## CODAGE INFORMATIQUE DES SYSTÈMES D'ÉCRITURE

Version 2.0.1

Gianni Vacca - giannieanna@infonie.fr

#### Introduction

Je vous propose une petite doc sans prétention car je me suis aperçu qu'un nombre préoccupant d'internautes ignoraient tout du codage informatique des systèmes d'écriture. Il s'agira uniquement d'une présentation générale basée sur mon intérêt extra-professionnel pour tout ce qui a trait à la linguistique; vous trouverez donc ici un aperçu des diverses normes, et non un exposé technique.

L'informatique ayant été conçue par des Américains, dont la langue utilise uniquement l'alphabet latin sans aucun accent et, qui plus est, par des radins (cf. le bogue de l'an 2000) qui ont voulu minimiser le nombre de bits employés, nous nous retrouvons aujourd'hui avec une foultitude de normes et de systèmes de codage divers pour les différents systèmes d'écriture utilisés au monde. En gros, dès qu'on sort du domaine anglophone, nous sommes perdus dans une véritable Babel du codage, qui empoisonne la vie de tous ceux qui transfèrent des fichiers d'un système d'exploitation à un autre, ou qui veulent utiliser des systèmes d'écriture différents au sein d'un même document.

Je vais essayer de faire une brève présentation des systèmes les plus utilisés (notamment sur Windows, ou pour le codage dans les fichiers **HTML**) en partant de mon expérience purement empirique. Notre tour d'horizon partira de l'**ascii** sur sept bits pour arriver à la norme **UCS-4**: du moins complet au plus complet.

## 1. ascii sur sept bits ou us-ascii ou ISO 646 ou ANSI X3.4

Cette norme code l'alphabet latin sur sept bits (de x20 à x7E). C'est la base à partir de laquelle toutes les autres normes ont été conçues, ce qui signifie que les caractères codés de 32 à 126, quelle que soit la norme, seront gardés tels quels en cas de changement de codage (ex. : transfert d'un fichier de Mac à PC : les caractères accentués sont affichés de façon incorrecte, les caractères non-accentués le sont de façon correcte).

20		21	22 11	23 #	24 \$	25 %	26 &	27	28 (	29	2A *	2B +	20	2D —	2E •	2F /
30		1	32 2	33	34 4	³ 5	36	37 7	* 8	" 9	3A :	зв ;	30	3D =	3E >	3F ?
40		ΑĤ	ч <u>г</u> В	43 C	TH D	ч <u>я</u> Е	<sup>46</sup> F	47 G	ч* Н	49 I	ч <sup>н</sup> J	чв К	<sup>4c</sup> L	<sup>чо</sup> М	чE N	4F ()
50 F		51 Q	<sup>52</sup> R	s S	54 T	55 U	56 V	57 W	58 X	59 Y	<sup>58</sup> Z	5B [	5C \	50 ]	5E ^	5F —
60 ,		a a	<sub>es</sub> p	63 C	đ	65 E	ee t	67 g	° h	69 İ	°° j	eB K	6C 1	W eb	ε n	6F
70  }	- 1	q	r	73 S	74 †	75 U	76 V	77 W	78 X	79 Y	78 Z	7B {	70	70	<sup>7E</sup> ∼	

Évidemment, cette norme (1968) ne satisfit pas les pays plus civilisés tels que la France, l'Allemagne ou l'Italie car elle ne permettait pas de coder les caractères accentués dont nous autres Européens sommes si friands. La suite logique de cette norme fut donc d'utiliser le 8ème bit laissé libre (au départ prévu pour coder le surlignage) pour des extensions nationales. Ce furent les diverses normes ISO 8859-*N*., où N vaut entre 1 et 15.

#### 2. les normes ISO 8859-N

Ces normes datent de la moitié des années quatre-vingts et sont encore de très loin les plus utilisées actuellement pour les langues européennes ou en général pour les langues qui s'écrivent avec un système alphabétique (cyrillique, grec, hébreu, arabe, thaï) par opposition aux langues qui s'écrivent avec un système d'idéogrammes (**CJK** pour les intimes : Chinese Japanese Korean). Elles ont l'avantage de n'utiliser que 8 bits (en ajoutant le "bloc" xA1 à xFF à l'us-ascii) pour un codage satisfaisant de tous les systèmes alphabetiques ; ce sont les normes "de base" utilisées pour le HTML (cf. plus bas).

## 2.1 - la norme ISO 8859-1 ou Latin 1 ou Europe Occidentale

Cette norme code tous les alphabets d'Europe Occidentale (à l'exception notable du gallois) en utilisant le 8ème bit pour les caractères accentués. C'est celle que nous utilisons quotidiennement sur notre PC (bureautique, e-mail, HTML, **polices True Type** cf. plus bas) et aussi sur les stations utilisant X11 et les DEC.

AO		aı İ		ф	A3	£	A4	Ħ	A5	¥	A6	1	A7	S	A8	"	A9	0	AA	a	AB	((	AC	7	AD	_	ΑE	®	AF _
BO	۰	B1 <u>+</u>	-	5 5	В3	3	вч	-	B5	μ	B6	1	В7		B8	,	В9	1	BA	0	BB	<b>&gt;&gt;</b>	BC	14	BD	X	BE	¥	BF ¿
CO	À	άÁ	1	Â	C3	Ã	СЧ	Ä	C5	Å	Ce	Æ	C7	Ç	C8	È	C9	É	CA	Ê	СВ	Ë	СС	Ì	CD	Í	CE	Î	ΓÏ
D0 .	Đ	D1 N	í	Õ	D3	Ó	D4	ô	D5	õ	De	ö	D7	×	D8	Ø	D9	Ù	DA	Ú	DB	Û	DC	Ü	DD	Ý	DE	Þ	В
E0	à	E1 É	. 1	â	E3	ã	EЧ	ä	E5	ă	E6	æ	E7	Ç	E\$	è	E9	é	EA	ê	EB	ë	EC	ì	ED	í	EE	î	EF
F0	ð	ř	í	ò	F3	ó	FЧ	ô	F5	õ	F6	ö	F7	÷	F8	Ø	F9	ũ	FA	ú	FB	û	FC	ü	FD	ý	FE	þ	ξÿ

Ceux qui suivent auront remarqué que le "bloc" x80 à x9F reste inutilisé. Sur votre PC, en allant dans le menu *Accessoires*, lancez la *Table de caractères*. Choisissez le sous-ensemble *Latin-I*; une ligne reste, en effet, vide. Choisissez maintenant le sous-ensemble *Caractères Windows*, la ligne en question se remplit de caractères étendus supplémentaires. C'est que Microsoft, qui "n'aime pas gâcher", apparemment, a utilisé ce vide entre deux normes pour coder quelques caractères supplémentaires, dont l'Euro (€) en x80. Voir plus bas le paragraphe sur les **jeux étendus Windows**.

## 2.2 - la norme ISO 8859-2 ou Latin 2 ou Europe Centrale ou Europe Orientale

Cette norme code tous les alphabets d'Europe Centrale (à l'exception de ceux qui utilisent l'alphabet cyrillique) en utilisant le 8<sup>ème</sup> bit pour les caractères accentués. C'est celle qu'utilisent normalement les PC vendus en Europe Centrale.

AO	Ą	A2 U	<sup>A3</sup> Ł	¥	L	ee Ś	<sup>A7</sup> S	A\$	<sup>A9</sup> Š	## Ş	Т	Z	AD —	ěŽ	aF .
ВО о	a a		ł		B5 ľ	se Ś	B7 U	B8	B9 Š	BA Ş	ť	ź	BD	ĔĔŽ	BF .
٣Ŕ	άÁ	٩Ê	Ä	Ä	° Ĺ	"Ć	"Ç	°Č	°É	űĘ	вË	"Ě	°ί	Î	Ď
™ Đ	Ñ	N	ő	Ö	ő	De .:	D7 ×	°Ř	D9 Ů	DA Ú	DB Ű	Ü	۳Ý	DE J	В
ř	á	â	a	ä	ES Î	Ć	е7 С	č	eُ	ę	ªë	ěě	ÉD ĺ	î	ĔĔĞ
°đ	'n	řň	F3 Ó	۴٩ Ô	FS Ő	F6 Ö	F7 ÷	۴ř	F9 Ů	FA Ú	FB Ű	FC Ü	ΕŪ	FE ţ	FF .

C'est pour cette raison que, si vous créez un site web perso qui accueille les internautes par un sympathique

Bonjour à tous, ici on fait la fête!

codé par exemple

<H1 align=center>Bonjour à tous, ici on fait la fête !</H1>
vos amis tchèques liront un mystérieux

## Bonjour r tous, ici on fait la fete!

compréhensible, certes, mais intriguant. Nous verrons plus loin comment éviter ce genre d'inconvénient.

## 2.3 - la norme ISO 8859-3 ou Latin3 ou Europe Méridionale

Cette norme code le maltais, le turc et, accessoirement, l'esperanto. Son nom est un peu tiré par les cheveux, vu que le portugais, l'espagnol, l'italien ou l'albanais sont codés par la norme ISO 8859-1, et que le grec et le bulgare sont codés par d'autres normes encore. J'imagine que c'était par un souci de symétrie avec "Europe Occidentale" et "Europe Orientale".

AO	ΉĦ	A2 U	#3 £	A4 X		ñ Ĥ	<sup>A7</sup> S	A8	<sup>A9</sup> İ	Ş	Ğ	J	AD —		<sup>AF</sup> Ż
B0 o	⁵ħ	82	83	B4 _	85 µ	۴ĥ	B7 •	B8 ,	1	BA Ş	₩ ĕ	ĵ	<sup>80</sup> ½		βF , Ż
٩	άÁ	â		C4 Ä	° Ċ	"Ĉ	°Ç	° È	°É	Ê	Ë	" Ì	ΩÍ	Î	Ϊ
	ñ	°ò	° Ó	۳Ô	Ğ	Ö	D7 ×	°°Ĝ	D9 Ù	DA Ú	DB Û	Ü	°Ŭ	°ŝ	B
à	á	≅â		ä	Ċ	Ĉ	E7 Ç	₽ê	eُ	ê	ĕë	ĩ	Í	î	EF
	ñ	۴²	٥	۴٩ Ô	<sup>FS</sup> ġ	F6 Ö	F7 ÷	F* ĝ	F9 Ù	<sup>FA</sup> Ú	FB Û	<sup>FC</sup> Ü	FD U	ÊŜ	FF .

## 2.4 - la norme ISO 8859-4 ou Latin 4 ou Europe Septentrionale

Cette norme code les langues des Pays Baltes, le lapon, et le groënlandais. Même remarque que cidessus, le nom me paraît franchement tiré par les cheveux vu que les langues scandinaves sont du ressort de la norme ISO 8859-1.

AO	Ą	ΑZ K	Ŗ	Ħ	<sup>A5</sup> Ĩ	Ļ	<sup>A7</sup> S	A\$	<sup>A9</sup> Š	**Ē	ab G	AC 不	AD —	ěŽ	AF _
ВО о	a a		Ì.	B4	1	ļ		,	Š	ē	g	BC Ž	'n	ĔĔŽ	βF
°Ā	άÁ	٩Ê	°Ã	Ä	Å	œ Æ	ïĮ	°Č	°É	űĘ	° Ë	"Ė	ΰĺ	Î	Ī
<sup>™</sup> Đ	Ņ	Ō	ΪĶ	Ô	os Õ	De .:	D7 ×	°° Ø	D9 Ų	DA Ú	DB Û	Ü	۳ĩ	Ū	В
ã	۵	â	ã	ä	ă	æ	į	Ě	é	ę	ëë	ė	ĺ	î	ĒF Ī
⁵°đ	ŗ	FΣ Ō	F3 Ķ	Ô	F5 Õ	F6 Ö	F7 ÷	F* Ø	F9 Ų	F# Ú	FB Û	FC Ü	FD ũ	FE Ū	FF .

## 2.5 - la norme ISO 8859-5 ou Cyrillique

Cette norme code les langues européennes utilisant l'alphabet cyrillique telles que le russe, le bulgare ou le serbe. C'est donc celle qu'utilisent normalement les PC vendus en ex-Union Soviétique, par exemple.

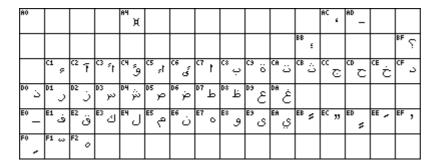
AO	Ë	ь Б	۴³ Ć	<b>6</b>	<sup>A5</sup> S	He I	<sup>A7</sup> Ï	<sup>e</sup> ³ J	<sup>нэ</sup> Љ	** Њ	<sup>AB</sup> Ћ	ac Ŕ	AD —	РΕЙ	ĥF
βO	Б	В	<sub>B3</sub> L	вч Д	85 E	ж	3	в	вэЙ	BA K	ВВ	ВС	ВВ	0	BF
°°Ρ	° C	T	зу	СЧФ	cs X	"Ц	<sup>7</sup> 4	°* Ш	°Щ	Ъ	СВ	ωР	ΰЭ	СЕЮ	۴Я
°a	ъб	B	D3	υчД	os e	ж	<sup>07</sup>	В	υэй	DA K	DB Л	М	Н	DE O	DF
р	E1 C	ΕΣ	<b>Б</b> З	Ф	E5 X	Е	E7 '4	E\$ Ш	E9 Щ	EA Ъ	EB Ы	Б	Э	EE Ю	FЯ
FO No.	F1 ë	F2 5	F³Γ	<b>F4</b> €	FS S	F6 İ	F7 Ï	f* j	F9 Љ	FA Њ	<sup>FB</sup> ħ	FC K	FD S	FΕΥ	FF Ų

Malheureusement, la Guerre Froide a encore frappé ici, et les Soviétiques avaient créé leurs propres normes **KOI** incompatibles avec la série ISO 8859-N avant que les statues de Lénine ne commencent à chuter un peu partout.

Sur internet, les normes ISO 8859-5 / KOI-8 / WinCyrillic (cf. plus bas) se répartissent de façon assez équitable les sites en langue russe, ce qui fait que même un paramétrage correct de votre fureteur n'empêchera pas que deux pages russes sur trois apparaissent comme du charabia incompréhensible.

## 2.6 - la norme ISO 8859-6 ou Arabe

Cette norme code l'arabe. Même si vous installez une police compatible avec cette norme sur votre PC, vous ne pourrez pas lire de documents écrits en arabe car il vous manquera l'outil qui intervertit la direction d'écriture.



#### 2.7 - la norme ISO 8859-7 ou Grec

Cette norme code le grec moderne (celui de Demis Roussos, donc, et pas celui d'Homère). On la nomme également ELOT-928.

À ma connaissance, contrairement à ce qui a lieu pour le cyrillique, elle n'a pas de norme concurrente.

AO	Ĥ	1,	A2 ,	#3 £			He :	<sup>87</sup> S	A\$	<sup>А3</sup> ©		AB 《	AC ¬	AD —		AF —
ВО (		±	82	83	ВЧ ,	B5 ,\	Be 'A	B7 •	** E	вэН	BA T	BB }>	Ö	BD 1/2	βEΥ	βF'Ω
C0 .	i c	A	° B	ЗΓ	<sup>сч</sup> Д	cs E	" Z	<sup>7</sup> H	° 8	° I	CA K	CB V	CC M	° N	Ξ	0
Do I	T D:	P		ρ3	Т	DS Y	Φ	D7 X	<sup>D8</sup> Ψ	Ω	DA Ï	рв Ü	ά	ΰĖ	¤ή	DF .
E0 į	j E:	α	β	E3 γ	δ.	E5 8	Ε6 ζ	ε7 η	θ	E9 L	EA K	ЕВ	EC	ED V	Ēξ	EF ()
F0 T	T F:	ρ	FΣ	F3 O	F4 T	FS U	F6 Ф	F7 X	F\$ Ψ	F9 ω	FA	FB Ü	FC ,	FD ,	FE , ώ	

#### 2.8 - la norme ISO 8859-8 ou Hébraïque

Cette norme code l'hébreu moderne et le yiddish. Même remarque que pour la norme ISO 8859-6 pour ce qui est de la direction d'écriture.

AO		φ	#3 £	A4 H	<sup>AS</sup> ¥	H6	<sup>87</sup> S	A8	<sup>нэ</sup> ©	AA ×	AB 《	AC ¬	AD —	AE ®	AF —
ВО о	B1 ±	5 88	83	84	85 µ	B6 ¶	B7 •	B8 ,	<sup>B9</sup> 1	BA ÷	BB }}	вс 74	BD 1/2	8E %	
															DF =
E0 K	E1 ]	E2 λ	E3 T	택	ES ]	J. Ee	E7 Π	E* []	E9 7	EA 7	即丁	ל "	ED 🗍		EF ]

## 2.9 - la norme ISO 8859-9 ou Latin 5

J'espère que vous suivez toujours... ici N=9 mais on parle bien de Latin 5, ce qui est logique puisque les normes ISO 8859-5 à 8 codent des alphabets non-latins.

Cette norme est la même que la norme ISO 8859-1 mais les lettres spécifiquement islandaises

des positions xF0, xFE et xFD ont été remplacées par les lettres spécifiquement turques

Cette norme a remplacé la norme ISO 8859-3 en ce qui concerne le turc. Elle a en effet l'avantage de pouvoir coder dans le même document le turc et les langues d'Europe Occidentale (qui sont tout de même plus fréquentes que le maltais ou l'esperanto!).

AO	A1 j	ф	#3 £	Ħ	¥		8		0	a	AB {{	AC ¬	AD —	AE ®	AF _
ВО о	B1 ±	82	83	вч "	85 µ	B6 ¶	B7 •	B8 _	<sup>B9</sup> 1	BA Q	)> 	вс 14	BD X	8E %4	BF ¿
٣À	άÁ	٩Â	°Ã	Ä	cs Å	Œ	"Ç	°° È	°É	° Ê	Ë	ũĩ	ΰÍ	Î	ΓÏ
°Ğ	D1 Ñ	Õ	D3 Ó	Ô۳	os Õ	Ö	D7 ×	°° Ø	D9 Ù	DA Ú	Û	Ü	Βİ	Ş	В
ñà	á	≅â	ã	ä	å	æ Ee	E7 Ç	₽ê	é	ê	≞ë	ì	í	î	EF
F°ÿ	ñ	F2 Ò	F3 Ó	۴٩ Ô	FS Õ	F6 Ö	F7 ÷	F≎ Ø	F9 Ù	<sup>FA</sup> Ú	FB Û	<sup>FC</sup> Ü	FD 1	FE Ş	ξÿ

## 2.10 - la norme ISO 8859-10 ou Latin 6

Cette norme, comme toutes les normes ISO 8859-N où N > 9, est récente puisqu'elle ne date que de 1992. Il s'agit en fait d'un réarrangement de la norme ISO 8859-4 dont certaines lettres peu utilisées ont été enlevées pour faire place aux fameuses lettres islandaises et mériter ainsi le nom de "Europe Septentrionale". Peine perdue puisque le nom reste attaché à la norme ISO 8859-4. Bad luck.

AO	Ĥ1	Ą	Ē	нз Ģ	Ϊ	<sup>A5</sup> Ĩ	¥e Ř	<sup>A7</sup> S	#* Ļ	нэ Đ	## Š	AB 下	ěŽ	AD —	<sup>AE</sup> Ū	ar D
ВО о	B1	ą	ē	₿ g	<sup>вч</sup> Ī	BS ĩ	ķ.	B7 •	* ]	"đ	BA Š	88 Ž	ĕζ	BD —	ū	BF Ŋ
۳Ā	C1	Á	âÂ	Ã	сч А	°s Å	œ Æ	ΊĮ	°Č	°É	ca Ę	œ ::	ΰĖ	ΰÍ	Ϊ	ΓËΪ
oo Đ	D1	Ņ	Ō	os Ó	Ô	os Õ	De	D7 Ũ	°° Ø	D9 U	DA Ú	DB Û	Ü	ΦÝ	Þ	ВΒ
°°Đ ē°ā	E1	Ņ	0	Ô	Ö	0	0	U	Ø	Ų	U	U	U	Υ	Þ	Β ΕF 1

#### 2.11 - la norme ISO 8859-11 ou Thaï

Cette norme code le thaï. Vu qu'il s'agit d'une norme récente, je n'ai aucune idée quant à sa diffusion réelle, et je n'ai pas de copains thaïs pour me renseigner. Je suis preneur de toute info à ce sujet.

		Ĥ1		A2		A3		A4		A5		86		87		A8		A9		AA		AB		AC		AD		ĤΕ		ĤΕ	
			ก		ข		ฃ		Я		A		ฆ		1		จ		a		ប		ซ		ณ		Ŋ		มู		ม
BO		B1		B2		В3		вч		B5		86		В7		В8		В9		BA		ВВ		BC		BD		BE		BF	
	ន		ሽ		M		ณ		Ø		Ø		ถ		n		б		и		ข		ป		W		W		W		ฟ
CO		C1		CZ		C3		СЧ		C5		C6		C7		C8		C9		CA		CB		cc		CD		CE		CF	
	ภ		IJ		ย		5		η		ล		Ŋ		3		คี		¥		ส		И		ឃ		Ð		ปี		٦
DO		D1	_	DZ	_	D3		04	$\overline{}$	D5	~	D6	~	D7	~	D8		D9		DA		Г		Г		Г		Г		DF	#
	2				.1		.1										•		~												Щ
ΕO		E1		E2	0	E3	9	EЧ	-71	E5		E6		E7	~	E8	-	E9	'n	ΕA	øv.	EB	+	EC	•	ED	•	EE	ε	EF	
	ı		Ш		٦		٦		l		٦		ໆ																		Θ
F0		F1		F2		F3		FЧ		F5		F6		F7		F8		F9		FΑ		FB		Г		Г					
	0		9		Ø		Ø		Œ		Œ		Ö		ø		ផ		ď		7		0-								

#### 2.12 - la norme ISO 8859-12 ou Indien

Cette norme n'existe pas encore. La place est réservée pour une hypothétique norme "Indienne". Vu le nombre d'alphabets différents utilisés en Inde, 8 bits me paraissent trop peu nombreux pour cette tâche herculéenne. On verra bien.

#### 2.13 - la norme ISO 8859-13 ou Latin 7

Cette norme est censée remplacer la norme ISO 8859-10 pour ce qui est des langues des Pays Baltes (en effet, il manque une lettre lettonne à cette dernière norme).

AO	Ĥ	¹ ,,	AZ	¢	<sup>A3</sup> £	A4	Ħ	A5	A6	<sup>A7</sup> S	** Ø	A9 (С)	## Ŗ	AB 《	AC ¬	AD —	e R	Æ
B0 (	, В	±	B2	2	83	вч	"	85 µ	B6 ¶	B7 •	<sup>88</sup> Ø	B9 1	ŗ	BB >>	вс 14	BD 1/2	BE ¾	æ
co F	į	Į	CZ	Ā	°ć	СЧ	Ä	Å	É	ΪĒ	°Č	°É	°É	° Ė	"Ģ	ωĶ	Ϊ	Ļ
90 5		Ñ	DZ	Ņ	٥	D4	Ō	° Õ	De C	D7 ×	" Ų	09 Ł	DA Ś	DB Ū	Ü	ΰż	ďŽ	βВ
E0 (	į E	į	E2	ā	É	E4	ä	å	é	ē7 ē	Ě	۴é	ÉÁŹ	₽ė	€ ģ	ķ	Ī	EF ]
FO S		'n	F2	ņ	F3 Ó	F4	ō	F5 Õ	F6 Ö	F7 ÷	F* Ų	F9 Ł	FA Ś	FB Ū	řü	FD ,	řŽ	FF,

#### 2.14 - la norme ISO 8859-14 ou Latin 8

Cette norme est une adaptation de la norme ISO 8859-1 avec les w accentués qui manquaient pour couvrir le gallois, et avec les consonnes surmontées d'un point pour indiquer l'amuissement en irlandais (ce qui n'était pas franchement indispensable, vu qu'on peut tout aussi bien l'indiquer avec un h après la consonne en question). Cette norme code donc toutes les langues celtiques.

AO	<b>в</b> і В	βģ	#3 £	Ċ	AS .	D We .	A7 S	<sup>e∗</sup> Ñ	аэ ()	<sup>88</sup> Ŵ	вd	ŘΥ	AD —	AE ®	¥Ϋ́
Β̈́Ė	βı	Ġ	ġ	М.	Ü.	<b>9</b> 6	Ρ̈́	<sup>₿₿</sup>	βġ	EΑ Ŵ	"Ś	۴ŷ	во	ΨÜ	BF.S
ñÃ	άÁ	٩Ê	°Ã	Ä	Å	Œ	ΰÇ	° È	°É	° Ê	° Ë	"ì	ΰÍ	Î	ΓÏ
° Ŵ	Ñ	Õ	D3 Ó	Ô	DS Õ	De ::	<sup>07</sup> †	° Ø	D9 Ù	DA Ú	Û	Ü	۳Ý	PΕŶ	В
ã	á	â	ã	ä	å	⊕ Ee	е7 С	ë è	é	ê	ë	EC Ĩ	ED 1	î	EF
Fº Ŵ	ñ	F2 Ò	F3 Ó	۰Ô	FS Õ	F6 Ö	F7 t	F≎ Ø	F9 Ù	FA Ú	FB Û	<sup>FC</sup> Ü	Бý	FE ŷ	ff :;

#### 2.15 - la norme ISO 8859-15 ou Latin 9

Cette norme est une adaptation de la norme ISO 8859-1 avec la lettre ligaturée "e dans l'o" ( $\alpha$ ) qui manquait inexpliquablement du jeu de cette dernière (c'est pour cette raison que les  $\alpha$  sont souvent ignorés par les fureteurs, mieux vaut utiliser oe dans vos pages web). Une énorme injustice vis-à-vis de notre belle langue française est enfin réparée!

Vu qu'il s'agit d'une norme récente (1998), elle comporte également le symbole de l'euro (€). Malgré ces petits avantages par rapport à la norme ISO 8859-1, elle n'a pas franchement eu de succès... Voici par exemple ce que pensait Reuters de la norme ISO 8859-15 (je cite) :

« We have just the place for ISO 8859-15 here in London. It is called the Science Museum and is full of charming historical relics. »

AO	A1 i	A	¢	A3	£	<sup>84</sup> €	A5	¥	A6	Š	A7	8	A8	š	A9 (C	) #	<sup>A</sup> a	AB	<b>«</b>	AC	7	AD	_	ΑE	®	AF _
ВО о	B1 <u>+</u>	- B	2	B3	3	ž	B5	μ	B6	9	В7		B8	ž	<sup>B9</sup> 1	В	<sup>0</sup> 0	BB	<b>&gt;&gt;</b>	ВС	Œ	BD	œ	BE	Ϋ	BF ¿
°Ã	°i A	C	Â	C3	Ã	C <sup>4</sup>	C5	Å	Ce	Æ	C7	Ç	C8	È	°É	· C	Ê	СВ	Ë	СС	Ì	CD	Í	CE	Î	ΓÏ
<sup>™</sup> Đ	D1 Ñ	Í	ò	D3	Ó	°°Ô	D5	õ	De	ö	D7	×	D8	Ø	D9 <u>`</u>	l Di	Ű	DB	Û	DC	Ü	DD	Ý	DE	Þ	В
à	E1 ĉ	E	â	E3	ã	ä	E5	å	E6	æ	E7	Ç	E8	è	E9 É	E	e ê	EB	ë	EC	ì	ED	í	EE	î	EF
řõ	F1 r	i F	ò	F3	ó	۰Ô	F5	õ	F6	ö	F7	÷	F8	Ø	F9 Ù	F	u ú	FB	û	FC	ü	FD	ý	FE	þ	ξÿ

#### 3. les autres normes sur 8 bits

La famille des normes ISO 8859-N n'est pas la seule à coder les systèmes d'écriture sur 8 bits. Puisque ces normes partaient du principe de codage du bloc xA0 à xFF, des langues comme le vietnamien n'y trouvaient pas leur compte (l'alphabet latin adapté au vietnamien utilise un très grand nombre de signes diacritiques).

Les normes sur 8 bits sont extrêmement populaires dans l'industrie de l'informatique car l'octet est une valeur sûre. Malgré les inconvénients dantesques pour tous les gens qui travaillent dans d'autres langues que l'anglais (ça fait quand même du monde ; il y a un milliard de Chinois, ne l'oublions pas), l'excuse selon laquelle l'adoption d'une norme universelle sur 16 bist (Unicode) accroîtrait la taille des fichiers bloque pour l'instant tout progrès significatif.

#### 3.1 - la norme VISCII

Cette norme est en gros la réponse vietnamienne à la norme ascii. Cependant, elle n'est reconnue par aucun fureteur, et il faut utiliser des polices particulières pour coder le vietnamien en HTML à partir des classiques ISO 8859-N.

		οž Å			°5 Ä	°° Ã									
				Ϋ́					19 Ÿ					1E Y	
20	21	22 11	23 #	24 \$	25 %	<sup>26</sup> &	27 1	28 (	29	2A *	2B +	20	2D —	2E •	2F /
30	<sup>31</sup> 1	2	"3	34 4	ຶ້ 5	36	37 7	* 8	" 9	3A :	зв ;	30	3D =	3E >	³F ?
40 @	<sup>41</sup> A	42 B	43 C	"D	45 E	<sup>46</sup> F	47 G	<sup>48</sup> H	49 I	Ч <sup>A</sup> J	<sup>чв</sup> К	чc L	<sup>чо</sup> М	чE N	<b>4F</b> O
50 P	51 Q	<sup>52</sup> R	⁵S	54 T	55 U	56 V	57 W	58 X	59 Y	<sup>58</sup> Z	5B [	5C \	50 ]	5E ^	5F —
60 t	a a	<sub>es</sub> p	e3 C	۴٩d	65 E	ee f	67 g	* h	<sup>69</sup> İ	° j	eB K	ec 1	W eb	ee n	6F 0
<sup>70</sup> р	q	72 T	73 S	74 †	75 U	76 V	77 W	78 X	79 Y	78 Z	7B {	7C	70	<sup>7E</sup> ~	
*• Ą	*1 Á	<sup>82</sup> À	°³Ă	٩٩Ã	*5 Ã	<sup>86</sup> Â	*7 Â	** Ĕ	ş, Ē	** Ê	** Ê	*° Ê	*° Ě	*Ē Ê	*F Õ
°° Ö	91 Ô	92 Õ	³³ Ô	94 Ü	95 Ű	96 Ù	<sup>97</sup> Ở	às Ì	°° Ò	эн О	³B Ì	Ü	90 Ũ	ΔĒ	9F Ŷ
° Õ	<sup>#1</sup> á	<sup>#2</sup> à	<sup>A3</sup> ă	ñ٩	as a	AF Â	<sup>A7</sup> â	* ẽ	₽3 E	" ẽ	<sup>AB</sup> ĕ	ê	<sup>AD</sup> ẽ	ê	AF Õ
ő	î	B2 Õ	B3 Õ	вч Oʻ	BS Ô	), Be	B7 Ů	<sup>B®</sup> į	μŲ	BA Ú	Ú	Ú	0,	Ő	BF U
٣Ã	άÁ	۵Â	αÃ	۲۹À	Ä	۴à	ã	° È	°É	ca Ê	вÈ	ũÌ	ωĺ	Έĩ	φ
oo Đ	Ú	Õ	° Ó	Ô	os a	ъ ў	D7 Ù	°ù	D9 Ù	DA Ú	βÿ	oc Y	٥٩	Õ	DF Ư
ĕ° à	á	≅â	ã	۴,	ă	E6 ữ	ε <sub>7</sub> ã	₽è	é	ê	₿.	EC Ì	ED ĺ	ĩ	EF 1
<sup>F⁰</sup> đ	F1 Ľ	F2 Ò	F3 Ó	۰Ô	FS Õ	F6 Ò	F7 Q	F∜ U	F9 Ù	<sup>FA</sup> Ú	FB Ũ	FC Ù	۴Ý	FE Ç	FFΩ

Il existe également des systèmes plus compliqués pour le vietnamien sur plus de 8 bits.

#### 3.2 - la norme ISO 6429

Ceux qui ont lu attentivement le paragraphe 2 auront remarqué un "trou" entre l'us-ascii et les diverses normes ISO 8859-N.

Ce trou (x80 à x9F) a été rempli par la norme ISO 6249 qui l'a utilisé pour coder des fonctions comme Carriage Return, Form Feed, Backspace, etc. J'imagine que tout cela avait pour but, à l'époque (1991),

de contrôler le comportement des imprimantes matricielles, des telex, et de toute cette sorte d'engins très bruyants.

Il paraît que cette norme est encore utilisée.

## 3.3 - les jeux étendus Windows

Le fameux "trou" x80 à x9F a également été exploité par Microsoft pour étendre le jeu de caractères utilisé par Windows. Ce jeu étendu apparaît, par exemple, au niveau des polices True Type. C'est ainsi qu'on a pu écrire  $\alpha$  en français dans un document Word ou Excel alors que cette lettre ligaturée avait été "oubliée" par la norme ISO 8859-1.

Selon les pays de localisation, les caractères x80 à x9F ne sont pas les mêmes.

#### 3.3.1 - le jeu de caractères CP1252 ou WinLatin 1

Il s'agit de la norme ISO 8859-1 avec les caractères étendus x80 à x9F.

80	€			82	,	83	f	84	,,	85		86	#	87	#	88	^	89	%	8A	Š	88	<	*C	Œ			\$E	ž	
		91	ť	92	,	93	cc	94	"	95	•	96	-	97	_	98	~	99	тн	9A	š	9B	>	9C	œ			9E	ž	9F Y
AO		A1	i	A2	ф	A3	£	A4	Ħ	A5	¥	A6	1	A7	S	A8	"	A9	0	AA	a	AB	((	AC	¬	AD	_	ΑE	®	AF _
BO	۰	B1	±	BZ	2	B3	3	В4	-	B5	μ	Be	1	В7		B8	,	B9	1	BA	0	BB	<b>&gt;&gt;</b>	BC	14	BD	X	BE	¥	BF ¿
CO	À	C1	Á	CZ	Â	C3	Ã	СЧ	Ä	C5	Å	Ce	Æ	C7	Ç	C8	È	C9	É	CA	Ê	СВ	Ë	СС	Ì	CD	Í	CE	Î	ΓËΪ
DO	Đ	D1	Ñ	DZ	ò	D3	Ó	D4	ô	D5	õ	De	ö	D7	×	D8	Ø	D9	Ù	DA	Ú	DB	Û	DC	Ü	DD	Ý	DE	Þ	В
ΕO	à	E1	á	E2	â	E3	ã	EЧ	ä	E5	å	E6	æ	E7	Ç	E\$	è	E9	é	ΕA	ê	EB	ë	EC	ì	ED	í	EE	î	EF
FO	ð	F1	ñ	F2	ò	F3	ó	F4	ô	F5	õ	F6	ö	F7	÷	F8	Ø	F9	ũ	FA	ú	FB	û	FC	ü	FD	ý	FE	Þ	₩.ÿ

C'est le jeu que nous utilisons avec nos PC français sous Windows NT. Cependant il vaut mieux s'en tenir, pour une page web par exemple, au jeu ISO 8859-1. En effet, si Netscape sous Windows reconnaît ces caractères étendus, Netscape sous Unix affiche de faux caractères.

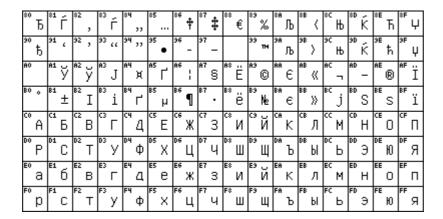
#### 3.3.2 - le jeu de caractères CP1250 ou WinLatin 2

Il s'agit de la norme ISO 8859-2 avec les caractères étendus x80 à x9F.

80 :	€			82	,			84	,,	85		86	‡	87	#			89	%	8A	š	88	<	<b>\$</b> C	Ś	\$D	Ť	\$E	ž	\$F .	ź
		91	í	92	,	93	"	94	"	95	•	96	-	97	-			99	тн	9A	š	9B	>	9C	ś	9D	ť	9E	ž	9F	ź
AO		A1	Ÿ	A2	Ü	A3	Ł	A4	Ħ	A5	Ą	A6	1	A7	S	A8	"	A9	0	AA	Ş	AB	((			AD	-	AE	®	AF .	ż
BO	۰	B1	±	BZ	Ţ	B3	ł	вч	-	B5	μ	B6	1	В7		B8	,	B9	ą	BA	Ş	BB	<b>&gt;&gt;</b>	BC	Ľ	BD	~	BE	ĭ	BF ;	ż
	Ŕ	C1	Á	CZ	Â	C3	Ă	СЧ	Ä	C5	Ĺ	C6	Ć	C7	Ç	C8	Č	С9	É	CA	Ę	СВ	Ë	СС	Ě	CD	Í	CE	Î	CF :	Ď
D0 {	Đ	D1	Ń	D2	Ň	D3	Ó	D4	ô	D5	ő	D6	ö	D7	×	D8	Ř	D9	Ů	DA	Ú	DB	Ű	DC	Ü	DD	Ý	DE	Ţ	DF 	В
E0	ŕ	E1	á	E2	â	E3	ă	EЧ	ä	E5	ĺ	E6	ć	E7	Ç	E\$	č	E9	é	EΑ	ę	EB	ë	EC	ě	ED	í	EE	î	EF I	ď
F0	đ	F1	ń	F2	ň	F3	ó	F4	ô	F5	ő	F6	ö	F7	÷	F8	ř	F9	ů	FA	ú	FB	ű	FC	ü	FD	ý	FE	ţ	FF	

#### 3.3.3 - le jeu de caractères CP1251 ou WinCyrillic

Nous eussions aimé qu'il s'agît de la norme ISO 8859-5 avec simplement des caractères étendus, mais voilà, toutes les lettres sont décalées (l'alphabet commence en xC0 au lieu de commencer en xB0). Elle est donc totalement incompatible à la fois avec l'ISO 8859-5 et la norme soviétique KOI-8.



#### 3.3.4 - le jeu de caractères CP1257 ou WinBaltic

Il s'agit de la norme ISO 8859-13 avec les caractères étendus x80 à x9F.

80	€			82	,			84	,,	85		86	#	87	#			89	%			8B	(			\$D		\$E	v	8F	_
		91	í	92	,	93	"	94	"	95	•	96	-	97	-			99	тн			9B	>			9D	-	9E	Ţ		
AO				A2	ф	A3	£	A4	Ħ			A6	1	A7	S	AS	Ø	A9	0	AA	Ŗ	AB	((	AC	7	AD	-	AE	®	AF	Æ
BO	۰	B1	±	B2	2	В3	3	вч	-	B5	μ	B6	1	B7		В8	Ø	В9	1	BA	ŗ	BB	<b>&gt;&gt;</b>	BC	14	BD	X	BE	¥	BF	æ
CO	Ą	C1	Į	CZ	Ā	C3	ć	СЧ	Ä	C5	Å	Ce	Ę	C7	Ē	C8	č	C9	É	CA	ź	СВ	Ė	СС	Ģ	CD	Ķ	CE	Ī	CF	Ļ
DO	Š	D1	Ń	DZ	Ņ	D3	Ó	D4	Ō	D5	õ	De	ö	D7	×	D8	Ų	D9	Ł	DA	Ś	DB	Ū	DC	Ü	DD	ż	DE	ž	DF	В
ΕO	ą	E1	į	E2	ā	E3	ć	EЧ	ä	E5	å	E6	ę	E7	ē	E\$	č	E9	é	ΕA	ź	EB	ė	EC	ģ	ED	ķ	EE	ī	EF	ļ
FO	š	F1	ń	F2	ņ	F3	ó	F4	ō	F5	õ	F6	ö	F7	÷	F8	ų	F9	ł	FA	ś	FB	ū	FC	ü	FD	ż	FE	ž	FF	•

#### 3.3.5 - et cætera

Vous l'avez compris, Microsoft couvre les mêmes alphabets que les diverses normes ISO 8859-N en exploitant les caractères supplémentaire x80-x9F. Nous avons ainsi CP1253 ou WinGreek, CP1254 ou WinTurkish, CP1255 ou WinHebrew qui sont compatibles avec les jeux ISO 8859-N correspondants, et CP1256 ou WinArabic qui ne l'est pas (à l'instar de WinCyrillic).

Il existe également le jeu CP1258 ou WinVietnamese qui n'est pas compatible avec le viscii, mais qui l'est avec l'ascii.

## 3.4 - la norme ISIRI 3342

Cette norme code l'écriture arabe modifiée employée pour écrire le persan. Il s'agit d'une norme nationale iranienne, et non d'une norme ISO.

## 4. les algorithmes UTF

UTF signifie Unicode Transformation Format. Les algorithmes UTF décrivent des ruses pour coder Unicode (qui est sur 2 octets) en utilisant, autant que possible, un seul octet.

Le système (je simplifie) est le suivant : 6 bits sont utilisés pour coder normalement les lettres de l'alphabet latin et les deux bits restants servent à indiquer, pour des caractères exotiques, le nombre de bits parmi ce qui suit utilisés pour leur codage.

#### 4.1 - UTF-1

C'est le plus simple des formats UTF. Son inconvénient est que, en cas de rupture du flux de données, on ne sait pas retrouver où on en était.

## 4.2 - UTF-2 rebaptisé plus tard UTF-8

Ce format utilise un nombre de bits standardisé par rapport à l'UTF-1 (toujours un multiple de 8), ce qui permet de reprendre un codage interrompu.

L'intérêt de l'UTF-8 est sa capacité à préserver les standards existants comme l'ASCII à 7 bits. Sa principale caractéristique est d'utiliser une méthode de codage variable, entre 1 et 6 octets. Les premiers octets constituent l'image de l'ensemble US-ASCII (7 bit), i.e. x00-x7F. Le premier octet du caractère (il peut y en avoir jusqu'à 6 !) définit la **séquence** à laquelle il appartient. Si cet octet commence par 0, il définit alors le codage US-ASCII et sa longueur est 1. Si les n (n > 1) premiers bits sont égaux à 1, le n+1 étant égal à 0, alors la longueur du caractère sera n. Les bits restants définissent la séquence de codage.

Normalement, tous les outils X11 et GNU développés après le 1<sup>er</sup> janvier 1999 doivent intégrer le format UTF-8.

Les désavantages principaux de ce format sont :

- le codage ISO 8859-1 n'est pas reconnu de façon transparente, et il n'est pas pensable de passer à la moulinette l'énorme quantité de documents déjà existants.
- tous les caractères n'ont pas la même longueur en octets.

#### 4.3 - UTF-7,5

Ce format utilise sept bits au lieu de six pour coder l'alphabet latin et permet la transparence par rapport à l'ISO 8859-1. L'inconvénient qui en résulte est la plus grande taille en octets nécessaire au codage des alphabets non-latins.

Ce format n'est pas reconnu comme étant une norme par l'ISO.

#### 4.4 - UTF-7

Ce format reprend la philosophie de l'UTF-8 mais n'utilise pas le 8<sup>ème</sup> bit pour être compatible avec certains mailers. Il semble être quasiment aussi efficace que l'UTF-8.

#### 4.5 - SCSU

Ce format (Standard Compression Scheme for Unicode) reprend la philosophie de l'UTF-8 mais permet la transparence par rapport à l'ISO 8859-1. Il nécessite des algorithmes très lourds.

## 5. les normes sur plus de 8 bits

Ces normes sont nées en Extrême-Orient, où les utilisateurs des *milliers* d'idéogrammes nécessaires à l'écriture des langues du groupe CJK ont bien ri quand ils ont vu que 2<sup>8</sup> ça ne faisait que 256. Exit donc les ISO 8859-N, et place à des normes spécifiques à chaque pays et à son système d'écriture.

#### 5.1 - la norme JIS X 0208

Comme d'habitude, les Japonais furent les plus rapides. La norme JIS X 0208, de 1976, est la toute première norme sur deux octets. Il s'agit d'une norme nationale japonaise, et non d'une norme ISO. Elle code les alphabets latin, cyrillique et grec, les syllabaires japonais (**kana**) et un grand nombre d'idéogrammes chinois utilisés en japonais (**kanji**).

Cette norme aurait pu, en fait, servir de base à la norme plus complète Unicode. Cependant, pour des raisons de compatibilité avec des imprimantes, seuls 8836 caractères peuvent être utilisés.

#### 5.2 - la norme GB 2312 ou Guo-Biao

Cette norme code les idéogrammes chinois simplifiés en usage en RPC et à Singapour. Il s'agit d'une norme nationale de la RPC.

La première mouture de cette norme codait 6763 idéogrammes ; la mouture actuelle code plus de 26000 idéogrammes.

## 5.3 - le protocole HZ (HanZi)

Il s'agit d'une adaptation de la norme précédente aux mailers. Elle n'utilise donc que deux fois sept bits.

#### 5.4 - la norme CNS 11643

Cette norme code les idéogrammes chinois traditionnels en usage à Taïwan et à Hong Kong. Il s'agit d'une norme nationale taïwanaise. Cette norme code 48027 idéogrammes.

## 5.5 - la norme KS C 5601

Cette norme code le syllabaire coréen (hangul) et les idéogrammes chinois (parfois) utilisés en coréen (hanja).

#### 5.6 - le format EUC

Extended Unix Code utilisé sur certaines stations de travail. Ce format regroupe en fait le JIS X 0208 japonais, le GB 2312 chinois simplifié, le CNS 11643 chinois traditionnel, et le KS C 5601 coréen.

#### 5.7 - le codage Big5

Il ne s'agit en fait pas d'une véritable norme, mais du standard industriel imposé de facto par l'industrie taïwanaise pour le codage des idéogrammes chinois traditionnels. Il est complètement différent de la norme taïwanaise CNS 11643 et sa méthode de codage ressemble à celle de la norme GB 2312 utilisée

en RPC pour les idéogrammes simplifiés, ce qui permet de passer de l'une à l'autre avec des scripts. Ce standard a le désavantage d'utiliser deux fois 8 bits au lieu de deux fois 7 bits (problèmes de mailers). Ce standard permet de coder 13051 idéogrammes traditionnels.

#### 5.8 - la norme Guohui

Elle est très semblable à la norme précédente et a l'avantage d'être une norme officielle taïwanaise.

#### 5.9 - UTF-16

Ce format code tous les caractères sur 16 bits. Point.

Un calcul rapide permet de s'apercevoir qu'il permet ainsi de coder 65536 caractères, c'est-à-dire presque tous les systèmes d'écriture du monde, y compris certains dont vous n'avez jamais entendu parler comme les hiéroglyphes méroïtiques ou les runes magyares.

Ce format suit fidèlement les directives de la norme Unicode ou ISO 10646 pour ce qui est de l'ordre des caractères, aboutissant ainsi à une norme *véritablement universelle*.

Ces caractères sont groupés en blocs indivisibles de 128, pour faciliter le classement par systèmes d'écriture.

Exemple : bloc des caractères de x3480 à x34FF

傏	偖	樲	然	僡	傓	僬	觙	僷	歰	燍	儝	傑	靡		
3480	3481	3482	3483	3484	3485	3486	3487	3488	3489	348A	348B	348C	348D	348E	348F
儮	儻	燵	燭	瘫	儶	儘			綿	濦	嚁	嬹	僟	儶	僡
3490	3491	3492	3493	3494	3495	3496	3497	3498	3499	349A	349B	349C	349D	349E	349F
<b>外</b> 34A0	34A1	鱂	儡	僼	俚细	離	傻	像 348	露3449	<b></b>	34AB	贶34AC	34AD	榫	$\frac{1}{34AF}$
3480	34B1	34B2	34B3	<u></u>	台 34B5	<u></u>	<b>家</b> 34B7	真真 3488	<b>送</b>	見 34BA	<b>苏</b> 34BB	署3480	施 34BD	<b>副</b> 第 34BE	34BF
	越	μU	皴	汀	0405	沃	冱	活	洞	洪	凎	0400	浸	凋	淕
3400	34C1	34C2	34C3	34C4	34C5	3406	34C7	34C8	34C9	34CA	34CB	34CC	34CD	34CE	34CF
<b>幸</b>	34D1	南3402	谢3403	滅	滭	潔學	<b>凰</b>	34D8	] 34D9	IJ 34DA	<b>幻</b> 34DB	荆 34DC	<b>韧</b> 34DD	34DE	34DF
3400	34DT	3402	3403	3404	3400	3400	34D7	34D0	3403	34DA	3400	3400	3400	34DE	34DF
<b>制</b> 34E0	<b>劉</b>	<b></b> 34E2	<u>計</u>	34E4	34E5	割34E6	刾	觓	34E9	34EA	刹34EB	<b>剔</b> 34EC	則 34ED	劉 34EE	34EF
創 34F0	<b>副</b> 34F1	<b>匙</b> 34F2	<b>则</b> 34F3	<u> </u>	34F5	割34F6	<b>盟</b> [	<b>厕</b> 34F8	<b>劇</b> 34F9	<b>凯</b> 34FA	刻34F8	割345公	<b>剛</b>	34FE	駅 34FF

#### 5.10 - UCS-4

Ce format est le même que le précédent, mais sur 31 bits. Il peut être compris par les systèmes informatiques actuels mais n'est guère utilisé en raison de sa lourdeur (plein de 0 inutiles pour les alphabets européens).

## 5.11 - le jeu de caractères étendu WGL4

Il s'agit d'un jeu de caractères étendu de chez Microsoft (du même type de ceux présentés au paragraphe 3.3) mais conforme à la norme ISO 10646. Il combine en un seul jeu tous les caractères des jeux 1250 à 1254. Certaines polices True Type pour Windows NT sont normalement WGL4 : Times New Roman, Lucida Unicode et Verdana notamment.

## 6. les normes descriptives ou de balisage

Les normes descriptives ne codent pas directement les caractères mais les "décrivent" au sein d'un document. HTML en est un exemple bien connu. Comme tous les systèmes ne codent pas de la même façon les guillemets, il vaut mieux remplacer, dans le code HTML d'une page web, les occurrences de

par

## "

qui sera universellement reconnu par tous les fureteurs. Il en va de même pour d'autres signes tipographiques comme &, @, etc.

#### 6.1 - la norme ISO 8879

Cette norme décrit les jeux de caractères Latin 1, Latin 2, Grec et Cyrillique, mais pour les systèmes **SGML**, un langage de balisage destiné à décrire la mise en forme de documents partageables par des systèmes informatiques différents.

#### 6.2 - la norme ERCS

Extended Reference Concrete Syntax for SGML. Cette norme décrit les jeux de caractères d'Extrême-Orient (CJK) pour les systèmes SGML.

#### 6.3 - Canonical XML

Il existe une sous-spécification de **XML** appelée Canonical XML qui doit normalement être compatible avec l'UTF-8.

## 7. les polices de caractère True Type

Jusqu'à Windows 95, les polices True Type étaient codées sur 8 bits et permettaient donc d'afficher 256 caractères dans une police donnée. Reprenant la philosophie des jeux de caractères Windows, les polices True Type affichent l'us-ascii sur les 128 premiers caractères, des caractères étendus Windows de x80 à x9F, et des caractères nationaux sur les caractères suivants.

Il existe ainsi des polices correspondant aux divers jeux Windows CP1250 à 1258 et, ce qui est intéressant, des polices correspondant à des systèmes d'écriture qui n'ont rien d'ISO ou pour lesquels il n'existe pas de codage en-dehors d'Unicode : le vietnamien, les langues amérindiennes, l'alphabet phonétique international, etc.

Pour créer une page web en vietnamien, par exemple, il suffit de la coder en utilisant une police True Type vietnamienne et de la mettre à disposition sur son site pour que les internautes puissent la télécharger.

Inconvénient : cela marchera très bien sous Windows, mais les internautes sous Unix ou MacOS ne comprendront rien...

À partir de Windows NT, les polices True Type sont codées sur 16 bits et permettent donc d'afficher une partie des caractères universels de la norme ISO 10646. On peut télécharger plusieurs polices de cette sorte sur le site de Microsoft (654 caractères), et également sur le site de la société Bitstream (plus de 8500 caractères).

## 8. le codage d'une page web (HTML)

Lorsque votre fureteur charge une page HTML, il va l'afficher avec le codage par défaut de votre système d'exploitation (pour nous, ISO 8859-1). D'où le genre de problèmes décrits au paragraphe 2.2.

La solution à ce problème est d'ajouter une ligne forçant le codage des caractères dans la partie comprise entre les tags <HEAD> et </HEAD> du code HTML de la page.

#### Exemples

UTF-7.

```
<meta http-equiv="Content-Type" content="text/html; charset=iso-8859-
1"> pour forcer l'affichage en ISO 8859-1.

<meta http-equiv="Content-Type" content="text/html; charset=windows-
1250"> pour forcer l'affichage en WinLatin 2.

<meta http-equiv="Content-Type" content="text/html; charset=KOI8-R">
pour forcer l'affichage en cyrillique (norme russe KOI8-R).

<meta http-equiv="Content-Type" content="text/html; charset=iso-2022-kr"> pour forcer le codage d'après la norme coréenne KS C 5601.

<meta http-equiv="content-type" content="text/html; charset=x-
UNICODE-2-0-UTF-7"> pour forcer l'affichage en Unicode avec l'algorithme de transformation
```

Il est également possible de choisir, à la main, l'affichage correspondant à la page chargée.

Sous Internet Explorer, il faut cliquer sur le menu *Affichage*, puis sur l'item *Codage*. Ma version (5.00) propose une quantité impressionnante de codages, à savoir : ISO 8859-1, 2, 5, 6, 7, 8, 9, 13 ; CP1250 à 1258 ; Chinois simplifié, Chinois traditionnel, Japonais, Coréen, Thaï (sans plus de détails) ; KOI8-R et U ; les anciens DOS Europe Centrale, arabe, cyrillique, hébreu ; ASMO 708 ; UTF-8.

Remarque : d'après mes sources, l'ISO 8859-6 et l'ASMO 708 sont la même chose, mais Internet Explorer propose quand même les deux !

Sous Netscape, il faut cliquer sur le menu *Afficher* et sur l'item *Encodage*. Ma version (4.5) propose les codages suivants : ISO 8859-1, 2, 5, 7, 9 ; CP1250, 1251, 1253 ; Japonais ou Coréen Auto-Detect (qu'es aquò?), Shift\_JIS et EUC-JP (deux implémentations légèrement différentes de JIS X 0208), Big5, EUC-TW (càd CNS 11643), GB 2312 ; KOI8-R ; UTF-7, UTF-8.

Le problème principal est que le codage ainsi défini s'applique à toute la page, d'où l'impossibilité d'afficher sur la même page un paragraphe en français et un en tchèque, par exemple, à moins de coder la page en Unicode, ce qui est beaucoup plus fastidieux, car au lieu de taper directement le caractère avec un éditeur de texte, il faut le décrire de la façon suivante

## & # D;

où  ${\it D}$  représente le code (unique) du caractère dans la norme ISO 10646, en notation décimale. Exemples :

# Annexe 1 - Codage des systèmes d'écriture pour les langues europeennes et les principales autres langues du monde

Langue	ISO 639-1	ISO 639- 2/T	ISO 639- 2/B	ISO 15924	codage (*)
albanais	sq	sqi	alb	la	ISO 8859-1, ISO 8859-2
allemand	de	deu	ger	la	ISO 8859-1
anglais	en	eng	eng	la	ISO 8859-1
arménien	hy	hye	arm	hy	UTF-7, UTF-8
azéri (†)	az	aze	aze	la	ISO 8859-9 (sauf ə), UTF-7, UTF-8
basque	eu	eus	baq	la	ISO 8859-1
biélorusse	be	bel	bel	су	ISO 8859-5
oreton	br	bre	bre	la	ISO 8859-1
bulgare	bg	bul	bul	су	ISO 8859-5
catalan	ca	cat	cat	la	ISO 8859-1
corse	co	cat	cos	la	ISO 8859-1
	hr	hrv		la	ISO 8859-2
croate danois	da	dan	dan	la	ISO 8859-1
				la	ISO 8859-1
erse	gd	gdh	gae	la	ISO 8859-1
espagnol	es	esp	spa		
estonien	et fo	est	est	la	ISO 8859-4, ISO 8859-10
féroïen		fao	fao	la	ISO 8859-1
finnois	fi	fin	fin	la	ISO 8859-1
français	fr	fra	fre	la	ISO 8859-1 (sauf œ), ISO 8859-15
frison	fy	fry	fry	la .	ISO 8859-1
galicien	gl	glg	glg	la	ISO 8859-1
gallois	cy	cym	wel	la	ISO 8859-1 (sauf w accentués), ISO
					8859-14
géorgien	ka	kat	geo	ka	UTF-7, UTF-8
grec	el	ell	gre	el	ISO 8859-7
hongrois	hu	hun	hun	la	ISO 8859-2
irlandais	ga	gai	iri	la, lg	ISO 8859-1
islandais	is	isl	ice	la	ISO 8859-1
italien	it	ita	ita	la	ISO 8859-1
letton	lv	lav	lav	la	ISO 8859-4
lituanien	lt	lit	lit	la	ISO 8859-4, ISO 8859-10
luxembourg eois	lb	ltz	ltz	la	ISO 8859-1
macédonien	mk	mkd	mac	су	ISO 8859-5
maltais	mt	mlt	mlt	la	ISO 8859-3
néerlandais	nl	nld	dut	la	ISO 8859-1
occitan	ос	oci	oci	la	ISO 8859-1
polonais	pl	pol	pol	la	ISO 8859-2
portugais	pt	por	por	la	ISO 8859-1
rhéto-roman	rm	roh	roh	la	ISO 8859-1
roumain	ro	ron	rum	la	ISO 8859-2
russe	ru	rus	rus	су	ISO 8859-5
serbe	sr	srp	scc	cy	ISO 8859-5
slovaque	sk	slk	slo	la	ISO 8859-2
slovàque	sl	slv	slv	la	ISO 8859-2
suédois	SV	swe	swe	la	ISO 8859-1
tchèque	cs		cze	la	ISO 8859-2
•		ces			ISO 8859-2
turc ukrainien	tr	tur	tur	la	ISO 8859-5
ukiaiiileli	uk	ukr	ukr	cy	130 0037-3

afrikaans	af	afr	afr	la	ISO 8859-1
arabe	ar	ara	ara	ar	ISO 8859-6
bengali	bn	ben	ben	bn	UTF-7, UTF-8
esperanto	eo	еро	еро	la	ISO 8859-3
chinois (1)	zh	zho	chi	ha	GB 2312, HZ
chinois (2)	zh	zho	chi	ha	Big5, CNS 11643, Guohui
coréen (3)	ko	kor	kor	hg	KS C 5601
coréen (4)	ko	kor	kor	kh	KS C 5601
haoussa	ha	hau	hau	la	UTF-7, UTF-8
hébreu	he	heb	heb	he	ISO 8859-8
hindi	hi	hin	hin	dv	UTF-7, UTF-8
indonésien	id	ind	ind	la	ISO 8859-1
japonais	ja	jpn	jpn	ja	JIS X 0208
malais	ms	msa	may	la	ISO 8859-1
marathe	mr	mar	mar	dv	UTF-7, UTF-8
ouolof	wo	wol	wol	la	ISO 8859-1 (sauf ŋ), UTF-7, UTF-8
ourdou	ur	urd	urd	ar	UTF-7, UTF-8
persan	fa	fas	per	ar	ISIRI 3342
souahéli	sw	swa	swa	la	ISO 8859-1
tagalog	tl	tgl	tgl	la	ISO 8859-1
tamoul	ta	tam	tam	ta	UTF-7, UTF-8
télougou	te	tel	tel	te	UTF-7, UTF-8
thaï	th	tha	tha	th	ISO 8859-11
vietnamien	vi	vie	vie	la	Aucun satisfaisant (5)
xhosa	xh	xho	xho	la	ISO 8859-1
yorouba	(6)	yor	yor	la	UTF-7, UTF-8
zoulou	zu	zul	zul	la	ISO 8859-1

## Remarques:

- (\*) il s'agit du codage par défaut ou officiel. Il est évident qu'on peut coder l'anglais d'après n'importe quelle autre norme ISO 8859-N, par exemple !
- $(\dagger)$  la voyelle azérie ə peut être remplacée par ä
- (1) RPC et Singapour
- (2) Taïwan et Hong Kong
- (3) Corée du Nord
- (4) Corée du Sud
- (5) Aucune séquence d'Unicode ne décrit le vietnamien ! Un vrai scandale, il n'y aucun moyen de le coder de façon universelle.
- (6) le yorouba n'apparaît pas dans la norme ISO 639-1.

## Annexe 2 - Glossaire

**hangul** : syllabaire phonétique pouvant noter tous les sons de la langue coréenne. En Corée du Nord, la langue est notée avec le hangul seul ; en Corée du Sud, avec un mélange de hangul et de hanja.

**hanja** : idéogrammes chinois utilisés en Corée du Sud pour noter les mots de façon étymologique plutôt que phonétique.

**HTML**: Hyper Text Mark-up Language. Code permettant d'afficher une page web de la même façon indépendamment du système d'exploitation et du fureteur utilisés (on y croit tous).

ISO: Organisation Internationale de Normalisation.

ISO 15924 : norme provisoire qui spécifie les digraphes pour les systèmes d'écriture du monde.

ISO 639-1 : norme ISO qui spécifie les digraphes pour les principales langues du monde.

**ISO 639-2**: norme ISO qui spécifie les trigraphes pour la plupart des langues du monde. Malheureusement, la norme est elle-même subdivisée en /T (terminologie) et /B (bibliographie). D'autres normes (comme la norme DVB-SI) n'ont pas pris la peine de spécifier s'il fallait utiliser la colonne T ou la colonne B!

**kana** : syllabaires phonétiques pouvant noter tous les sons de la langue japonaise. Il en existe deux : le *hiragana* pour les mots japonais, et le *katakana* pour les emprunts aux langues étrangères. Chaque série compte presque 70 caractères.

**kanji**: idéogrammes chinois utilisés au Japon pour noter les mots de façon étymologique plutôt que phonétique. Les réformes des années 50 et 70 ont restreint le nombre de kanji essentiels à 2000 : 881 idéogrammes de base (kyôiku kanji)

941 idéogrammes d'usage courant (tôyô kanji)

En tout 1822, donc, plus 188 autres plus rares utilisés par exemple pour écrire les noms de famille.

KOI: Code pour l'échange de données. Norme soviétique pour le codage de l'alphabet cyrillique.

**séquence** : les 65536 caractères potentiels d'Unicode sont divisés en 256 séquences (ou *rangées*) correspondant grosso modo chacune à un ensemble de systèmes d'écriture similaires. La séquence x00 correspond à l'ISO 8859-1, la séquence x04 correspond aux alphabets cyrilliques, les idéogrammes couvrant, eux, plusieurs séquences!

**SGML** : Standard Generalised Mark-up Language. Le language de référence pour les documents en ligne. Le HTML est une sous-partie du SGML.

UCS: le standard Unicode est basé sur le système de codage défini dans la norme ISO 10646, qui ellemême définit deux méthodes de codage à 16 et 31 bits, respectivement : 10646.UCS-2 (16 bits) et 10646.UCS-4 (31 bits). UCS signifie Universal Character Set (ensemble universel de caractères).

**Unicode** : la norme qui permet de coder tous les systèmes d'écriture du monde. Actuellement, elle répertorie officiellement 38885 caractères issus de 25 systèmes d'écriture différents.

WG2: Working Group 2 - Les gens de l'ISO qui travaillent sur la norme ISO 10646.

**WG3**: Working Group 3 - Les gens de l'ISO qui travaillent sur les normes à 7 et à 8 bits. La norme ISO 2375 donne la liste des normes de codage reconnues (une norme de normes... ça devient puissant).

**XML** : Extensible Mark-up Language. Version simplifiée du SGML destinée au world wide web. On peut trouver une FAQ à l'URL suivante : http://www.ucc.ie/xml/

## **Annexe 3 - Sources**

J'ai surtout mis à contribution les deux sites suivants pour le contenu technique de ce document : www.czyborra.com

et

http://www.gms.lu/~fr\_jap/unic\_std.htm