

# Architecture des ordinateurs

## 4 - Représentation de l'information en machine Définitions de base Les codes alphanumériques

Philippe Darche  
IUT Paris Descartes

## Le problème

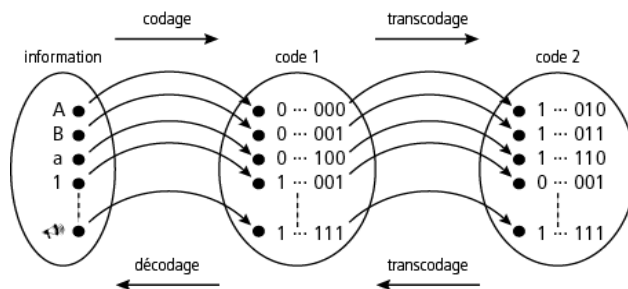
- Pour écrire l'anglais, 26 lettres, 10 chiffres et quelques caractères de ponctuation sont suffisants
- Pour les langue européennes comme le français, il y a en plus d'autres caractères ou variations sur les caractères
- Il existe d'autres alphabets ou certains pays utilisent des idéogrammes
- **Comment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

Philippe Darche

2

IUT Paris Descartes

## La notion de code



Philippe Darche

3

IUT Paris Descartes

## 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.

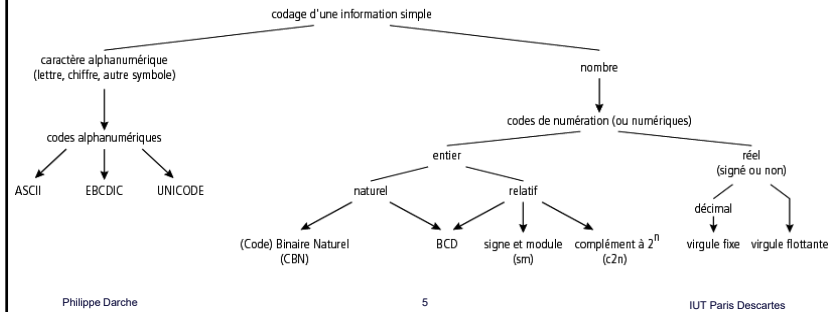
Philippe Darche

4

IUT Paris Descartes

## 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.

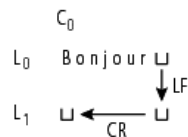
Philippe Darche

6

IUT Paris Descartes

## Utilisation d'un code de contrôle

- Pour gérer un périphérique
  - écran + clavier = un terminal
  - imprimante
- Exemple : le retour à la ligne ou inversement
  - retour chariot puis saut de ligne



Philippe Darche

7

IUT Paris Descartes

## Remarque importante

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

Philippe Darche

8

IUT Paris Descartes

## 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

The table shows the original 7-bit character set. The first 16 rows represent control characters, with bit patterns (b7 to b1) and names like NULL, SOM, EOA, EDM, EOT, WRU, RU, BELL, FEG, LF, VTAB, FF, CR, SO, SI. The remaining rows show printable characters from space to tilde (~). A red box highlights the control characters from NULL to SI, and another red box highlights the bit patterns b7 to b1.

contrôle

## Évolution de la table originale ASCII

- Prise en compte des minuscules

The table shows the evolution of the ASCII table. The first 16 rows represent control characters, with bit patterns (b7 to b1) and names like NULL, SOH, STX, ETX, EOT, ENQ, ACK, BEL, BS, HT, LF, VT, FF, CR, SO, SI. The remaining rows show printable characters from space to tilde (~). A red box highlights the control characters from NULL to SI, and another red box highlights the bit patterns b7 to b1. The table is titled "Table E.1. AMERICAN STANDARD CODE FOR INFORMATION INTERCHANGE".

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

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

## L'EBCDIC

- Origine : IBM
- format natif n = 8 bits

hex	0-	1-	2-	3-	4-	5-	6-	7-	8-	9-	A-	B-	C-	D-	E-	F-
00	NUL	DEL			ctrl	&	-								0	
01	SOH	DC1				a	j				A	J			1	
02	STX	DC2	FS	SIN		b	k	s			B	K	S	2		
03	ETX	DC3				c	l	t			C	L	T	3		
04						d	m	u			D	M	U	4		
05	HT	LF				e	n	v			E	N	V	5		
06		BS	ETB			f	o	w			F	O	W	6		
07			ESC	EDT			g	p	x		G	P	X	7		
08							h	q	y		H	Q	Y	8		
09							i	r	z		I	R	Z	9		
0A																
0B																
0C																
0D																
0E																
0F																
10																
11																
12																
13																
14																
15																

contrôle

## Limitation des codes alphanumériques classiques

- Format limité à n = 7 ou 8 bits
  - nombre de caractères codables limité (2<sup>n</sup>)
    - ⇒ 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

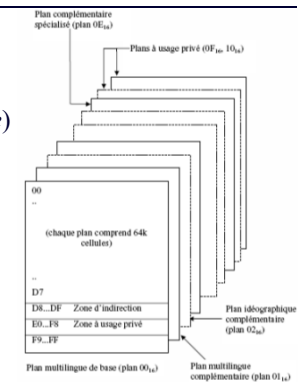
- Caractère : forme abstraite d'un composant atomique d'un langage
  - lettre, signe de ponctuation, etc.
  - exemple : U+0041 *latin capital letter a*
- Glyphe : forme graphique d'un caractère
  - police : ensemble de glyphes
  - de l'exemple précédent : A, ▲, A, etc.

## 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*)

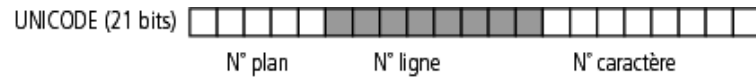
## Plans et zones de codage

- 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)



## 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



## D'autres codes alphanumériques

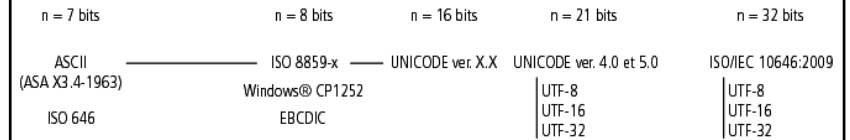
- Code Page 850 de l'IBM PC (latin 1, Europe occidentale)
  - codes 128 à 255 en complément de l'ASCII
- Windows®-1252 (latin 1, Europe occidentale) →
- Le code Hollerith
  - origine : IBM
  - pour la carte perforée
- Le jeu de caractères vidéotex
- Etc.

Philippe Darche

29

IUT Paris Descartes

## Résumé des codes courants

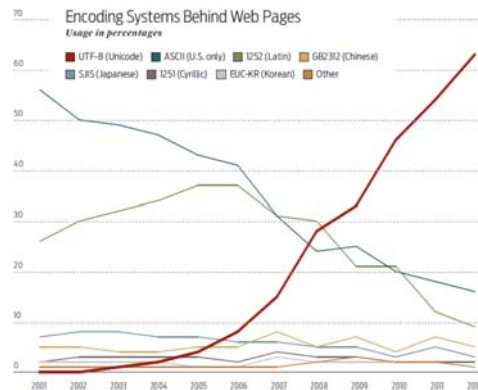


Philippe Darche

30

IUT Paris Descartes

## And the winner is ...

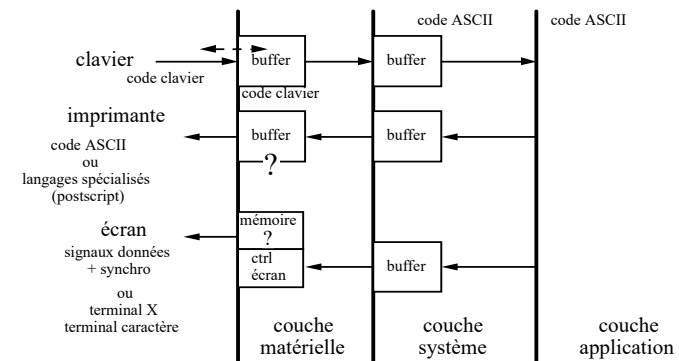


[King 12]

Philippe Darche

IUT Paris Descartes

## Conclusion : le chemin des données (datapath)



Philippe Darche

32

IUT Paris Descartes