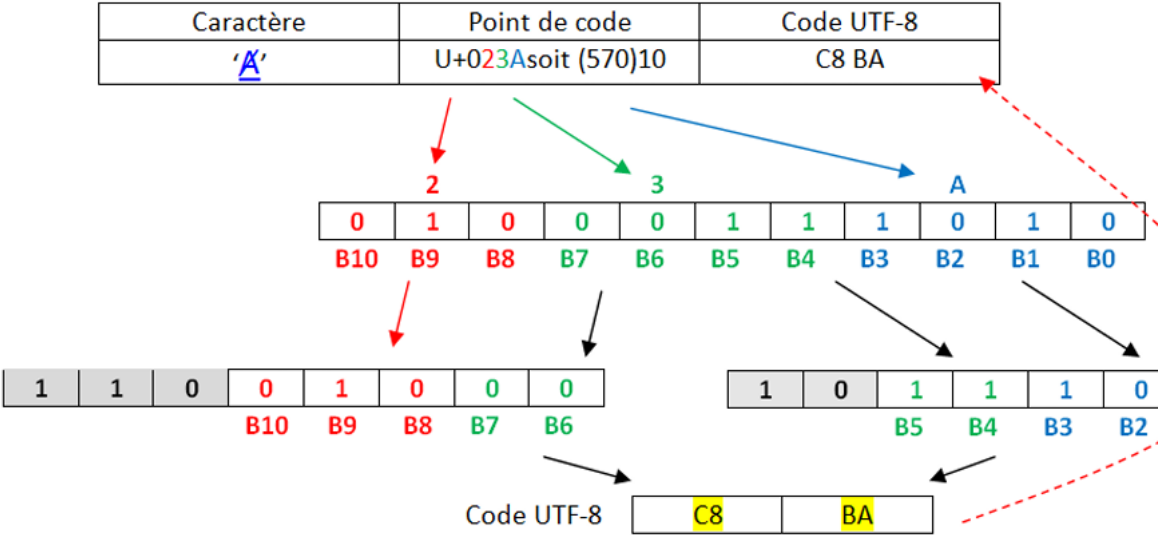
HEX		0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
	DEC	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
180	384	à	á	â	ã	ä	å	ç	ĉ	ċ	ċ	ċ	ċ	ċ	ċ	ċ	ċ
190	400	ē	ĕ	ĝ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ	ĥ
1A0	416	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ	Œ
1B0	432	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ	ŭ
1C0	448	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
1D0	464	İ	İ	İ	İ	İ	İ	İ	İ	İ	İ	İ	İ	İ	İ	İ	İ
1E0	480	Ā	ā	Ā	ā	Ā	ā	Ā	ā	Ā	ā	Ā	ā	Ā	ā	Ā	ā
1F0	496	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ	Ĵ
200	512	Ä	ä	Ä	ä	Ä	ä	Ä	ä	Ä	ä	Ä	ä	Ä	ä	Ä	ä
210	528	Ř	ř	Ř	ř	Ř	ř	Ř	ř	Ř	ř	Ř	ř	Ř	ř	Ř	ř
220	544	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ	ŋ
230	560	Ō	ō	Ō	ō	Ō	ō	Ō	ō	Ō	ō	Ō	ō	Ō	ō	Ō	ō
240	576	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ	ƶ

Exemple :

A partir de l'extrait de la table Unicode(version hexa), rappelez le Point de Code du caractère Å (valeur numérique).

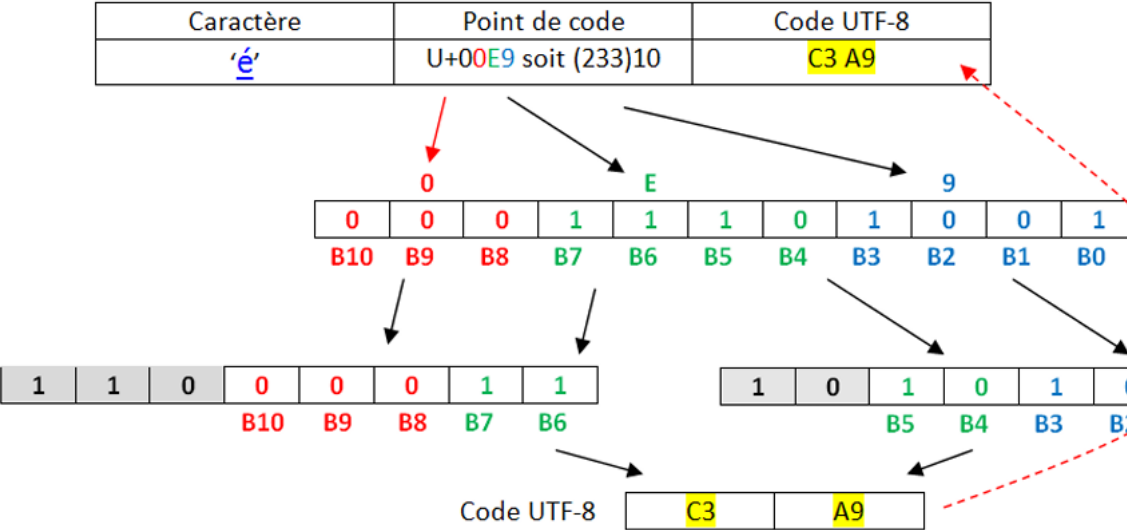
Le caractère appartient à l'intervalle U+0080 à U+07FF, donc les 11 bits seront répartis sur deux octets en code UTF-8.

- Convertir le code en binaire sur 11 bits
- 11 bits sur deux octets selon la disposition spécifiée dans la norme



HEX		0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
	DEC	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
80	128	PAD	HOP	BPH	NBH	IND	NEL	SSA	ESA	HTS	HTJ	VTS	PLD	PLU	RI	SS2	SS3
90	144	DCS	PU1	PU2	STS	CCH	MW	SPA	EPA	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
A0	160	NBSP	¡	¢	£	¤	¥	¦	§	¨	©	ª	«	¬	SHY	®	¯
B0	176	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
C0	192	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
D0	208	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
E0	224	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
F0	240	ð	ñ	ò	ó	ô	õ	÷	ø	ù	ú	û	ü	ý	þ	ÿ	

Exercice Faire de même avec la lettre é



Exercice Quel est le code binaire de "défi" codé avec UTF-8 ?

Exercice Quels mots se cachent sous les codes UTF-8 suivants ?

1. `01101000 01100101 01101100
01101100 01101111`

2. `01101001 01101110
01100110 01101111
01110010 01101101
01100001 01110100
01101001 01110001
01110101 01100101`

3. `01100010 01101001
01101110 01100001
01101001 01110010 01100101`

4.1 Utilisation grandissante de l'encodage UTF-8

La majorité des sites internet utilisent maintenant l'UTF-8, tout comme les systèmes d'exploitation récents.

Growth of UTF-8 on the Web

