8-TTU

plications cun de vous doit être capable de faire et d'expliquer les exercices d'ap-Cette activité de groupe doit être comprise et réalisée en autonomie. Cha-

1. Encodage de textes

T WACHINES

la différencier de la lettre 'A'? Comment la lettre 'Z' est-elle stockée dans la mémoire de l'ordinateur? Comment

caractère et sa représentation en binaire. Pour répondre à ces questions, il a fallu créer une table de conversion entre un

Première table utilisée (début 1960) : **ASCII**

[DEF]	4L	127	-	45	S6	ż	3E	63	(ROTARA932 TINU)	JE.	τε
~	37	126	~	35	76	<	3E	79	(RECORD SEPARATOR)	31	30
- {	ΔL	ISS	[2D	89	=	3D	19	[ROTARA932 9UOR3]	ID	57
	27	124	١.	2C	76	>	3C	09	(FILE SEPARATOR)	JC	28
}	87	123	1	88	16		38	69	[ESCAPE]	IB	77
Z	A۲	122	Z	AZ	06		ΑE	89	(SUBSTITUTE)	ΥT	97
٨	64	121	٨	65	68	6	68	LS	(END OF MEDIUM)	61	52
X	87	120	Х	85	88	8	38	99	[CANCEL]	18	74
w	LL	611	W	LS	78	L	75	22	[ENG OF TRANS. BLOCK]	۷.τ	23
٨	94	118	٨	99	98	9	98	75	(SLNCHBONONS IDFE)	91	22
n	SZ	LTT	n	22	82	9	32	23	[NEGATIVE ACKNOWLEDGE]	ST	7.7
4	tι	911	1	75	1/8	t	34	25	[DEVICE CONTROL 4]	ÞΤ	20
S	٤2	STT	S	23	83	3	33	TS	[DEVICE CONTROL 3]	ΙЗ	6T
	7.5	114	Я	25	28	2	32	20	[DEVICE CONTROL 2]	15	18
b	TΖ	113	b	τs	18	τ	3.7	61	[DEVICE CONTROL 1]	ττ	77
d	04	112	d	20	08	0	30	81/	[DATA LINK ESCAPE]	OT	91
0	49	TTT	0	∃₽	64	/	2F	Lt	(NI L±IHS)	3	ST
u	39	TTO	N	3Þ	87		SE	91	[TUO THIR]	3	Įτ
ш	Q9	60T	M	ďΦ	LL	-	SD	St	(CARRIAGE RETURN)	O	13
- 1	29	108	- 1	7C	94	,	SC	ヤヤ	[FORM FEED]	2	15
K	89	LOT	K	4B	SZ	+	28	43	[VERTICAL TAB]	8	TT
į	∀9	90T	ſ	٧Þ	7.4	*	AS	45	(CINE LEED)	A	TO
	69	TOS	- 1	67	73	(67	Ιt	[8AT JATNOSIROH]	6	6
ч	89	T04	H	84	7.5)	82	01⁄2	(BYCKSbyce)	8	8
6	L9	103	9	LÞ	τz		77	68	(1739)	L	L
ł	99	102	4	91	04	79	97	38	[ACKNOWLEDGE]	9	9
Э	59	TOT	3	St	69	%	SZ	7.5	[ENÓNIBA]	S	S
р	† 9	T00	D	ヤヤ	89	\$	77	98	[END OF TRANSMISSION]	Þ	Þ
0	٤9	66	2	43	L9	#	23	32	[TX3T 90 GN3]	3	3
q	79	86	8	45	99		22	45	[TX3T 90 TRAT2]	7	7
В	19	L6	A	ΙÞ	59	1.1	7.7	33	(START OF HEADING)	τ	T
	09	96	0	01⁄2	79	[SPACE]	20	32	[אחרד]	0	0
Char	хән	Decimal	Char	хән	Decimal	Char	хэн	Decimal	Сһаг	ХЭН	Decimal

Cas 4 : point de code de longueur supérieur à 16 bits :

- encodage sur 4 x 8 bits

- utilisation des préfixes : 11110 puis 10, 10 et 10

1111 Ounu 10uu zzzz 10yy yyyy 10xx xxxx	000 nuuu zzzz yyyy yyxx xxxx			
valeur en UTF-8	point de code en binaire			

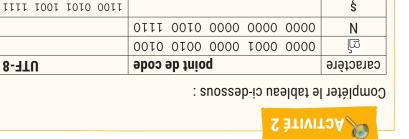
2. Exemples

stid 8×4 (4 sec)	0010 0001 0111 0000 1000 1000 0111	0001 0000 0011 1000 0100	
stid 8×8 (8 sso)	1110 0110 1001 0110 0111	1110 0001 1010 0110	X
stid 8×2 (2 seo)	1100 1110 1010 1001	1001 1010 1001	Ω
stid 8×f (f 2so)	1000 0010	1000 0010	Α
remarque	8-TTU ne valeur en UTF-8	point de code en binaire	

I octet) mais pas pour les caractères asiatiques (utilisation de 3 octets). Remarques UTF-8 est optimisé pour les caractères occidentaux (utilisation de

3. Exercice d'application

Ħ



0001 1101 0011 0001 1101

MACHINES



-**D**MACHINES

Quelle est la représentation en décimal du caractère Z? en binaire?

Quelle sont celles de la caractère z? Comment représenter le caractère é?

CORRECTION

Z est représenté par le nombre décimal $(90)_{10}$. Ce qui donne en binaire $(1011010)_2$ car $90 = 64 + 16 + 8 + 2 = 2^6 + 2^4 + 2^3 + 2^1$

z est représenté par le nombre décimal $(122)_{10}$. Ce qui donne en binaire $(1111010)_2$ car $122 = 64 + 32 + 16 + 8 + 2 = 2^6 + 2^5 + 2^4 + 2^3 + 2^1$

Impossible avec la table ASCII. C'est pour cela qu'à la fin des années 1960 ont été introduites d'autres tables (par exemple ISO-8859).

Conversions texte \leftrightarrow binaire

La norme Unicode (début 1990) permet une correspondance entre chaque caractère de l'humanité utilisé dans le monde (passé et présent) et un nombre entier appelé **point de code**.

Le point de code d'une taille maximale de 32 bits peut être représenté par des nombres binaires de tailles différentes :

- 1 à 4 paquets de taille 8 bits : UTF-8

- 1 à 2 paquets de taille 16 bits : UTF-16

- 1 paquet de taille 32 bits : UTF-32.

Dans la suite de ce document, nous allons étudier *uniquement* l'encodage en UTF-8 des caractères suivants :

caractère	origine	point de code	point de code binaire
		décimal	
Α	₩	65	0100 0001
Ω	lettre grecque	937	0011 1010 1001
文	idéogramme chinois	25 991	0110 0101 1000 0111
¥¥¥	ougaritique (Syrie) des années -1400 à -1100	66 436	0001 0000 0011 1000 0100

Méthode UTF-8 4 cas en fonction de la longueur du point de code :

Cas 1: point de code de longueur 7 bits ou moins :

- encodage sur 7 bits identique au point de code
- complété par des 0 à gauche pour atteindre une taille de 8 bits (padding)

point de code en binaire	valeur en UTF-8		
0xxx xxxx	0xxx xxxx		

Cas 2 : point de code de longueur 8 bits à 11 bits :

- encodage sur 2 × 8 bits

- utilisation des préfixes : 110 puis 10

point de code en binaire	valeur en UTF-8					
Оууу уухх хххх	110y yyyy 10xx xxxx					

Cas 3 : point de code de longueur 12 bits à 16 bits :

encodage sur 3 × 8 bits

- utilisation des préfixes : 1110 puis 10 et 10

point de code en binaire			valeur en UTF-8						
ZZZZ	уууу	уухх	xxxx	1110	ZZZZ	10уу	уууу	10xx	xxxx