

# Charsets & encodings

Jacquelin Charbonnel

Journées Mathrice de Grenoble, mars 2024

---

# Micmac & embrouillamini

Les données manipulées par un ordinateur sont représentées par une suite d'octets. Les textes ne font pas exception à cette règle : ils sont stockés eux aussi sous forme numérique. Quelles transformations sont nécessaires pour passer du texte à sa représentation numérique, et inversement ?

Les données manipulées par un ordinateur sont représentées par une suite d'octets. Les textes ne font pas exception à cette règle : ils sont stockés eux aussi sous forme numérique. Quelles transformations sont nécessaires pour passer du texte à sa représentation numérique, et inversement ?

# Micmac & embrouillamini

configurations spécifiques et de la méthodologie ENERGY STAR®, cela suppose que le produit est utilisé dans le pays et/ou la région du compte du client.

- Les valeurs de transport sont basées sur le pays et/ou la région du compte du client.

Les estimations de l'empreinte carbone sont basées sur l'ACV, et non sur les données d'empreinte réelle fournies par le fabricant du composant. Les estimations de l'empreinte carbone du produit peuvent différer lorsque les composants du produit et l'adresse de livraison sont mises à jour par le client. Toutes les

---

**poste de travail 2022 standard**

estimations d'émission de carbone sont incertaines, ce qui résulte en grande partie des limitations et de la qualité des données. Pour contrebalancer cette incertitude, HP a développé des outils spécifiques qui utilisent une combinaison de processus et de données sur les produits HP, ainsi que des données d'évaluation du cycle de vie de haute qualité. HP s'efforce de fournir les résultats les plus précis en matière d'impact environnemental, mais l'incertitude ne pourra jamais vraiment disparaître. Les résultats doivent être considérés en conséquence.

Page 4 of 4  
DEVIS CLIENT 4072614  
Les renseignements contenus dans ce devis sont confidentiels et restent la propriété exclusive de HP. Ne pas diffuser.

# Données textuelles

- Définitions

- texte composé de signes typographiques (les *caractères*) issus d'un alphabet (le *charset*)
- sur support numérique, chaque caractère d'un charset est représenté par un code (un entier)
- la correspondance *caractère* ↔ *code* pour un charset donné est un codage (encodage, *encoding*)

- Historique

- Les années 60 aux US
  - l'alphabet de la langue anglaise comporte peu de lettres (26)
  - le nombre de caractères pour rédiger un texte (min, MAJ, chiffres, signes) ne dépasse pas 90
  - 7 bits sont suffisants pour les coder
  - reste même de la place pour les caractères de contrôle
  - l'ASCII : un charset de 128 caractères

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

# Codage ASCII

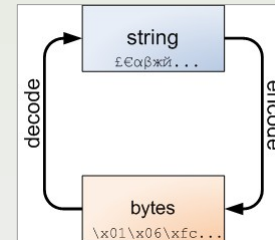
- chaque caractère est codé par son indice dans la table

```
$ cat <<EOT | xxd -g1
Lorem ipsum dolor sit amet, consectetur adipiscing elit,
sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.
EOT
00000000: 4c 6f 72 65 6d 20 69 70 73 75 6d 20 64 6f 6c 6f  Lorem ipsum dolo
00000010: 72 20 73 69 74 20 61 6d 65 74 2c 20 63 6f 6e 73  r sit amet, cons
00000020: 65 63 74 65 74 75 72 20 61 64 69 70 69 73 63 69  ectetur adipisci
00000030: 6e 67 20 65 6c 69 74 2c 0a 73 65 64 20 64 6f 20  ng elit,.sed do
00000040: 65 69 75 73 6d 6f 64 20 74 65 6d 70 61 72 20 69  eiusmod tempor i
00000050: 6e 63 69 64 69 64 75 6e 74 20 75 74 20 6c 61 62  ncididunt ut lab
00000060: 6f 72 65 20 65 74 20 64 6f 6c 6f 72 65 20 6d 61  ore et dolore ma
00000070: 67 6e 61 20 61 6c 69 71 75 61 2e 0a                gna aliqua..
```

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

code compris entre 00 et 7F

- opérations de codage/décodage à chaque lecture/écriture



# Années 80 revendicatives

- et les caractères nationaux ?
- idée simple : utiliser toute la tables (en utilisant le 8ème bit)

```
** iso8859-6 **
 0 1 2 3 4 5 6 7 8 9 A B C D E F
0:  @ P p
1:  ! 1 A Q a
2:  " 2 B R r
3:  # 3 C S s
4:  $ 4 D T t
5:  % 5 E U u
6:  & 6 F V v
7:  ' 7 G W w
8:  ( 8 H X x
9:  ) 9 I Y y
A:  * : J Z z
B:  + ; K [ k
C:  , < L \ l
D:  . = M ] m
E:  / ? N ^ n
F:  _ o
```

```
$ cat <<EOT | iconv -f utf8 -t latin1 | xxd -g1
Lorem ipsum dolor sit âmèt, consectetur adipisçing élit,
sêd do eiusmod tempor incididünt ut labore et dolore magna aliqua.
EOT
00000000: 4c 6f 72 65 6d 20 69 70 73 75 6d 20 64 6f 6c 6f  Lorem ipsum dolo
00000010: 72 20 73 69 74 20 e2 6d e8 74 2c 20 63 6f 6e 73  r sit .m.t, cons
00000020: 65 63 74 65 74 75 72 20 61 64 09 70 69 73 e7 69  ectetur adipis.i
00000030: 6e 67 20 e9 6c 69 74 2c 0a 73 ea 64 20 64 6f 20  ng .lit,.s.d do
00000040: 65 69 75 73 6d 6f 64 20 74 65 6d 70 6f 72 20 69  eiusmod tempor i
00000050: 6e 63 69 64 69 64 fc 6e 74 20 75 74 20 6c 61 62  ncidid.nt ut lab
00000060: 6f 72 65 20 65 74 20 64 6f 6c 6f 72 65 20 6d 61  ore et dolore ma
00000070: 67 6e 61 20 61 6c 69 71 75 61 2e 0a  gna aliqua..
```

code compris entre 00 et FF

# Charsets nationaux

## ISO8859-5 (cyrillic)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0:				@	P	Q	R	a	p	q	r	s	t	u	v	w
1:		!	1	A	Q	a	b	c	d	e	f	g	h	i	j	
2:		"	2	B	R	b	c	d	e	f	g	h	i	j	k	
3:		#	3	C	S	c	d	e	f	g	h	i	j	k	l	
4:		\$	4	D	T	d	e	f	g	h	i	j	k	l	m	
5:		%	5	E	U	e	f	g	h	i	j	k	l	m	n	
6:		&	6	F	V	f	g	h	i	j	k	l	m	n	o	
7:		'	7	G	W	g	h	i	j	k	l	m	n	o	p	
8:		(	8	H	X	h	i	j	k	l	m	n	o	p	q	
9:		)	9	I	Y	i	j	k	l	m	n	o	p	q	r	
A:		*	:	J	Z	j	k	l	m	n	o	p	q	r	s	
B:		+	;	K	[	k	l	m	n	o	p	q	r	s	t	
C:		,	<	L	\	l	m	n	o	p	q	r	s	t	u	
D:		-	=	M	]	m	n	o	p	q	r	s	t	u	v	
E:		.	>	N	^	n	o	p	q	r	s	t	u	v	w	
F:		/	?	O	_	o	p	q	r	s	t	u	v	w	x	

## iso8859-1 (latin1)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0:				@	P	Q	a	p	q	r	s	t	u	v	w	x
1:		!	1	A	Q	a	b	c	d	e	f	g	h	i	j	
2:		"	2	B	R	b	c	d	e	f	g	h	i	j	k	
3:		#	3	C	S	c	d	e	f	g	h	i	j	k	l	
4:		\$	4	D	T	d	e	f	g	h	i	j	k	l	m	
5:		%	5	E	U	e	f	g	h	i	j	k	l	m	n	
6:		&	6	F	V	f	g	h	i	j	k	l	m	n	o	
7:		'	7	G	W	g	h	i	j	k	l	m	n	o	p	
8:		(	8	H	X	h	i	j	k	l	m	n	o	p	q	
9:		)	9	I	Y	i	j	k	l	m	n	o	p	q	r	
A:		*	:	J	Z	j	k	l	m	n	o	p	q	r	s	
B:		+	;	K	[	k	l	m	n	o	p	q	r	s	t	
C:		,	<	L	\	l	m	n	o	p	q	r	s	t	u	
D:		-	=	M	]	m	n	o	p	q	r	s	t	u	v	
E:		.	>	N	^	n	o	p	q	r	s	t	u	v	w	
F:		/	?	O	_	o	p	q	r	s	t	u	v	w	x	

## Iso8859-7 (grec)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0:				@	P	Q	a	p	q	r	s	t	u	v	w	x
1:		!	1	A	Q	a	b	c	d	e	f	g	h	i	j	
2:		"	2	B	R	b	c	d	e	f	g	h	i	j	k	
3:		#	3	C	S	c	d	e	f	g	h	i	j	k	l	
4:		\$	4	D	T	d	e	f	g	h	i	j	k	l	m	
5:		%	5	E	U	e	f	g	h	i	j	k	l	m	n	
6:		&	6	F	V	f	g	h	i	j	k	l	m	n	o	
7:		'	7	G	W	g	h	i	j	k	l	m	n	o	p	
8:		(	8	H	X	h	i	j	k	l	m	n	o	p	q	
9:		)	9	I	Y	i	j	k	l	m	n	o	p	q	r	
A:		*	:	J	Z	j	k	l	m	n	o	p	q	r	s	
B:		+	;	K	[	k	l	m	n	o	p	q	r	s	t	
C:		,	<	L	\	l	m	n	o	p	q	r	s	t	u	
D:		-	=	M	]	m	n	o	p	q	r	s	t	u	v	
E:		.	>	N	^	n	o	p	q	r	s	t	u	v	w	
F:		/	?	O	_	o	p	q	r	s	t	u	v	w	x	

# Charsets nationaux

## iso8859-1 (latin1)

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

## iso8859-15

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



# Années 80 : les marques aussi se distinguent

## Apple - MacRoman

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0:				0	@	P	`	p	Ä	é	†	∞	¿	—	±	
1:		!	"	1	A	Q	a	q	Å	ë	°	±	¡	"	·	Ò
2:				2	B	R	b	r	Ç	í	¢	±	¡	"	,	Ó
3:		#	\$	3	C	S	c	s	È	î	£	±	¡	"	,"	Ô
4:		%	%	4	D	T	d	t	É	ï	§	±	¡	"	%	Ù
5:		&	&	5	E	U	e	u	Ê	ï	¥	±	¡	"	~	Ú
6:		'	'	6	F	V	f	v	Ë	ï	µ	±	¡	"	·	Û
7:		(	(	7	G	W	g	w	Ï	ï	¶	±	¡	"	·	Ü
8:		)	)	8	H	X	h	x	Ï	ï	·	±	¡	"	·	Ý
9:		*	*	9	I	Y	i	y	Ï	ï	©	±	¡	"	·	ÿ
A:		+	+		J	Z	j	z	Ï	ï	™	±	¡	"	·	
B:		,	,		K	[	k	{	Ï	ï	·	±	¡	"	·	
C:		-	-		L	\	l		Ï	ï	·	±	¡	"	·	
D:		.	.		M	]	m	}	Ï	ï	·	±	¡	"	·	
E:		/	/		N	^	n	~	Ï	ï	·	±	¡	"	·	
F:					O	_	o	~	Ï	ï	·	±	¡	"	·	

## Microsoft - windows\_1250

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0:				0	@	P	`	p	€	◊	˘	°	Á	Đ	ř	ď
1:		!	"	1	A	Q	a	q	◊	◊	˘	±	Á	Đ	ř	ď
2:				2	B	R	b	r	◊	◊	˘	±	Á	Đ	ř	ď
3:		#	\$	3	C	S	c	s	◊	◊	˘	±	Á	Đ	ř	ď
4:		%	%	4	D	T	d	t	◊	◊	˘	±	Á	Đ	ř	ď
5:		&	&	5	E	U	e	u	◊	◊	˘	±	Á	Đ	ř	ď
6:		'	'	6	F	V	f	v	◊	◊	˘	±	Á	Đ	ř	ď
7:		(	(	7	G	W	g	w	◊	◊	˘	±	Á	Đ	ř	ď
8:		)	)	8	H	X	h	x	◊	◊	˘	±	Á	Đ	ř	ď
9:		*	*	9	I	Y	i	y	◊	◊	˘	±	Á	Đ	ř	ď
A:		+	+		J	Z	j	z	◊	◊	˘	±	Á	Đ	ř	ď
B:		,	,		K	[	k	{	◊	◊	˘	±	Á	Đ	ř	ď
C:		-	-		L	\	l		◊	◊	˘	±	Á	Đ	ř	ď
D:		.	.		M	]	m	}	◊	◊	˘	±	Á	Đ	ř	ď
E:		/	/		N	^	n	~	◊	◊	˘	±	Á	Đ	ř	ď
F:					O	_	o	~	◊	◊	˘	±	Á	Đ	ř	ď

- tous ces charsets ont 256 caractères
- les codes tiennent sur 1 octet
- tous ces charsets englobent l'ASCII
- mais ils sont mutuellement exclusifs
  - par exemple, impossible d'écrire :

Le dicton russe « *Азбука -- к мудрости ступенька.* » signifie « *Nous devons d'abord apprendre à marcher avant d'apprendre à courir.* »

# Fin des années 80 : ça part en vrille

- arrivée de charsets de plus de 256 caractères :
  - EUC-JP, pour le japonais (caractères codés sur 2 octets)
  - Big5, pour le chinois traditionnel à Taïwan (codes sur 1 ou 2 octets suivant les caractères)
- mais changer la taille des codes a un impact colossal
  - sur le matériel et le logiciel :
    - sur les langages de programmation, les compilateurs, les bibliothèques, les systèmes d'exploitation
  - **rupture technologique**
  - c'est pourquoi, les charsets codés sur 1 octet ont été exploités aussi longtemps que possible
- quitte à changer, autant repartir à zéro et voir loin, pour :
  - uniformiser le codage des signes de toutes les cultures (actuelles et passées), des symboles, smileys, etc.,
  - pouvoir utiliser tous ces signes au sein d'un même texte
  - obtenir des documents utilisables à l'identique sur tout système informatique

# Années 90 : Unicode remet à plat

- un signe typographique est défini comme la plus petite composante d'un alphabet ayant une valeur sémantique :
  - l'accent circonflexe est un signe à part entière
  - la cédille en est un autre
  - Unicode référence plus d'un million de signes
- plus que des alphabets, Unicode gère des écritures
  - chaque caractère Unicode possède des attributs :
    - un nom : solidus – Greek capital letter delta – left curly bracket
    - un **code point**
      - identifiant entier (sur 4 octets) canonique et unique
      - sans vocation à coder numériquement le caractère
      - noté U+xxxx, U+xxxxx ou U+xxxxxx
        - exemple : U+0000, U+FFFF, U+10FFFF
    - et des méta-données :
      - est-ce une majuscule ? une minuscule ? une ponctuation ? un chiffre (et si oui, sa valeur numérique) ?
      - se combine-t-il à d'autres caractères (accent, cédille, tilde, tréma, etc.) ?
      - son orientation (droite → gauche, gauche → droite, haut → bas, etc.)

# Unicode attributs

Par exemple via Python

```
>>> import unicodedata

>>> unicodedata.name('/')
'SOLIDUS'

>>> unicodedata.name('{')
'LEFT CURLY BRACKET'

>>> unicodedata.lookup('solidus')
'/'

>>> unicodedata.category('A')
'Lu' # Letter uppercase

>>> unicodedata.category('a')
'Ll' # Letter lowercase

>>> unicodedata.category('€')
'Sc' # Symbol currency

>>> c = "\u0664"
>>> c
'٤'

>>> unicodedata.name(c)
'ARABIC-INDIC DIGIT FOUR'

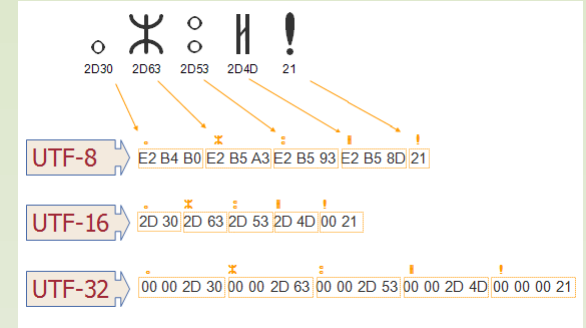
>>> unicodedata.category(c)
'Nd' # Number decimal

>>> unicodedata.numeric(c)
4.0
```

# Les codages d'Unicode

Unicode propose plusieurs encodage

- UTF-32
  - chaque caractère est codé directement avec son code point (sur 4 octets)
  - cette simplicité est son principal avantage : le codage/décodage est simple
  - par contre, il est particulièrement gourmand en taille
  - évidemment incompatible avec l'ASCII
  - peu utilisé
- UTF-16
  - la plupart des caractères (notamment ceux des langues modernes) est représentée sur 2 octets, les autres le sont sur 4 octets.
  - utilisé dans l'environnement Windows
  - évidemment incompatible avec l'ASCII
- UTF-8
  - encodage le plus utilisé en occident
  - son codage de taille variable (entre 1 et 4 octets) lui permet d'être parcimonieux en taille pour les alphabets latin, au détriment de la vitesse de codage/décodage.
  - compatible avec l'ASCII
  - inconvénient : la longueur variable de son codage complexifie le codage/décodage



# UTF-8

- principe du codage :
  - si le code point est inférieur à 128, le code est égal au code point,
  - sinon, le code est une séquence de 2, 3 ou 4 octets, où chacun est compris entre 128 et 255.
- détail du codage :
  - les caractères ASCII sont codés sur un octet dont le bit de poids fort est 0 : 0xxxxxxx
  - les caractères non-ASCII utilisent au moins deux octets :
    - le premier octet commence par autant de 1 que d'octets dans la séquence, suivis par un 0.
    - les autres octets commencent par 10.
- ce qui donne :
  - code point ayant entre 8 à 11 bits : 110xxxxx 10xxxxxx
  - code point ayant entre 12 à 16 bits : 1110xxxx 10xxxxxx 10xxxxxx
  - code point ayant entre 17 à 21 bits : 11110xxx 10xxxxxx 10xxxxxx 10xxxxxx

```
$ cat <<EOT | xxd -g1
Lorem ipsum dolor sit âmet, consectetur adipiscing elit,
sêd do eiusmod tempor incididunt ut labore et dolore magna aliqua.
EOT
00000000: 4c 6f 72 65 6d 20 69 70 73 75 6d 20 64 6f 6c 6f  Lorem ipsum dolo
00000010: 72 20 73 69 74 20 c3 a2 6d c3 a8 74 2c 20 63 6f  r sit ..m..t, co
00000020: 6e 73 65 63 74 65 74 75 72 20 61 64 69 70 69 73  nsectetur adipis
00000030: c3 a7 69 6e 67 20 c3 a9 6c 69 74 2c 0a 73 c3 aa  ..ing ..lit,..s..
00000040: 64 20 64 6f 20 65 69 75 73 6d 6f 64 20 74 65 6d  d do eiusmod tem
00000050: 70 6f 72 20 69 6e 63 69 64 69 64 c3 bc 6e 74 20  por incidid..nt
00000060: 75 74 20 6c 61 62 6f 72 65 20 65 74 20 64 6f 6c  ut labore et dol
00000070: 6f 72 65 20 6d 61 67 6e 61 20 61 6c 69 71 75 61  ore magna aliqua
00000080: 2e 0a
```

4c est inférieur à 128  
donc le code est sur 1 octet

c3==11000011  
donc le code est sur 2 octets

# Écueils & curiosités

Faux-amis

```
>>> import unicodedata

>>> a, b = "\u0061", "\u0430"
>>> a, b, a==b
('a', 'a', False)

>>> [ unicodedata.name(x) for x in [a,b] ]
['LATIN SMALL LETTER A', 'CYRILLIC SMALL LETTER A']
```

Python

```
>>> import unicodedata

>>> a, b = "\u0050", "\u03a1"
>>> a, b, a==b
('P', 'P', False)

>>> [ unicodedata.name(x) for x in [a,b] ]
['LATIN CAPITAL LETTER P', 'GREEK CAPITAL LETTER RHO']
```

Python

```
>>> import unicodedata

>>> s = "A\u216c\u2160B\u0391\u0392\u0391.COM"
>>> s
'ALIBABA.COM'

>>> [ unicodedata.name(c) for c in s ]
['LATIN CAPITAL LETTER A',
 'ROMAN NUMERAL FIFTY',
 'ROMAN NUMERAL ONE',
 'LATIN CAPITAL LETTER B',
 'LATIN CAPITAL LETTER A',
 'GREEK CAPITAL LETTER BETA',
 'GREEK CAPITAL LETTER ALPHA', ... ]
```

Python

Unicode modélise des écritures

```
>>> import unicodedata

>>> a, b = "C", "Ç"
>>> c, d = "\u00B8", "\u0327"
>>> a, b, c, d
('C', 'Ç', ' ', ' ')

>>> [ unicodedata.name(x) for x in [a,b,c,d] ]
['LATIN CAPITAL LETTER C',
 'LATIN CAPITAL LETTER C WITH CEDILLA',
 'CEDILLA',
 'COMBINING CEDILLA']

>>> b, a+c, a+d
('Ç', 'C, ', 'Ç')

>>> b == a+d
False

>>> len(b), len(a+d)
(1, 2)

>>> unicodedata.name(a+d)
TypeError: name() argument 1 must be a unicode character, not str
```

Python

```
$ echo C - $'\u0327' - C$'\u0327' - C$'\u0327'\u0327'
C - , - Ç - Ç,

$ echo A$'\u0327'A - Live in Paris
AÇA - Live in Paris
```

bash