

Architecture des ordinateurs

4 - Représentation de l'information en machine

Définitions de base

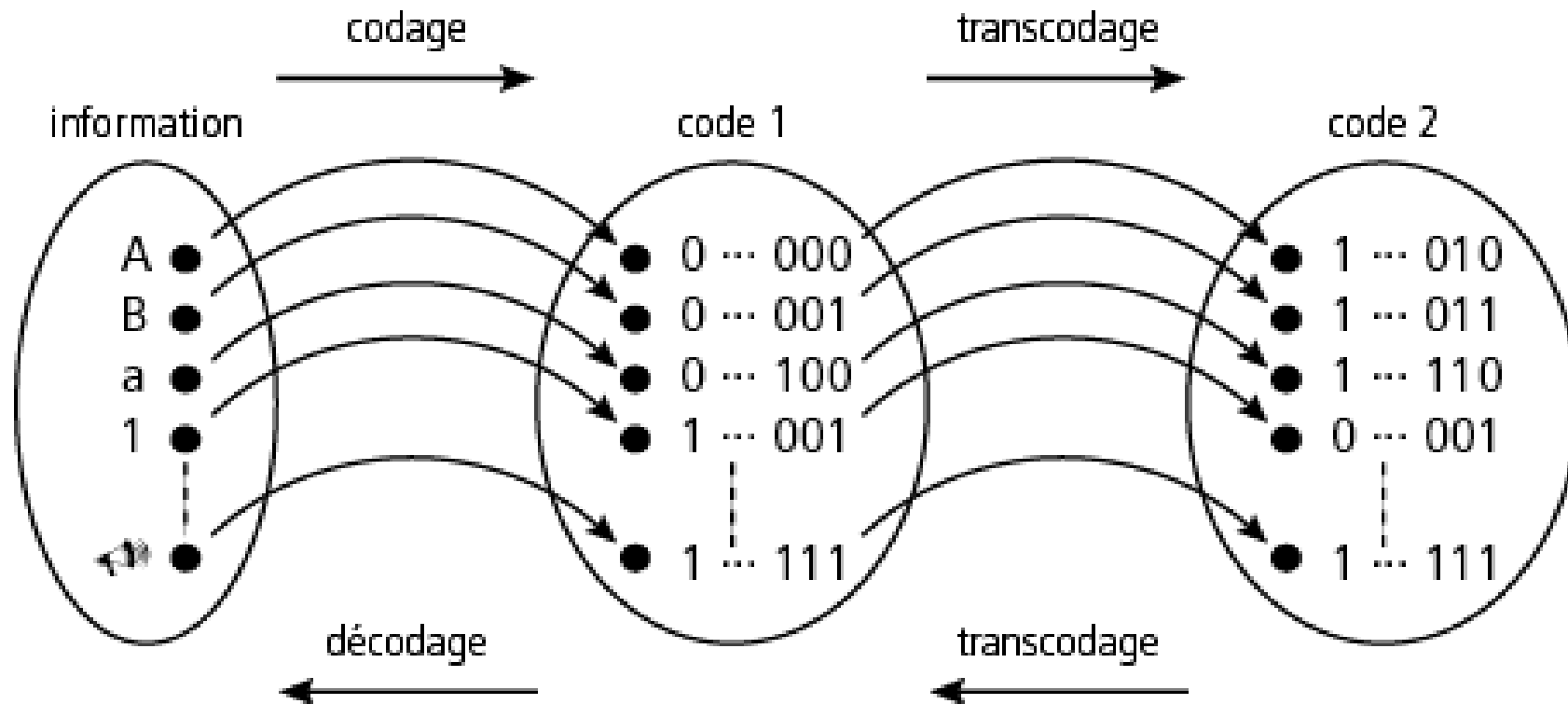
Les codes alphanumériques

Philippe Darche
IUT Paris Descartes

Le problème

- ❑ Pour écrire l'anglais, 26 lettres, 10 chiffres et quelques caractères de ponctuation sont suffisants
- ❑ Pour les langue européennes comme le français, il y a en plus d'autres caractères ou variations sur les caractères
- ❑ Il existe d'autres alphabets ou certains pays utilisent des idéogrammes
- ❑ **Comment coder ces symboles en machine ?**

La notion de code

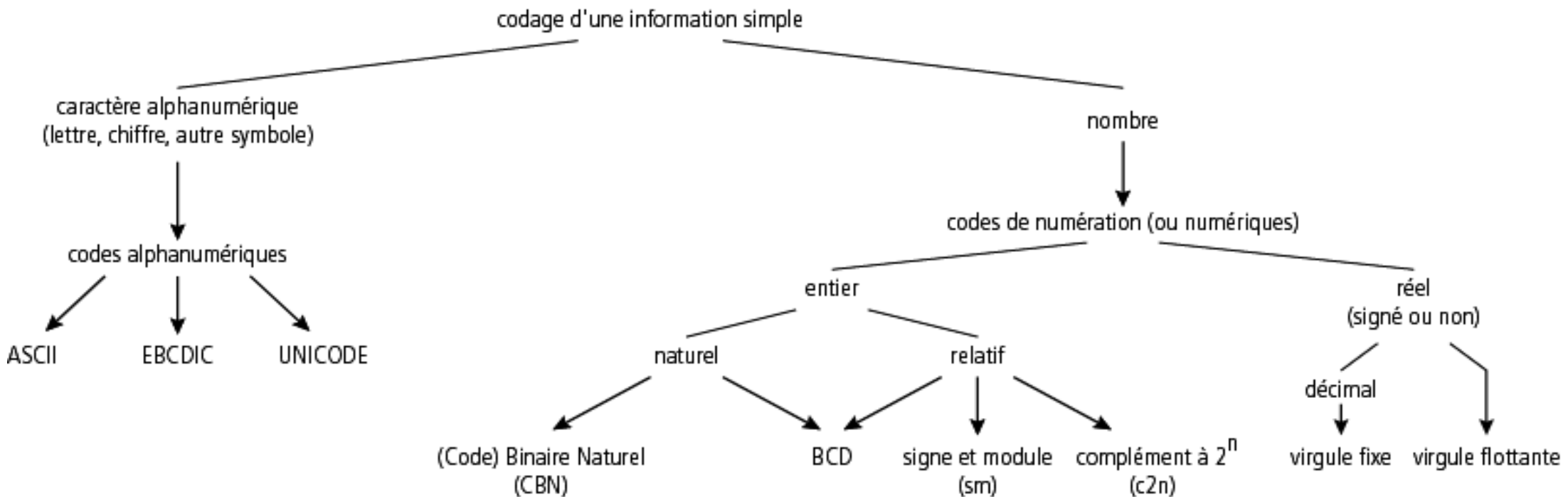


Définitions

- Code = {mot de code au format n}
 - nombre de mots-code possible = 2^n
 - rendement $\eta = \frac{\text{Card}_C}{2^n}$
- Codage = loi de correspondance arbitraire
- Opération inverse : le décodage
- Types de code :
 - alphanumérique
 - de numération ou numérique
 - instruction
 - de canal
 - détecteur (et correcteur) d'erreur(s)
 - de chiffrement
 - etc.

Représentation de l'information en machine

- Codage des deux principales familles avec codes usités



Code alphanumérique

- Caractères alphanumériques
 - contraction des mots *alphabétique* et *numérique*
 - ⇒ caractère alphabétique, numérique ou, éventuellement, caractère spécial (ponctuation, technique (monnaie, etc.) ou caractère espace (*space*))
(définition de la norme ISO/IEC 2382-5)
- Codage des caractères alphanumériques (*graphic character*)
 - En anglais, *character set*
 - loi de correspondance biunivoque
 - et éventuellement des codes de contrôle de périphérique (*control function*)
 - retour chariot (*Carriage Return* ou CR), saut à la ligne (*Line Feed* ou LF), etc.

Remarque importante

- ❑ Les codes de contrôle ne s'affichent pas.
- ❑ Ils sont interprétés par le périphérique.

Ne pas confondre

- Caractère graphique
 - information à coder
 - exemple : lettre capitale latine a
- Symbole graphique = forme graphique
 - représentation visuelle d'un caractère graphique ou d'une fonction de contrôle
 - autre appellation : glyphe
- Police de caractères = {symboles graphiques}

Le code ASCII d'origine

- *American Standard Code for Information Interchange*
- Inventé par l'américain Bob Bemer [Bemer 61]
- Code au format $n = 7$ bits
 - si octet, un bit de parité pour le contrôle d'erreur
- Pas de minuscule
- Codage des fonctions de contrôle

La table originale standard ASA X3.4-1963

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

contrôle

Évolution de la table originale ASCII

- Prise en compte des minuscules

Table E-1. AMERICAN STANDARD CODE FOR INFORMATION INTERCHANGE

		128-SYMBOL PRINTING SET								
COLUMN →		0	1	2	3	4	5	6	7	ROW ↓
b7 b6 b5		0	0	0	0	1	1	1	1	
b4 b3 b2 b1		NON-PRINTING (9)				96-SYMBOL PRINTING SUBSET				
0 0 0 0	NUL \	DLE @	SP	0	@	P		p	0	
0 0 0 1	SOH ¶	DC1 @	!	!	A	Q	a	q	1	
0 0 1 0	STX ¶	DC2 @	"	2	B	R	b	r	2	
0 0 1 1	ETX ¶	DC3 @	#	3	C	S	c	s	3	
0 1 0 0	EOT ¶	DC4 @	\$	4	D	T	d	t	4	
0 1 0 1	ENQ ¶	NAK @	%	5	E	U	e	u	5	
0 1 1 0	ACK ¶	SYN @	8	6	F	V	f	v	6	
0 1 1 1	BEL ¶	ETB @	'	7	G	W	g	w	7	
1 0 0 0	BS ¶	CAN @	(8	H	X	h	x	8	
1 0 0 1	HT ¶	EM @)	9	I	Y	i	y	9	
1 0 1 0	LF ¶	SUB @	*	:	J	Z	j	z	10	
1 0 1 1	VT ¶	ESC @	+	;	K	[k	{	11	
1 1 0 0	FF ¶	FS @	.	<	L	\	l	!	12	
1 1 0 1	CR ¶	GS @	-	=	M]	m	}	13	
1 1 1 0	SO ¶	RS @	.	>	N	^	n	~	14	
1 1 1 1	SI ¶	US @	/	?	O	_	o	DEL	15	

contrôle

La norme ISO 8859-x

- Le code des micro-ordinateurs
- Code ASCII complété pour les valeurs 128 à 255
 - format $n = 8$ bits
 - prise en compte des besoins des langues européennes
 - ⇒ nécessité de passer au format octet
 - normalisé par l'*International Organization for Standardization* (ISO)
 - ensemble latin 1 : ISO 8859-1:1998
 - prise en compte des caractères accentués d'Europe occidentale
 - version avec signe euro : ISO 8859-15

La norme ISO/CEI 8859-1

- 191 caractères graphiques codés
- Version avec codes de contrôle
 - ISO-8859-1

				b8	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	
				b7	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1
				b6	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
				b5	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
					00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15
b4	b3	b2	b1	00			SP	0	@	P	`	p		NBSP	°	À	Ð	à	ø	0
0	0	0	1	01			!	1	A	Q	a	q		i	±	Á	Ñ	á	ñ	1
0	0	1	0	02			"	2	B	R	b	r		¢	²	Â	Ò	â	ò	2
0	0	1	1	03			#	3	C	S	c	s		£	³	Ã	Ó	ã	ó	3
0	1	0	0	04			\$	4	D	T	d	t		¤	'	Ä	Ô	ä	ô	4
0	1	0	1	05			%	5	E	U	e	u		¥	µ	Å	Õ	å	õ	5
0	1	1	0	06			&	6	F	V	f	v		¦	¶	Æ	Ö	æ	ö	6
0	1	1	1	07			'	7	G	W	g	w		§	-	Ç	×	ç	÷	7
1	0	0	0	08			(8	H	X	h	x		"	,	È	Ø	è	ø	8
1	0	0	1	09)	9	I	Y	i	y		©	¹	É	Ù	é	ù	9
1	0	1	0	10			*	:	J	Z	j	z		º	º	Ê	Ú	ê	ú	A
1	0	1	1	11			+	;	K	Ç	k	ç		«	»	Ë	Û	ë	û	B
1	1	0	0	12			,	<	L	\	l			¬	¼	Ì	Ü	ì	ü	C
1	1	0	1	13			-	=	M	J	m	¸			½	Í	Ý	í	ý	D
1	1	1	0	14			.	>	N	^	n	~		®	¾	Î	Þ	î	þ	E
1	1	1	1	15			/	?	O	_	o			¯	¿	Ï	ß	ï	ÿ	F
				0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	hex

L'EBCDIC

- ❑ Origine : IBM
- ❑ format natif n = 8 bits

HEX DIGITS 1ST → 2ND ↓	0-	1-	2-	3-	4-	5-	6-	7-	8-	9-	A-	B-	C-	D-	E-	F-
-0	NUL 0010000	DLE SE170000			SP SP010000	& M030000	- SP100000					^ SD150000	{ SM110000	}	\	0 ND100000
-1	SOH SE020000	DC1 SE180000				/ SP120000			a LA010000	j LJ010000	~ SD190000		A LA020000	J LJ020000		1 ND010000
-2	STX SE030000	DC2 SE190000	FS SE350000	SYN SE230000					b LB010000	k LK010000	s LS010000		B LB020000	K LK020000	S LS020000	2 ND020000
-3	ETX SE040000	DC3 SE200000							c LC010000	l LL010000	t LT010000		C LC020000	L LL020000	T LT020000	3 ND030000
-4									d LD010000	m LM010000	u LU010000		D LD020000	M LM020000	U LU020000	4 ND040000
-5	HT SE100000		LF SE110000						e LE010000	n LN010000	v LV010000		E LE020000	N LN020000	V LV020000	5 ND050000
-6		BS SE190000	ETB SE240000						f LF010000	o LO010000	w LW010000		F LF020000	O LO020000	W LW020000	6 ND060000
-7			ESC SE280000	EOT SE050000					g LG010000	p LP010000	x LX010000		G LG020000	P LP020000	X LX020000	7 ND070000
-8		CAN SE250000							h LH010000	q LQ010000	y LY010000		H LH020000	Q LQ020000	Y LY020000	8 ND080000
-9		EM SE260000						^ SD130000	i LI010000	r LR010000	z LZ010000		I LI020000	R LR020000	Z LZ020000	9 ND090000
-A					! SP020000	:						[SM060000			FN2 SE400000	FN3 SE410000
-B	VT SE120000				. SP110000	\$ SC030000	, SP080000	# SM010000] SM080000				
-C	FF SE130000			DC4 SE1210000	< SA020000	* SA030000	% SA020000	@ SA0510000								
-D	CR SE140000	GS SE160000	ENQ SE060000	NAK SE220000	(SP060000) SP0710000	= SP090000	' SP0510000								
-E	SO SE150000	RS SE270000	ACK SE070000		+ SA010000	; SP140000	> SA0510000	= SA080000					FN4 SE420000			
-F	SI SE160000	US SE180000	BEL SE190000	SUB SE270000	 SO130000		? SP150000	" SP040000	FN1 SE190000							DEL SE330000

contrôle

Code Page 01303

Limitation des codes alphanumériques classiques

- Format limité à $n = 7$ ou 8 bits
 - nombre de caractères codables limité (2^n)
 - ⇒ impossibilité de prise en compte de toutes les lettres, signes de ponctuation et autres symboles d'une langue
- Si prise en compte de quelques langues
 - ⇒ problème pour l'échange mondiale d'informations non résolu
 - ⇒ transcodage nécessaire mais souvent impossible

La réponse aux limitations

- UNICODE pour *Universal Code*
- Codage de tous les caractères de toutes les langues
- Association pour chaque caractère abstrait d'un numéro de code (valeur scalaire, *code point*) et d'un nom standard
 - comme pour les codes précédents
 - $n = 16$ bits à l'origine, 21 bits actuellement
 - identique à la norme ISO/CEI 10646-1:1993
- Et représentation par un format d'encodage (*Unicode Transformation Format*)
 - UTF-32 (32 bits), UTF-16 (16 bits) et UTF-8 (8 bits)
 - UTF-8 le plus efficace en terme de stockage

Vocabulaire d'UNICODE

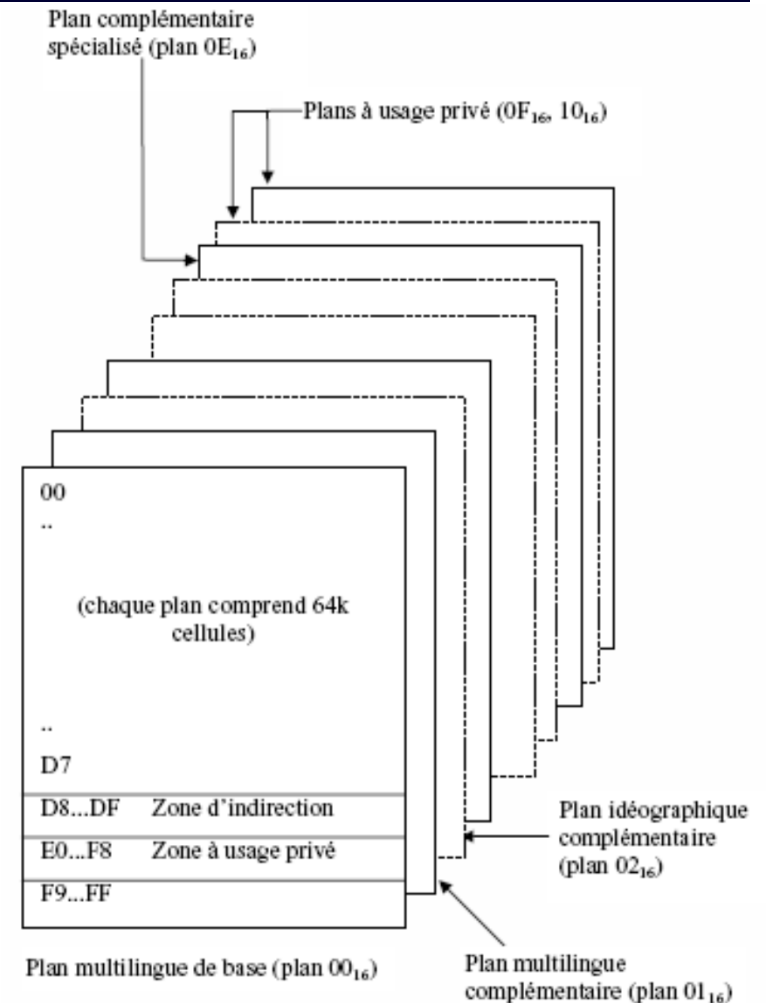
- Caractère : forme abstraite d'un composant atomique d'un langage
 - lettre, signe de ponctuation, etc.
 - exemple : U+0041 *latin capital letter a*
- Glyphe : forme graphique d'un caractère
 - police : ensemble de glyphes
 - de l'exemple précédent : A, ▲, A, etc.

UNICODE

- Conformité avec le standard ISO/IEC 10646 pour l'affectation du numéro de code et du nom standard
 - UCS-2 et UCS-4 (*Universal Character Set*)

Plans et zones de codage

- Plans de 64 K cellules
 - plan de base
(PMB pour *Basic Multilingual Plane*)
 - caractères usuels alphabétiques



Espace d'adressage

- 1 114 112 caractères possibles (*code point*)



Schémas d'encodage

- UTF-x (x = 8, 16 et 32)
 - *Unicode or UCS Transformation Format*
- Objectifs : adapter la valeur du code au format courant des données (format n = 8, 16 ou 32 bits) en optimisant l'encombrement

Schéma d'encodage UTF-32

- Correspondance directe (remplissage avec 0)

UNICODE (21 bits) 



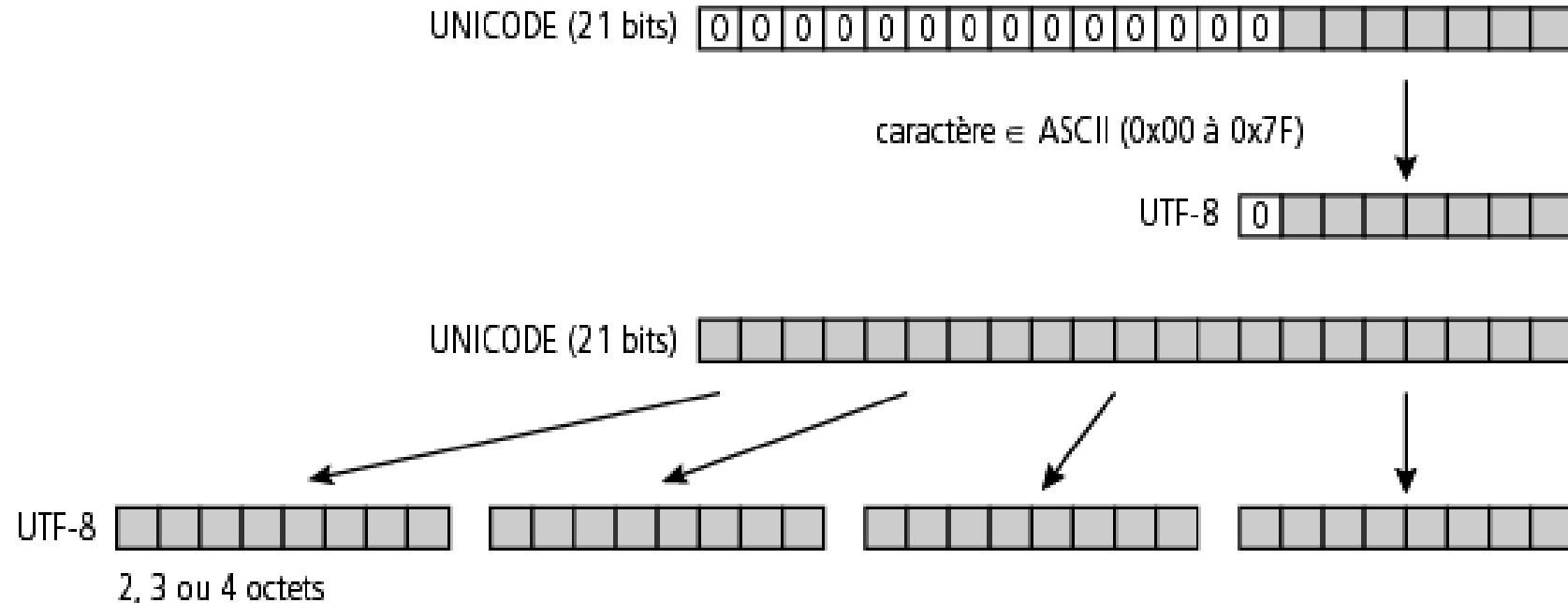
UTF-32 

Schéma d'encodage UTF-16

- Valeur tronquée à 16 bits pour les caractères ASCII sinon séquence de deux mots de 16 bits
- Optimisé pour les caractères du plan multilingue de base (BMP)
- Tableau de codage dans le TD

Schéma d'encodage UTF-8

- ❑ Facilité d'usage avec les systèmes basés sur l'ASCII
- ❑ Format variable
- ❑ Tableau de codage dans le TD



Schémas d'encodage

- Attention, ce n'est pas qu'un simple découpage

Table 3-5. UTF-16 Bit Distribution

Scalar Value	UTF-16
XXXXXXXXXXXXXXXXXX	XXXXXXXXXXXXXXXXXX
000uuuuuXXXXXXXXXXXXXXXXXX	110110wwwXXXXXX 110111XXXXXXXXXX

Note: www = uuuu - 1

Table 3-6. UTF-8 Bit Distribution

Scalar Value	First Byte	Second Byte	Third Byte	Fourth Byte
0000000 0xxxxxxx	0xxxxxxx			
00000yyy yyxxxxxx	110yyyyy	10xxxxxx		
zzzyyyyy yyxxxxxx	1110zzzz	10yyyyyy	10xxxxxx	
000uuuuu zzzzyyyy yyxxxxxx	11110uuu	10uuzzzz	10yyyyyy	10xxxxxx

Schémas de sérialisation

- Pour UTF-32 et UTF-16 :
 - forme *Little-Endian* (LE) : UTF-16LE et UTF-32LE
⇒ LSB en premier
 - forme *Big-Endian* (BE) : UTF-16BE et UTF-32BE
⇒ MSB en premier
- Version *endian-ness*
 - *Byte Order Mark* (BOM)

D'autres codes alphanumériques

- Code Page 850 de l'IBM PC (latin 1, Europe occidentale)
 - codes 128 à 255 en complément de l'ASCII
- Windows[®]-1252 → (latin 1, Europe occidentale)
- Le code Hollerith
 - origine : IBM
 - pour la carte perforée
- Le jeu de caractères vidéotex
- Etc.

	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
00	<u>NUL</u> 0000	<u>STX</u> 0001	<u>SOT</u> 0002	<u>ETX</u> 0003	<u>EOT</u> 0004	<u>ENQ</u> 0005	<u>ACK</u> 0006	<u>BEL</u> 0007	<u>BS</u> 0008	<u>HT</u> 0009	<u>LF</u> 000A	<u>VT</u> 000B	<u>FF</u> 000C	<u>CR</u> 000D	<u>SO</u> 000E	<u>SI</u> 000F
10	<u>DLE</u> 0010	<u>DC1</u> 0011	<u>DC2</u> 0012	<u>DC3</u> 0013	<u>DC4</u> 0014	<u>NAK</u> 0015	<u>SYN</u> 0016	<u>ETB</u> 0017	<u>CAN</u> 0018	<u>EM</u> 0019	<u>SUB</u> 001A	<u>ESC</u> 001B	<u>FS</u> 001C	<u>GS</u> 001D	<u>RS</u> 001E	<u>US</u> 001F
20	<u>SP</u> 0020	<u>!</u> 0021	<u>"</u> 0022	<u>#</u> 0023	<u>\$</u> 0024	<u>%</u> 0025	<u>&</u> 0026	<u>'</u> 0027	<u>(</u> 0028	<u>)</u> 0029	<u>*</u> 002A	<u>+</u> 002B	<u>,</u> 002C	<u>-</u> 002D	<u>.</u> 002E	<u>/</u> 002F
30	<u>0</u> 0030	<u>1</u> 0031	<u>2</u> 0032	<u>3</u> 0033	<u>4</u> 0034	<u>5</u> 0035	<u>6</u> 0036	<u>7</u> 0037	<u>8</u> 0038	<u>9</u> 0039	<u>:</u> 003A	<u>;</u> 003B	<u><</u> 003C	<u>=</u> 003D	<u>></u> 003E	<u>?</u> 003F
40	<u>@</u> 0040	<u>A</u> 0041	<u>B</u> 0042	<u>C</u> 0043	<u>D</u> 0044	<u>E</u> 0045	<u>F</u> 0046	<u>G</u> 0047	<u>H</u> 0048	<u>I</u> 0049	<u>J</u> 004A	<u>K</u> 004B	<u>L</u> 004C	<u>M</u> 004D	<u>N</u> 004E	<u>O</u> 004F
50	<u>P</u> 0050	<u>Q</u> 0051	<u>R</u> 0052	<u>S</u> 0053	<u>T</u> 0054	<u>U</u> 0055	<u>V</u> 0056	<u>W</u> 0057	<u>X</u> 0058	<u>Y</u> 0059	<u>Z</u> 005A	<u>[</u> 005B	<u>\</u> 005C	<u>]</u> 005D	<u>^</u> 005E	<u>_</u> 005F
60	<u>`</u> 0060	<u>a</u> 0061	<u>b</u> 0062	<u>c</u> 0063	<u>d</u> 0064	<u>e</u> 0065	<u>f</u> 0066	<u>g</u> 0067	<u>h</u> 0068	<u>i</u> 0069	<u>j</u> 006A	<u>k</u> 006B	<u>l</u> 006C	<u>m</u> 006D	<u>n</u> 006E	<u>o</u> 006F
70	<u>p</u> 0070	<u>q</u> 0071	<u>r</u> 0072	<u>s</u> 0073	<u>t</u> 0074	<u>u</u> 0075	<u>v</u> 0076	<u>w</u> 0077	<u>x</u> 0078	<u>y</u> 0079	<u>z</u> 007A	<u>{</u> 007B	<u> </u> 007C	<u>}</u> 007D	<u>~</u> 007E	<u>DEL</u> 007F
80	<u>€</u> 20AC		<u>f</u> 201A	<u>"</u> 0192	<u>…</u> 201E	<u>†</u> 2026	<u>‡</u> 2020	<u>~</u> 02C6	<u>%</u> 2030	<u>Š</u> 0160	<u><</u> 2039	<u>œ</u> 0152			<u>ž</u> 017D	
90		<u>`</u> 2018	<u>´</u> 2019	<u>¨</u> 201C	<u>ˆ</u> 201D	<u>•</u> 2022	<u>—</u> 2013	<u>—</u> 2014	<u>~</u> 02DC	<u>™</u> 2122	<u>š</u> 0161	<u>></u> 203A	<u>oe</u> 0153		<u>ž</u> 017E	<u>Ÿ</u> 0178
A0	<u>nbsp</u> 00A0	<u>ı</u> 00A1	<u>ç</u> 00A2	<u>£</u> 00A3	<u>¤</u> 00A4	<u>¥</u> 00A5	<u>¦</u> 00A6	<u>§</u> 00A7	<u>¨</u> 00A8	<u>©</u> 00A9	<u>ª</u> 00AA	<u>«</u> 00AB	<u>¬</u> 00AC	<u>–</u> 00AD	<u>®</u> 00AE	<u>¯</u> 00AF
B0	<u>°</u> 00B0	<u>±</u> 00B1	<u>²</u> 00B2	<u>³</u> 00B3	<u>´</u> 00B4	<u>µ</u> 00B5	<u>¶</u> 00B6	<u>·</u> 00B7	<u>¸</u> 00B8	<u>¹</u> 00B9	<u>º</u> 00BA	<u>»</u> 00BB	<u>¼</u> 00BC	<u>½</u> 00BD	<u>¾</u> 00BE	<u>¿</u> 00BF
C0	<u>À</u> 00C0	<u>Á</u> 00C1	<u>Â</u> 00C2	<u>Ã</u> 00C3	<u>Ä</u> 00C4	<u>Å</u> 00C5	<u>Æ</u> 00C6	<u>Ç</u> 00C7	<u>È</u> 00C8	<u>É</u> 00C9	<u>Ê</u> 00CA	<u>Ë</u> 00CB	<u>Ì</u> 00CC	<u>Í</u> 00CD	<u>Î</u> 00CE	<u>Ï</u> 00CF
D0	<u>Ð</u> 00D0	<u>Ñ</u> 00D1	<u>Ò</u> 00D2	<u>Ó</u> 00D3	<u>Ô</u> 00D4	<u>Õ</u> 00D5	<u>Ö</u> 00D6	<u>×</u> 00D7	<u>Ø</u> 00D8	<u>Ù</u> 00D9	<u>Ú</u> 00DA	<u>Û</u> 00DB	<u>Ü</u> 00DC	<u>Ý</u> 00DD	<u>Þ</u> 00DE	<u>ß</u> 00DF
E0	<u>à</u> 00E0	<u>á</u> 00E1	<u>â</u> 00E2	<u>ã</u> 00E3	<u>ä</u> 00E4	<u>å</u> 00E5	<u>æ</u> 00E6	<u>ç</u> 00E7	<u>è</u> 00E8	<u>é</u> 00E9	<u>ê</u> 00EA	<u>ë</u> 00EB	<u>ì</u> 00EC	<u>í</u> 00ED	<u>î</u> 00EE	<u>ï</u> 00EF
F0	<u>ð</u> 00F0	<u>ñ</u> 00F1	<u>ò</u> 00F2	<u>ó</u> 00F3	<u>ô</u> 00F4	<u>õ</u> 00F5	<u>ö</u> 00F6	<u>÷</u> 00F7	<u>ø</u> 00F8	<u>ù</u> 00F9	<u>ú</u> 00FA	<u>û</u> 00FB	<u>ü</u> 00FC	<u>ý</u> 00FD	<u>þ</u> 00FE	<u>ÿ</u> 00FF

Résumé des codes courants

n = 7 bits

n = 8 bits

n = 16 bits

n = 21 bits

n = 32 bits

ASCII
(ASA X3.4-1963)

ISO 646

ISO 8859-x

Windows® CP1252

EBCDIC

UNICODE ver. X.X

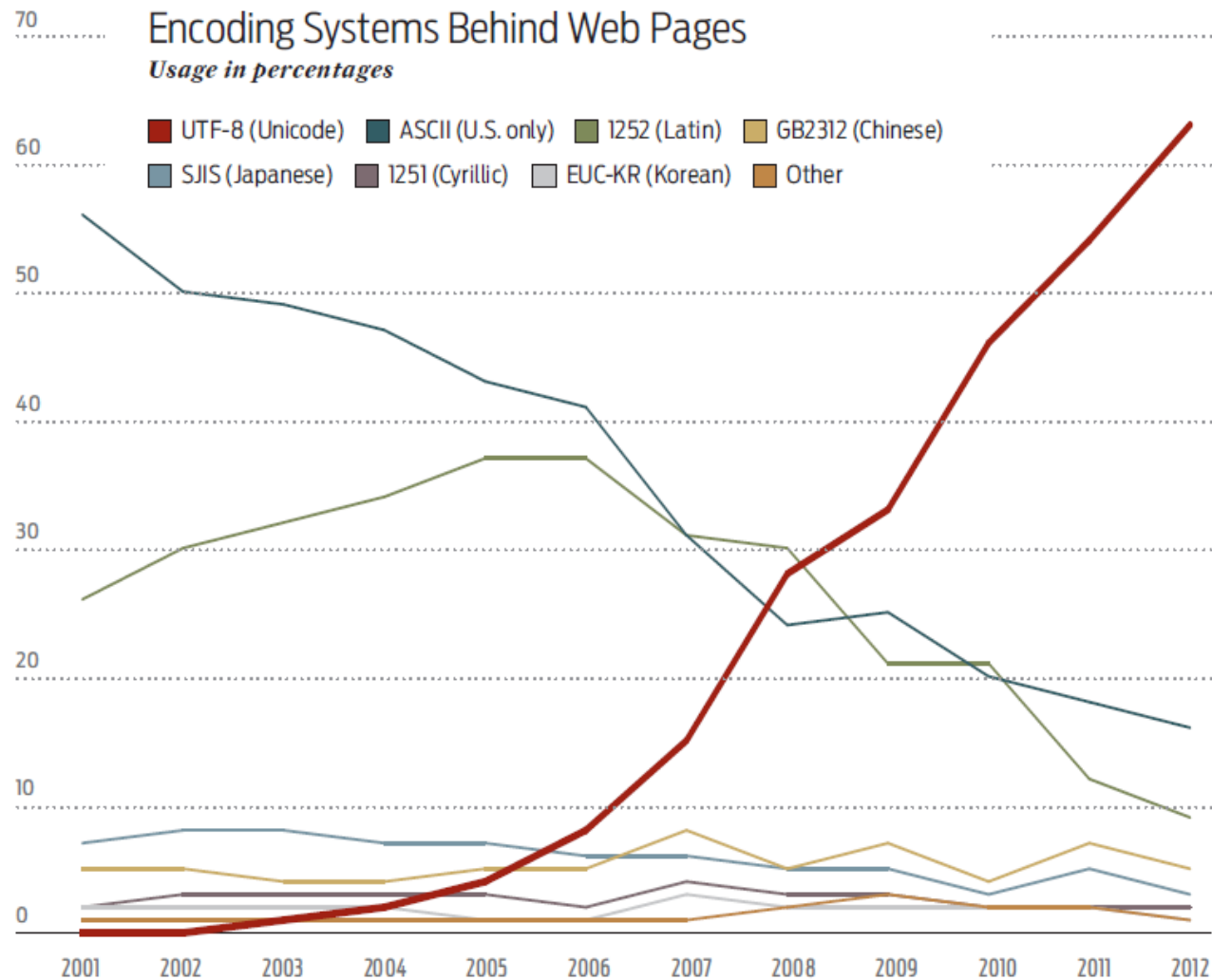
UNICODE ver. 4.0 et 5.0

UTF-8
UTF-16
UTF-32

ISO/IEC 10646:2009

UTF-8
UTF-16
UTF-32

And the winner is ...



[King 12]

Conclusion : le chemin des données (*datapath*)

