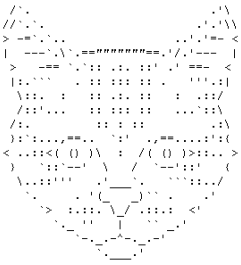


5

Représentation des textes

Extrait du programme

Thème : TYPES ET VALEURS DE BASE



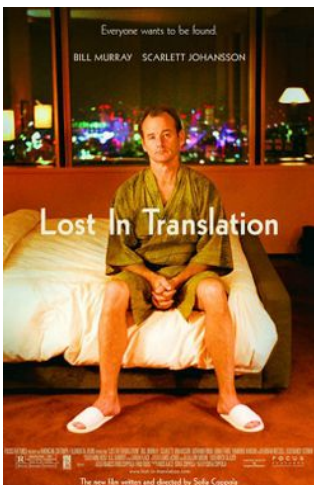
Contenus : Représentation d'un texte en machine.

Exemples des encodages ASCII, ISO-8859-1, Unicode

Capacités attendues : Identifier l'intérêt des différents systèmes d'encodage.

Convertir un fichier texte dans différents formats d'encodage.

Commentaires : Aucune connaissance préalable des normes d'encodage n'est exigible.



Mais pourquoi ai-je reçu le message suivant ? :

Ceci est un texte accentué enregistré en utf8 avec des Å des Å des Å\$ des Å et même des Å¹

I. Codage ASCII

a) Code ASCII

Dans les années 1950, il existait un nombre important d'encodages de caractères dans les ordinateurs, 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 dans un autre. Au début des années 1960, l'*ANSI* (American National Standards Institute) propose une norme de codage de caractères appelée **ASCII** (American Standard Code for Information Interchange). Cette norme définit un jeu de 128 caractères. En effet, à l'époque les ordinateurs fonctionnaient en 8 bits, et, 1 bit de parité étant conservé pour la détection des erreurs de transmission, soit $2^7 = 128$ caractères, ce qui était très largement suffisant pour coder les lettres majuscules, minuscules, chiffres et ponctuations.

Le document ci-dessous date de 1972 (*source Wikipédia*)

USASCII code chart

<div> <div> <div>b₇</div> <div>b₆</div> <div>b₅</div> </div> <div> <div>b₄</div> <div>b₃</div> <div>b₂</div> <div>b₁</div> </div> <div> <div>Column</div> <div>Row</div> </div> </div>	0 0		0 0		0 1		0 1		1 0		1 0		1 1		1 1	
	0		1		2		3		4		5		6		7	
0 0 0 0	0		1		2		3		4		5		6		7	
0 0 0 0	NUL		DLE		SP		0		@		P		\		p	
0 0 0 1	SOH		DC1		!		1		A		Q		a		q	
0 0 1 0	STX		DC2		"		2		B		R		b		r	
0 0 1 1	ETX		DC3		#		3		C		S		c		s	
0 1 0 0	EOT		DC4		\$		4		D		T		d		t	
0 1 0 1	ENQ		NAK		%		5		E		U		e		u	
0 1 1 0	ACK		SYN		&		6		F		V		f		v	
0 1 1 1	BEL		ETB		'		7		G		W		g		w	
1 0 0 0	BS		CAN		(8		H		X		h		x	
1 0 0 1	HT		EM)		9		I		Y		i		y	
1 0 1 0	LF		SUB		*		:		J		Z		j		z	
1 0 1 1	VT		ESC		+		;		K		[k		{	
1 1 0 0	FF		FS		,		<		L		\		l			
1 1 0 1	CR		GS		-		=		M]		m		}	
1 1 1 0	SO		RS		.		>		N		^		n		~	
1 1 1 1	SI		US		/		?		O		—		o		DEL	

Pour lire le codage de chaque caractère, on commence par récupérer le code binaire de la colonne correspondante, auquel on concatène¹ le code binaire de la ligne correspondante.

Ainsi, pour le caractère A, on est sur la colonne 100 et sur la ligne 0001. Le code binaire ASCII du A est donc 100 0001. Soit en décimale : $1 \times 2^6 + 0 \times 2^5 + \dots + 0 \times 2^1 + 1 \times 2^0 = 65$.

On peut faire de même pour le codage en hexadécimal du caractère A. On est sur la colonne 4 et sur la ligne 1. Le codage hexadécimal du caractère A est donc 41, soit en décimale : $4 \times 16^1 + 1 \times 16^0 = 65$.

Exercice 1

Donner le codage binaire, hexadécimal puis décimal des caractères suivants :

1. a :

- binaire
- hexadécimal
- décimal

1. Si, si, ce mot existe.

2. 9 :

- binaire
- hexadécimal
- décimal

Dans la suite, nous ne donnerons que les codages en hexadécimal, ce qui est beaucoup plus court et pratique.

On remarque que la table **ASCII** contient plusieurs catégories de caractères :

- les chiffres de 0 à 9 (entre 30 et 39)
- les lettres de l'alphabet latin en majuscules (entre 41 et 5A)
- les lettres de l'alphabet latin en minuscules (entre et)
- des signes de ponctuations (comme la virgule qui vaut 2C, le crochet [qui vaut)
- des signes opérateurs arithmétiques (comme le + qui vaut 2B, le = qui vaut)
- des caractères spéciaux (entre 00 et 20)

b) Caractères spéciaux

Voici quelques détails sur certains caractères spéciaux :

Caractère	Code hexa	Signification
HT	09	Tabulation horizontale
LF	0A	Nouvelle ligne
CR	0D	Retour chariot
FF	0C	Nouvelle page
SP	20	Espace
BS	...	Suppression
DEL	7F	Effacement

Exercice 2

Décoder le message suivant écrit en code ASCII suivants donnés en hexadécimal :

4F 75 69 20 43 48 45 46 21 (0D) 0A 42 69 65 6E 20 43 48 45 46 21

.....

.....

c) Python

En Python, on peut obtenir directement le code **ASCII** en décimal d'un caractère avec la fonction `ord()` et le caractère correspondant à un code avec la fonction `chr()` :

```

1 >>> ord('A')
2 65
3 >>> bin(ord('A')) #code en binaire
4 '0b1000001'
5 >>> hex(ord('A')) #code en hexadécimal
6 '0x41'
7 >>> chr(65)
8 'A'
9 >>> chr(0b1000001) #en binaire
10 'A'
11 >>> chr(0x41) #en hexadécimal
12 'A'
```

On peut aussi écrire directement en hexadécimal en rajoutant `\x` devant le code en hexadécimal. On appelle cette technique la technique du caractère échappé car le caractère `\` s'appelle aussi le caractère d'échappement.

```
1 >>> print( '\x4F\x75\x69\x20\x43\x48\x45\x46\x21\x0A\x42\x69\x65\x6E\x20\x43\x48\x45\x46\x21' )
2 Oui CHEF!
3 Bien CHEF!
```

Il existe quelques raccourcis pour les caractères spéciaux. Le tableau ci-dessous en propose quelques-uns :

Raccourcis	Caractère	Signification
<code>\t</code>	HT	Tabulation horizontale
<code>\n</code>	LF	Nouvelle ligne
<code>\r</code>	CR	Retour chariot
<code>\f</code>	FF	Nouvelle page
<code>\b</code>	BS	Suppression

Ainsi, on a :

```
1 >>> print( 'Biem\n!\nnpetit\tscarabée' )
2 Bien!
3 petit    scarabée
```

II. Les normes ISO 8859

Les caractères imprimables de la table *ASCII* se sont vite avérés insuffisants pour transmettre des textes dans les autres langues que l'anglais. En effet, rien qu'en considérant les langues reposant sur un alphabet latin, il manque dans la table *ASCII* de nombreux caractères comme les lettres accentuées (À, Â, Ã, Ä, Å, ...), les symboles monétaires (€, ...). Pour remédier à ce problème, l'ISO (Organisation de Normalisation) a proposé la norme **ISO 8859**, une extension de l'*ASCII* qui définit 128 premiers caractères de la norme *ASCII*; les 128 suivants sont ceux spécifiques à la table. Les caractères identiques

Code ISO	Nom	Zone
8859-1	latin-1	Europe occidentale
8859-2	latin-2	Europe centrale ou de l'est
8859-3	latin-3	Europe du sud
8859-4	latin-4	Europe du nord
8859-5		Cyrillique
8859-6		Arabe
8859-7		Grec
8859-8		Hébreu
8859-9	latin-5	Turc, Kurde
8859-10	latin-6	Révision du latin-4
8859-11		Thaï
8859-12		Devanagari ³ (projet abandonné)
8859-13	latin-7	Balte
8859-14	latin-8	Celtique
8859-15	latin-9	Révision du latin-1 (avec €)
8859-16	latin-10	Europe du sud-est

2. Il y en a 16 en tout, dont 10 uniquement pour les langues latines.

3. Écriture utilisée pour le sanskrit, le prākṛit, le hindi, le népālī, le marathi et plusieurs autres langues indiennes.