



Représentation de l'information

**Spécialité numérique et sciences informatiques en classe de
première**

Lycée Rabelais
Saint Briec
2023-2024

Chapitre 1

Bases de numération

« Partons sur de bonnes bases. »

On note \mathbf{N} l'ensemble des *entiers naturels* : $\mathbf{N} = \{0; 1; 2; \dots\}$.

Nous avons l'habitude d'utiliser la base 10 pour représenter les entiers naturels, c'est-à-dire qu'on utilise 10 symboles, appelés *chiffres* pour les écrire : 0, 1, 2, ..., 9. Or il n'en a pas toujours été ainsi :

- au I^{er} millénaire av. J.-C., les Babyloniens utilisaient la base soixante pour mesurer le temps et les angles;
- durant le I^{er} millénaire, les Mayas et les Aztèques se servaient de la base vingt (et d'ailleurs en France, 80 se lit « quatre-vingts »);
- entre le VII^e et le XV^e siècle, les astronomes arabes utilisaient la base cent cinquante pour élaborer des tables permettant de trouver la position d'un astre dans le ciel à un moment donné.

De nos jours, en Informatique, on utilise beaucoup la base deux, dite *binaire* et la base seize, appelée *hexadécimale*.

L'objectif de ce chapitre est de donner les méthodes permettant d'écrire un entier naturel dans une base donnée, plus précisément dans les bases 2, 10 et 16. Nous verrons également comment passer facilement du binaire à l'hexadécimal et vice-versa.

1 Écriture binaire d'un entier naturel

1.1 Pourquoi le binaire ?



Pour simplifier, disons qu'au niveau le plus « bas » d'un ordinateur, se trouvent des (millions de) transistors qui jouent chacun un rôle d'interrupteur. De multiples points de l'ordinateur peuvent alors être soumis à une tension (état 1) ou non (état 0). En considérant 2 de ces points, on voit que l'état de ce système peut être 00, 01, 10 ou 11. Cela fait 4 possibilités et le binaire est né !

1.2 Comprendre l'écriture en base 2

Puisqu'il n'y a que deux chiffres en binaire, compter est simple mais nécessite rapidement plus de chiffres qu'en base 10 :

Écriture décimale	0	1	2	3	4	5	6	7	8	...
Écriture binaire	0	1	10	11	100	101	110	111	1000	...

Notation

On écrira $(11)_{10}$ pour insister sur le fait qu'on parle du nombre 11 *en base 10*, et $(11)_2$ pour dire que c'est une écriture binaire.

Lorsque ce n'est pas précisé cela veut dire que l'écriture est en base 10.

Ainsi $(11)_2 = 3$, et de même, $(111)_2 = 7$.

Tout entier naturel admet une unique écriture décimale (c'est-à-dire en base 10), il en va de même en binaire :

Propriété : écriture binaire d'un entier naturel

Tout entier naturel possède une unique écriture en base 2, dite *écriture binaire*. Plus précisément, soit $n \in \mathbf{N}$, alors il existe un unique entier $k \in \mathbf{N}$ et $k+1$ nombres a_i , uniques et valant 0 ou 1 et tels que

$$n = a_0 2^0 + a_1 2^1 + \dots + a_k 2^k$$

ce qui s'écrit aussi

$$n = \sum_{i=0}^k a_i 2^i$$

Exemple

Lorsqu'on regarde le tableau précédent, on voit que $6 = (110)_2$.
Cela s'interprète ainsi :

Chiffre binaire	1	1	0
Valeur	2^2	2^1	2^0

et on obtient $6 = 0 \times 2^0 + 1 \times 2^1 + 1 \times 2^2$.

Méthode 1 : passer de la base 2 à la base 10

Que vaut $(11101)_2$?

Chiffre binaire	1	1	1	0	1
Valeur	2^4	2^3	2^2	2^1	2^0

$$\begin{aligned}
 (11101)_2 &= 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 \\
 &= 16 + 8 + 4 + 1 \\
 &= 29
 \end{aligned}$$

Méthode 2 : passer de la base 10 à la base 2

Comprenons ce que veut dire une écriture décimale :

$$\begin{aligned}
 203 &= 200 + 3 \\
 &= 2 \times 10^2 + 0 \times 10^1 + 3 \times 10^0
 \end{aligned}$$

Faisons la même chose en base 2 :

$$\begin{aligned}
 203 &= 128 + 64 + 8 + 2 + 1 \\
 &= 2^7 + 2^6 + 2^3 + 2^1 + 2^0 \\
 &= 1 \times 2^7 + 1 \times 2^6 + 0 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 \\
 &= (11001011)_2
 \end{aligned}$$

Cette méthode est pratique quand l'entier est petit et que l'on connaît bien les premières puissances de deux.

Quand ce n'est pas le cas, une autre méthode (la méthode 3) peut être employée.

Évidemment, PYTHON connaît le binaire et travaille avec des valeurs entières de type `int` (*integer* veut dire « entier » en Anglais, pour plus de précisions, voir le chapitre « Valeurs et types », section `int`).

Voici comment écrire un `int` en binaire et comment obtenir l'écriture binaire d'un `int`.

Python

```
>>> 0b11001011 # faire précéder le nombre de 0b
203
>>> bin(29)
'0b11101'
```

1.3 Un algorithme pour déterminer l'écriture binaire d'un entier naturel**Méthode 3 : les divisions successives**

Voici comment on trouve les chiffres de l'écriture *décimale* de 203 :

On divise 203 par 10, cela fait 20, il reste 3, c'est le chiffre des unités.
 On recommence avec 20, on le divise par 10, cela fait 2, reste 0, chiffre des dizaines.
 On continue, on divise 2 par 10, cela fait 0, reste 2, chiffre des centaines.
 Puisqu'on a trouvé un quotient de 0, on s'arrête.
 On peut écrire cela simplement :

$$\begin{array}{r|l} 203 & 10 \\ \hline 3 & 20 \\ \hline & 0 \\ & 2 \\ & 10 \\ \hline & 0 \end{array}$$

Voici maintenant comment on trouve son écriture binaire. On procède comme en base 10 mais en divisant par 2 :

$$\begin{array}{r|l} 203 & 2 \\ \hline 1 & 101 \\ \hline & 1 \\ & 50 \\ \hline & 0 \\ & 25 \\ \hline & 1 \\ & 12 \\ \hline & 0 \\ & 6 \\ \hline & 0 \\ & 3 \\ \hline & 1 \\ & 1 \\ \hline & 1 \\ & 1 \\ \hline & 0 \end{array}$$

On a donc successivement établi :

$$\begin{aligned} 203 &= 101 \times 2 + 1 \\ &= (50 \times 2 + 1) \times 2 + 1 \\ &= ((25 \times 2 + 0) \times 2 + 1) \times 2 + 1 \\ &= (((12 \times 2 + 1) \times 2 + 0) \times 2 + 1) \times 2 + 1 \end{aligned}$$

$$\begin{aligned}
&= (((((6 \times 2 + 0) \times 2 + 1) \times 2 + 0) \times 2 + 1) \times 2 + 1) \times 2 + 1 \\
&= (((((3 \times 2 + 0) \times 2 + 0) \times 2 + 1) \times 2 + 0) \times 2 + 1) \times 2 + 1 \\
&= ((((((1 \times 2 + 1) \times 2 + 0) \times 2 + 0) \times 2 + 1) \times 2 + 0) \times 2 + 1) \times 2 + 1 \\
&= 1 \times 2^7 + 1 \times 2^6 + 0 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 \\
&= (11001011)_2
\end{aligned}$$

Cette succession d'égalités n'est (heureusement) pas à écrire à chaque fois.

Les trois méthodes précédentes se programment facilement et la dernière est de loin la plus courte à écrire.

1.4 Vocabulaire



Un chiffre décimal peut être 0, 1, 2, 3, 4, 5, 6, 7, 8 ou 9. Un chiffre binaire peut être seulement 0 ou 1. En Anglais, *chiffre binaire* se traduit par *binary digit*, que l'on abrège en *bit*. On garde cette dénomination en Français. Le bit est « le plus petit morceau d'information numérique ».

Pour les écrire, on regroupe les chiffres décimaux par paquets de 3, comme dans 1230 014 par exemple. En binaire on groupe les bits par 4, on écrira donc $17 = (1\ 0001)_2$.

La plupart du temps, en machine, les bits sont groupés par 8 (deux paquets de 4). Un tel paquet s'appelle un *octet*, et on écrit donc des *mots binaires* de longueur 8 tels que 0000 0011 : l'octet représente ici le nombre 3. Les bits à zéros ne sont pas inutiles.

Lorsqu'on considère un nombre écrit en binaire, on parle souvent de *bit de poids fort* et de *bit de poids faible* pour parler respectivement du bit associé à la plus grande puissance de 2, et du bit d'unités.

Considérons l'octet $(0010\ 0101)_2$. Son bit de poids fort est 0, son bit de poids faible est 1.

2 Écriture hexadécimale d'un entier naturel

La base « naturelle » de l'informatique est la base 2, mais elle n'est pas très pratique car elle donne lieu à des écritures trop longues. La base 10 nous paraît bien meilleure parce que nous avons l'habitude de l'utiliser, mais elle ne fait pas bon ménage avec la base 2 : il n'y a pas de méthode simple pour passer du décimal au binaire, et vice versa.

La base 16, ou base *hexadécimale*, est en revanche très adaptée à l'écriture des paquets de 4 bits, et par extension à celle des octets et autres écritures binaires.

En hexadécimal, on dispose de 16 chiffres : 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E et F.

Propriété : écriture binaire d'un entier naturel

Tout entier naturel possède une unique écriture en base 16, dite *écriture hexadécimale*. Plus précisément, soit $n \in \mathbf{N}$, alors il existe un unique entier $k \in \mathbf{N}$ et $k + 1$ nombres a_i , uniques et valant 0, 1, 2, ..., ou F et tels que

$$n = a_0 16^0 + a_1 16^1 + \dots + a_k 16^k$$

ce qui s'écrit aussi

$$n = \sum_{i=0}^k a_i 16^i$$

Remarque

On a vu une propriété similaire en base 2 et en fait elle est valable *dans toutes les bases* b (où b est un entier naturel supérieur ou égal à 2). Cela justifie par exemple l'utilisation de la base 20 ou de la base 150.

Les méthodes que l'on a vu en base 2 et 10 se transposent en base 16.

Méthode 4 : passer de la base 16 à la base 10

Déterminons l'écriture décimale de $(D4A)_{16}$:

$$\begin{aligned} (D4A)_{16} &= 13 \times 16^2 + 4 \times 16 + 10 \times 16^0 \text{ car D vaut 13 et A vaut 10.} \\ &= 3402 \end{aligned}$$

Méthode 5 : passer de la base 10 à la base 16

Déterminons maintenant l'écriture hexadécimale de 503 en utilisant la méthode des divisions successives par 16 :

$$\begin{array}{r|l} 503 & 16 \\ \hline 7 & 31 \\ \hline & 15 \\ & 1 \\ & 1 \\ & 0 \end{array}$$

$$503 = 31 \times 16 + \underline{7}.$$

$$31 = 1 \times 16 + \underline{15} \text{ et } 15 \text{ s'écrit } \underline{F}.$$

$1 = 0 \times 16 + \underline{1}$ et on arrête car le quotient est nul.

$$503 = (1F7)_{16}.$$

```

0x000340: C6 10 80 E3 02 10 01 E0 B0 10 C3 E1 00 30 0F E1 Å.ä...ä°.Ää.0.ä
0x000350: DF 30 C3 E3 1F 30 83 E3 03 F0 29 E1 3C 10 9F E5 B0Ää.0ä.ä)ä<.ä
0x000360: 0C 10 81 E0 00 00 91 E5 00 40 2D E9 00 E0 8F E2 .. ä..'ä.@-é.ä ä
0x000370: 10 FF 2F E1 00 40 BD E8 00 30 0F E1 DF 30 C3 E3 .ÿ/ä.@ä.0.äB0Ää
0x000380: 92 30 83 E3 03 F0 29 E1 0F 40 BD E8 B0 20 C3 E1 '0ä.ä)ä.@ä° Ää
0x000390: B8 10 C3 E1 00 F0 69 E1 1E FF 2F E1 28 73 00 03 ,.Ää.äiä.ÿ/ä(s..
0x0003A0: 90 34 00 03 F0 B5 47 46 80 B4 FF 20 E1 F1 FA FD '4...äµGF'ÿ äñúÿ
0x0003B0: A0 21 C9 04 27 4A 10 1C 08 80 00 F0 D3 FA 26 49 !É.'J...ä.óú&I
0x0003C0: 26 4A 10 1C 08 80 00 F0 F9 F8 00 F0 5B F9 DB F1 &J...ä.öüø.ä[ùÛñ
0x0003D0: 1B FA 00 F0 DF F8 F8 F0 1F FB 4B F0 73 FE 00 F0 .ú.äBøøä.ùKäsp.ä
0x0003E0: 6F F8 71 F0 BB FA 00 F0 07 FC 00 F0 1B FE 1C 48 øøqä»ú.ä.ü.ä.ä.p.H
0x0003F0: E0 21 49 02 02 F0 7A FB F7 F0 90 FC 19 48 00 24 à!I...äzú+ä ü.H.ä
0x000400: 04 70 19 48 04 70 F5 F0 1F F8 18 48 04 70 18 4F .p.H.pää.ø.H.p.O
0x000410: 00 21 88 46 06 1C 00 F0 E5 F8 12 48 00 78 00 28 .!äF...ääø.H.x.(
0x000420: 0E D1 39 8D 01 20 08 40 00 28 09 D0 0E 20 08 40 .Ñ9 . .@.(.ä. .@
0x000430: 0E 28 05 D1 DE F1 B8 FE DE F1 48 FE 00 F0 4A FA .(.Ññ,äñHä.äJú
0x000440: 57 F0 BE FF 01 28 15 D1 30 70 00 F0 2F F8 00 20 Wäÿ.(.Ñ0p.ä/ø.
0x000450: 30 70 22 E0 FF 7F 00 00 04 02 00 04 14 40 00 00 0p"äÿ!.....@..
0x000460: 00 00 00 02 80 34 00 03 1C 5E 00 03 34 30 00 03 ....ä....^...40..
0x000470: 40 30 00 03 0C 4D 00 20 28 70 00 F0 17 F8 57 F0 @0...M. (p.ä.øWä
0x000480: 69 FF 04 1C 01 2C 08 D1 00 20 F8 85 06 F0 16 FF iÿ....,Ñ. øä.ä.ÿ
0x000490: 2C 70 00 F0 0B F8 42 46 2A 70 54 F0 57 FA 71 F0 ,p.ä.øBF*pTäWúqä
0x0004A0: 67 FA 00 F0 F3 F9 B6 E7 34 30 00 03 00 B5 0A F0 gú.äóúä40...µ.ä
0x0004B0: 19 FE 00 06 00 28 01 D1 00 F0 28 F8 01 BC 00 47 .ä....(.Ñ.ä(ø.ä.G
0x0004C0: 10 B5 0B 48 00 24 04 62 44 62 04 60 09 48 00 F0 .µ.H.ä.bDä.ä.H.ä
0x0004D0: 37 F8 09 48 09 49 01 60 09 4A 0A 48 10 60 F2 20 7ø.H.I.ä.J.H.ä`ò
0x0004E0: 00 01 09 18 0C 60 08 48 04 70 10 BC 01 BC 00 47 .....ä.H.p.ä.ä.G

```

Un éditeur hexadécimal montre le contenu d'un fichier, d'un disque dur, ou de la RAM d'un ordinateur. La première colonne indique l'adresse, puis 16 octets écrits en hexadécimal et enfin les caractères correspondants.

3 Hexadécimal et binaire : un mariage heureux

Le grand avantage qu'apporte l'hexadécimal s'illustre facilement :

Méthode 6 : passer de la base 2 à la base 16

$$\begin{aligned}
 (101101000011101)_2 &= (0101\ 1010\ 0001\ 1101)_2 \\
 &= \left(\underbrace{0101}_5 \underbrace{1010}_A \underbrace{0001}_1 \underbrace{1101}_D \right)_2 \\
 &= (5A1D)_{16}
 \end{aligned}$$

x	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	10
1	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	10
2	2	4	6	8	A	C	E	10	12	14	16	18	1A	1C	1E	20
3	3	6	9	C	F	12	15	18	1B	1E	21	24	27	2A	2D	30
4	4	8	C	10	14	18	1C	20	24	28	2C	30	34	38	3C	40
5	5	A	F	14	19	1E	23	28	2D	32	37	3C	41	46	4B	50
6	6	C	12	18	1E	24	2A	30	36	3C	42	48	4E	54	5A	60
7	7	E	15	1C	23	2A	31	38	3F	46	4D	54	5B	62	69	70
8	8	10	18	20	28	30	38	40	48	50	58	60	68	70	78	80
9	9	12	1B	24	2D	36	3F	48	51	5A	63	6C	75	7E	87	90
A	A	14	1E	28	32	3C	46	50	5A	64	6E	78	82	8C	96	A0
B	B	16	21	2C	37	42	4D	58	63	6E	79	84	8F	9A	A5	B0
C	C	18	24	30	3C	48	54	60	6C	78	84	90	9C	A8	B4	C0
D	D	1A	27	34	41	4E	5B	68	75	82	8F	9C	A9	B6	C3	D0
E	E	1C	2A	38	46	54	62	70	7E	8C	9A	A8	B6	C4	D2	E0
F	F	1E	2D	3C	4B	5A	69	78	87	96	A5	B4	C3	D2	E1	F0
10	10	20	30	40	50	60	70	80	90	A0	B0	C0	D0	E0	F0	100

Table de multiplication hexadécimale.

Méthode 7 : passer de la base 16 à la base 2

$$(F7B)_{16} = \left(\underbrace{1111}_F \underbrace{0111}_7 \underbrace{1011}_B \right)_2$$

4 Additions

On pose l'opération à la main : c'est la même chose qu'en base 10.

En base 2

La seule différence avec la base 10 c'est que deux 1 donnent $(2)_{10}$ donc $(10)_2$, donc un zéro et une retenue de 1. Quand il y a deux 1 et une retenue de 1 en plus, cela donne $(3)_{10}$ donc $(11)_2$, donc un 1 et une retenue de 1.

Exemple

$$\begin{array}{rcccccccc}
 & & & 1 & & 1 & & 1 & \\
 & 1 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\
 + & & & & & 1 & 1 & 1 & 0 & 0 \\
 \hline
 = & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 1 &
 \end{array}$$

En base 16

C'est encore la même chose. Il faut bien se souvenir de la valeur de A, B, C, D, E et F.

Ajouter 8 et 3 ne provoque pas de retenue puisque $8 + 3 = 11$ et que 11 est B en base 16.

Dès que l'addition de 2 chiffres dépasse 15, il y aura une retenue : par exemple 9 et A donnent $(19)_{10}$, donc $(13)_{16}$. Ainsi on note 3 et une retenue de 1.

Exemple

$$\begin{array}{rcccc}
 & & & 1 & \\
 & 2 & 4 & 9 & 8 \\
 + & 1 & 7 & A & 3 \\
 \hline
 = & 3 & C & 3 & B
 \end{array}$$

5 Exercices**Exercice 1**

Donner l'écriture décimale des onze premières puissances de deux.

Exercice 2

1. Calculer $2^6 + 2^4 + 2^3 + 2^0$.
2. En déduire l'écriture binaire de 89.

Exercice 3

1. Calculer $2^7 + 2^3 + 2^2 + 2^1$.
2. En déduire l'écriture décimale de $(10001110)_2$.

Exercice 4

En utilisant la méthode 2, donner l'écriture binaire de

1. 56
2. 35
3. 13

Exercice 5

En utilisant la méthode 3, donner l'écriture binaire de

1. 142
2. 273
3. 1000

Exercice 6

- Donner l'écriture décimale de $(1101\ 1010)_2$.
- Donner l'écriture binaire de 2016.
- Donner l'écriture hexadécimale de 2016.

Exercice 7

- Donner les écritures décimales de $(11)_2$, $(111)_2$, $(1111)_2$.
- Soit $n \in \mathbf{N}$, conjecturer la valeur de $\left(\underbrace{1 \dots 1}_{n \text{ chiffres}} \right)_2$.

Exercice 8

Pour multiplier par dix un entier naturel exprimé en base dix, il suffit d'ajouter un 0 à sa droite, par exemple, $12 \times 10 = 120$.

Quelle est l'opération équivalente pour les entiers naturels exprimés en base deux ?

Exercice 9

1. Donner l'écriture binaire de 174.
2. Donner celle de 17.
3. Poser l'addition de 174 et 17 en binaire.

4. Donner l'écriture décimale du résultat et vérifier.

Exercice 10

1. Donner l'écriture hexadécimale de 1022.
2. Donner celle de 3489.
3. Poser l'addition de 1022 et 3489 en hexadécimal.
4. Donner l'écriture décimale du résultat et vérifier.

Exercice 11*

Le Roi d'un pays imaginaire fait frapper sa monnaie par 8 nains : chacun d'entre eux produit des pièces d'or de 10g chacune.

Un jour, son Mage lui annonce : « Majesté, mon miroir magique m'a prévenu que certains de vos nains vous volent. Ils prélèvent 1g d'or sur chaque pièce qu'ils frappent. Pour vous aider à trouver les voleurs, voici une balance magique. Elle est précise au gramme près et peut peser autant que vous voulez. Malheureusement elle ne peut être utilisée qu'une fois. »

Le lendemain, le Roi convoque les 8 nains en demandant à chacun d'apporter un coffre rempli de pièces d'or qu'il a frappées.

On suppose que

- chaque nain dispose d'autant de pièces que nécessaire;
- un nain honnête n'a que des pièces de 10g;
- un nain voleur n'a que des pièces de 9g;

Peux-tu aider le Roi pour démasquer les voleurs ?

Chapitre 2

Représentation des entiers

« Pour tout comprendre, lire ce chapitre en entier. »

Nous avons vu au chapitre précédent comment écrire les entiers naturels en binaire ou en hexadécimal. Maintenant nous allons étudier comment les entiers *relatifs*, c'est-à-dire positifs ou négatifs (on dit aussi *signés*) sont représentés en machine.

On va d'abord se limiter aux entiers naturels et on va voir qu'il n'y a pas de difficulté majeure à comprendre leur représentation en machine.

Remarque importante

Il existe des *centaines* de langages de programmation. Les principaux sont : C, C++, C#, JAVA, PYTHON, PHP et JAVASCRIPT (cette liste est non exhaustive). Chaque langage utilise ses propres *types de variable* mais en général il y a beaucoup de ressemblances. On travaillera donc sur des exemples de types de variables utilisés dans tel ou tel langage...sachant que ce type n'existe pas nécessairement en PYTHON.

1 Représentation des entiers naturels : l'exemple du unsigned char

En Informatique on dit souvent qu'un entier naturel est *non signé* (ou *unsigned* en anglais). Le type **unsigned char** se rencontre en C et en C++ (entre autres).

Un **unsigned char** est stocké sur un octet, c'est à dire 8 bits :

- l'octet 0000 0000 représente l'entier 0;
- 0000 0001 représente 1;
- et ainsi de suite jusqu'au plus grand entier représentable sur un octet : $(1111\ 1111)_2 = 255$.

On peut donc représenter les 256 premiers entiers avec un **unsigned char** et c'est logique : un octet, c'est 8 bits, chaque bit peut prendre 2 valeurs et $2^8 = 256$.

Si on a besoin de représenter des entiers plus grands, on pourra utiliser l'**unsigned short** : c'est la même chose mais ce type est représenté sur 2 octets. Donc on peut représenter les 2^{16} premiers entiers naturels, c'est à dire les nombres compris entre 0 et 65 535 inclus.

Cela continue avec l'**unsigned int** (sur 4 octets) et l'**unsigned long** (8 octets).

Remarque

En PYTHON c'est différent : le type **int** (abréviation de *integer*, qui veut dire « entier » en Anglais) permet de représenter des entiers arbitrairement grands, les seules limitations étant la mémoire de la machine. Il n'y a qu'à évaluer 2^{100000} dans un *shell* PYTHON pour s'en convaincre.

2 L'exemple du type char

Le type **char** (qui n'existe pas en PYTHON) utilise un octet et l'on veut représenter des entiers relatifs (donc plus seulement positifs).

2.1 Une première idée... qui n'est pas si bonne

On pourrait décider que le bit de poids fort est un bit de signe : 0 pour les positifs et 1 pour les négatifs, par exemple. Les 7 autres bits serviraient à représenter la valeur absolue du nombre. Puisqu'avec 7 bits on peut aller jusqu'à $(111\ 111)_2 = 127$, ce format permettrait de représenter tous les nombres entiers de -127 à 127.

Par exemple 1000 0011 représenterait -3 et 0001 1011 représenterait 27.

Il est un peu dommage que zéro ait 2 représentations : 0000 0000 et 1000 0000, mais ce qui est encore plus dommage c'est que lorsqu'on ajoute les représentations de -3 et de 27 voici ce qui se passe :

$$\begin{array}{rcccccccc}
 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\
 + & 0 & 0 & 0 & 1 & 1 & 0 & 1 & 1 \\
 \hline
 = & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 0
 \end{array}$$

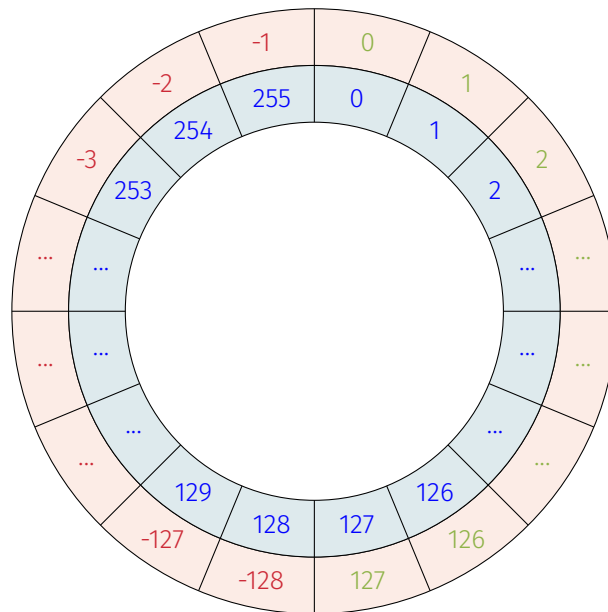
On obtient 10001 1110, qui représente -30... On aurait bien sûr préféré que cela nous donne 24.

2.2 La bonne idée : le complément à deux

Ce format, qui est utilisé avec le type **char**, permet de représenter les entiers de -128 à 127 :

Propriété : Représentation en complément à 2 sur un octet

- Soit x un entier positif plus petit ou égal à 127, alors on représente x par son écriture binaire (qui comprend 7 bits) et donc le bit de poids fort de l'octet (le 8^e) est égal à zéro.
- Sinon si x est un entier strictement négatif plus grand que -128, on le représente par l'écriture binaire de $256 + x$, qui est toujours représenté par un octet avec un bit de poids fort égal à 1.



Correspondance entre les valeurs « brutes » d'un octet (en bleu) et les entiers relatifs associés.

Exemples

Comment représenter 97 ? Ce nombre est positif, on le représente par son écriture binaire sur 8 bits : 0110 0001

Comment représenter -100 ? Ce nombre est négatif, il est donc représenté en machine par $256 - 100 = 156$, c'est-à-dire 1001 1100

Que représente 0000 1101 ? Le bit de poids fort est nul donc cela représente $(0000 1101)_2$, c'est à dire 13.

Que représente 1000 1110 ? Le bit de poids fort est non nul. $(1000\ 1110)_2 = 142$ représente x avec donc $256 + x = 142$, c'est-à-dire $x = -114$.

Méthode

Pour passer d'un nombre à son opposé en complément à 2, en binaire on procède de *la droite vers la gauche* et

- on garde tous les zéros et le premier 1;
- on « inverse » tous les autres bits.

Exemple

Si on veut l'écriture en complément à 2 de -44 on commence par écrire 44 en base 2 :

$$44 = (0010\ 1100)_2$$

Puis on applique la méthode précédente :

$$\begin{array}{ccccccc} \underline{00101} & \underline{100} \\ \text{on change} & \text{on garde} \end{array}$$

Ce qui nous donne

$$\begin{array}{ccccccc} \underline{11010} & \underline{100} \\ \text{on a changé} & \text{on a gardé} \end{array}$$

Ainsi la représentation de -44 en complément à 2 sur 8 bits est 1101 0100.

Ce qui est agréable, c'est que **l'addition naturelle est compatible avec cette représentation** dans la mesure où

- on ne dépasse pas la capacité : on n'ajoutera pas 100 et 120 car cela dépasse 127;
- si on ajoute deux octets et que l'on a une retenue à la fin de l'addition (ce serait un 9^e bit), alors celle-ci n'est pas prise en compte.

Exemple

Ajoutons les représentations de 97 et -100 :

$$\begin{array}{rcccccccc} & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ + & 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ \hline = & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \end{array}$$

On a $(1111\ 1101)_2 = 253$, il représente donc x sachant que $256 + x = 253$, c'est à dire $x = -3$.

On retrouve bien $97+(-100)=-3$.

Attention à ne pas dépasser la capacité du format : en ajoutant 120 et 20 on obtient 140 qui, étant plus grand que 128, représente $140-156 = -116$. C'est ce qu'affiche le programme suivant.

```
#include <iostream> // nécessaire pour utiliser cout

int main() // début de la fonction main
{
    char c1 = 120; // on définit une première variable
    char c2 = 20; // puis une deuxième
    char c3 = c1+c2; // on les ajoute
    std::cout << (int) c3; // on affiche le résultat en base 10
    return 0; // la fonction main renvoie traditionnellement zéro
}
```

3 Les principaux formats

En général, dans la majorité des langages (C, C++, C#, Java par exemple) les types suivants sont utilisés pour représenter les entiers relatifs (les noms peuvent varier d'un langage à l'autre) :

- **char** : Pour représenter les entiers compris entre -128 et +127. Nécessite un octet.
- **short** : Pour des entiers compris entre -2^{15} et $2^{15} - 1$. Codage sur 2 octets.
- **int** : (4 octets) entiers compris entre -2^{31} et $2^{31} - 1$.
- **long** : (8 octets) entiers compris entre $-9\,223\,372\,036\,854\,775\,808$ et $+9\,223\,372\,036\,854\,775\,807$ (-2^{63} et $2^{63} - 1$).

4 Quelques ordres de grandeur

- 1 kilooctet = 1 ko = 1 000 octets (fichiers textes)
- 1 mégaoctet = 1 Mo = 1 000 ko (fichiers .mp3)
- 1 gigaoctet = 1 Go = 1 000 Mo (RAM des PC actuels, jeux)
- 1 téraoctet = 1 To = 1 000 Go (Disques durs actuels)
- 1 pétaoctet = 1 Po = 1 000 To
- 1 exaoctet = 1 Eo = 1 000 Po = 10^{18} octets (trafic internet mondial mensuel en 2022 : 376 Eo)
- 1 zétaoctet = 1 Zo = 1 000 Eo (stockage des données prévu en 2025 sur Terre : 175 Zo)

5 Exercices

Exercice 12

Donner les représentations en complément à deux sur un octet de :
88, 89, 90, -125, -2 et -3.

Exercice 13

Donner les représentations en complément à 2 (sur un octet) de -1 et de 92.
Ajouter ces deux représentations (en ignorant la dernière retenue). Quel nombre représente cette somme ?

Exercice 14

On considère le code C++ suivant :

```
#include <iostream> // bibliothèque d'affichage
int main() // début de la fonction principale
{
    unsigned char c = 0; // on définit la variable c
    for (int i = 0; i < 300; i++) // on fait une boucle pour
    {
        std::cout << "valeur de i : " << i;
        // on affiche la valeur de i
        std::cout << " et valeur de c : " << (int)c;
        // on affiche la valeur de c en base 10
        std::cout << endl; // on revient à la ligne
        c++; // on augmente c
    }
    return 0; // la fonction principale renvoie zéro
}
```

Voilà ce que la console affiche :

```
valeur de i : 0 et valeur de c : 0
valeur de i : 1 et valeur de c : 1
et cætera
valeur de i : 254 et valeur de c : 254
valeur de i : 255 et valeur de c : 255
valeur de i : 256 et valeur de c : 0
valeur de i : 257 et valeur de c : 1
```

et cætera

valeur de i : 298 et valeur de c : 42

valeur de i : 299 et valeur de c : 43

Comment expliquer ceci?

Exercice 15

On considère le code C++ suivant :

```
#include <iostream> // bibliothèque d'affichage
int main() // début de la fonction principale
{
    char c = 0; // on définit la variable c
    for (int i = 0; i < 257; i++) // on fait une boucle pour
    {
        std::cout << "valeur de i : " << i;
        // on affiche la valeur de i
        std::cout << " et valeur de c : " << (int)c;
        // on affiche la valeur de c en base 10
        std::cout << endl; // on revient à la ligne
        c++; // on augmente c
    }
    return 0; // la fonction principale renvoie zéro
}
```

Voilà ce que la console affiche :

valeur de i : 0 et valeur de c : 0

valeur de i : 1 et valeur de c : 1

et cætera

valeur de i : 126 et valeur de c : 126

valeur de i : 127 et valeur de c : 127

valeur de i : 128 et valeur de c : -128

valeur de i : 129 et valeur de c : -127

et cætera

valeur de i : 254 et valeur de c : -2

valeur de i : 255 et valeur de c : -1

valeur de i : 256 et valeur de c : 0

Comment expliquer ceci?

Chapitre 3

Représentation des réels

« Tout cela est-il bien réel ? »

1 De la calculatrice à l'ordinateur

Dans une machine on *ne peut pas* représenter tous les nombres réels car la majorité a une écriture décimale illimitée et parmi celle-ci la majorité a une écriture décimale illimitée sans qu'aucun motif se répète, comme c'est le cas pour $\pi \simeq 3,141592\,653\,589\,793$.

Ceci dit on peut donner une valeur approchée d'un nombre réel r en écriture décimale en utilisant l'*écriture scientifique* vue au collège :

$$r \approx (-1)^s \times d \times 10^e$$

- s vaut 0 ou 1.
- d est un nombre décimal entre 1 inclus et 10 exclu.
- e est un entier relatif.

Ainsi, avec la calculatrice, on obtient :

$$5,4^{-5} \approx (-1)^0 \times 2,177866231 \times 10^{-4}$$

Dans ce cas, le nombre d comporte 10 chiffres décimaux, appelés *chiffres significatifs*.

Un ordinateur ne travaille pas avec des écritures décimales, mais avec des écritures *dyadiques* (l'équivalent des nombres décimaux en base 2).

Exemple : nombre décimal / nombre dyadique

- Considérons le nombre $x = 53,14$ (écrit en base 10). C'est un *nombre décimal* et on peut écrire :

$$\begin{aligned} x &= 53 + 1 \times 0,1 + 4 \times 0,04 \\ &= 5 \times 10^1 + 3 \times 10^0 + 1 \times 10^{-1} + 4 \times 10^{-2} \end{aligned}$$

- Les nombres dyadiques sont l'équivalent des nombres décimaux en base 2 : considérons $y = (101, 011)_2$, on peut écrire :

$$\begin{aligned} y &= 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} \\ &= 4 + 1 + 0,25 + 0,125 \\ &= 5,375 \end{aligned}$$

Exemple : écriture scientifique pour un nombre dyadique

- En reprenant l'exemple précédent on a $53,14 = 5,314 \times 10^1$, c'est une écriture scientifique.
- Pour $(101, 011)_2$, en se rappelant que multiplier par 2 décale la virgule d'un cran vers la droite, on a $(101, 011)_2 = (1, 01011)_2 \times 2^2$.

Il faut donc décider d'un *format de représentation* des nombres dyadiques dans l'ordinateur :

- Combien de bits significatifs pour le nombre dyadique ?
- Quelle est la plage de valeurs pour l'exposant ?

2 Le format IEEE 754

Ce format est une norme pour représenter les nombres dyadiques (notamment par PYTHON, en format 64 bits).

Par simplicité commençons par le format 32 bits (4 octets, donc). Voici comment cela fonctionne :

Définition

On considère un mot de 32 bits.

- Le premier bit s , indique le signe.
- Les 8 bits suivants $e_7 e_6 \dots e_0$ servent à coder l'exposant e de 2, qui vaut $e = (e_7 \dots e_0)_2 - 127$.
- Les 23 bits restants $m_1 \dots m_{23}$ servent à coder la *mantisse*, notons $m = (1, m_1 \dots m_{23})_2$.
- En définitive, ce mot de 32 bits représente

$$\boxed{(-1)^s \times m \times 2^e}$$

Ce qu'on peut noter

$$s e_7 \dots e_0 m_1 \dots m_{23} \mapsto (-1)^s \times (1, m_1 \dots m_{23})_2 \times 2^{(e_7 \dots e_0)_2 - 127}$$

$s e_7 \dots e_0 m_1 \dots m_{23}$ s'appelle une écriture *normalisée*.

Exemple : des 32 bits au nombre

Que représente 1 0100 0011 1110 0000 0000 0000 0000 000 ?

- $s = 1$.
- $e = (0100\ 0011)_2 - 127 = 67 - 127 = -60$
- $m = (\boxed{1}, 111000000000000000000000)_2 = 1 + 2^{-1} + 2^{-2} + 2^{-3} = 1,875$
- $1\ 0100\ 0011\ 1110\ 0000\ 0000\ 0000\ 0000\ 000 \rightsquigarrow (-1)^1 \times 1,875 \times 2^{-60}$

Cela fait environ $-1,6263032587282567 \times 10^{-18}$

Attention

En PYTHON au format 64 bits, les nombres sont représentés sur 8 octets :

- 1 bit de signe;
- 11 bits d'exposant (donc une plage de -1023 à 1024);
- 52 bits de mantisse.

Le principe est le même qu'en 32 bits.

3 Les limitations du format IEEE 754

Lorsqu'un format de représentation en virgule flottante est choisi (32 ou 64 bits), on ne peut pas représenter tous les nombres réels : il y a un plus grand nombre représentable (et son opposé pour les nombres négatifs) et un « plus petit nombre positif représentable » (le plus proche de zéro possible). De plus *on a automatiquement des valeurs approchées si le nombre que l'on veut représenter n'est pas de la forme $\frac{a}{2^n}$, avec $a \in \mathbf{Z}$ et $n \in \mathbf{N}$.*

Par exemple, un nombre aussi simple que 0,1 (c'est-à-dire $\frac{1}{10}$) n'est pas de la forme $\frac{a}{2^n}$, avec $a \in \mathbf{Z}$ et $n \in \mathbf{N}$, donc son écriture dyadique ne se termine pas.

On peut montrer que :

$$(0,1)_{10} = (0,0001\ 1001\ 1001\ 1001\ \dots)_2$$

C'est l'équivalent dyadique de

$$\frac{1}{3} = (0,3333\ \dots)_{10}$$

Et puisque PYTHON ne peut pas représenter intégralement le nombre 0,1 il l'approche du mieux qu'il peut : en fait pour PYTHON la valeur de 0,1 est :

0.1000000000000000055511151231257827021181583404541015625

Mais celui-ci a la gentillesse d'afficher 0.1.

Il faut se résigner à accepter les erreurs d'arrondis :

Python

```
>>> 0.1 + 0.1 + 0.1
0.30000000000000004
```

Cet exemple prouve que *tester l'égalité de 2 `float` n'a pas d'utilité*. On aura plutôt intérêt à *tester si leur différence est très petite*.

De même, l'addition de plusieurs `float` donne un résultat qui peut dépendre de l'ordre dans lequel on les ajoute. Elle *n'est pas non plus associative* :

$a + (b + c)$ et $(a + b) + c$ peuvent avoir des valeurs différentes.

Les erreurs d'arrondis se cumulent. Pour les minimiser on aura intérêt à appliquer la règle suivante :

Propriété : Règle de la photo de classe

Dans une somme de `float`, l'erreur est minimisée quand on commence par ajouter les termes de plus petite valeur absolue.

4 Exercices

Exercice 16

Donner l'écriture décimale des nombres suivants

- a. $(101, 1)_2$
- b. $(1, 011)_2$
- c. $(0, 1111\ 111)_2$ en remarquant que c'est « $(111\ 1111)_2$ divisé par 2^7 ».

Exercice 17

Écrire en base 2 les nombres suivants :

- a. 3,5
- b. 7,75
- c. 27,625

Exercice 18 : étendue du «gruyère»

On considère le format IEEE 754 64 bits : 1 bit de signe, 11 bits d'exposant (donc une plage de -1023 à 1024) et 52 bits de mantisse.

1. Quel est le plus grand nombre positif représentable ?
2. Quel est le plus petit nombre positif représentable ?

Exercice 19 : taille des trous du «gruyère»

On considère le format IEEE 754 64 bits : 1 bit de signe, 11 bits d'exposant (donc une plage de -1023 à 1024) et 52 bits de mantisse.

1. Quel est le flottant immédiatement plus grand que 1 ? Quelle distance les sépare ?
2. Quel est le flottant immédiatement plus grand que le plus petit nombre positif représentable ? Quelle distance les sépare ?
3. Quel est le flottant immédiatement plus petit que le plus grand nombre positif représentable ? Quelle distance les sépare ?

Exercice 20

Écrire un programme déterminant le plus petit entier n pour lequel PYTHON considère que $1 + 2^{-n} = 1$.

Faire le lien avec le format IEEE 754 64 bits.

Exercice 21

Faire calculer $1+2^{**}(-53)-1$ puis $1-1+2^{**}(-53)$.
Que remarque-t-on ?

Exercice 22

Écrire un programme déterminant le plus petit entier n pour lequel PYTHON considère que $2^{-n} = 0$.

En faisant le lien avec le format IEEE 754 64 bits, quelle valeur devrait-on trouver ?

Quelle explication peut-on imaginer (sachant qu'en PYTHON, les `float` sont bien codés sur 64 bits) ?

Exercice 23

On peut prouver (c'est dur) que

$$\sum_{n=1}^{+\infty} \frac{1}{n^4} = \frac{\pi^4}{90}$$

Appelons c cette constante.

1. Calculer à l'aide d'un script $S = \sum_{n=1}^{10^6} \frac{1}{n^4}$.

Ajoute-t-on des termes de plus en plus petits ou de plus en plus grands ?

Continuer le script pour afficher $c - S$ (on pourra utiliser `from math import pi`).

2. Calculer S « dans l'autre sens ».

Afficher $c - S$.

3. Qu'illustre cet exercice ?

Exercice 24 (quand nous aurons vu les fonctions)

1. Montrer qu'un triangle (3,4,5) est rectangle, ainsi qu'un triangle ($\sqrt{11}$, $\sqrt{12}$, $\sqrt{23}$).

2. Écrire une fonction `est_rectangle(a,b,c)` : qui renvoie `True` si `c**2 == a**2 + b**2` et, sinon, qui renvoie la différence entre `c**2` et `a**2+b**2`.

3. Tester la fonction `est_rectangle` avec les deux triangles précédents.

Que remarque-t-on ? Comment modifier la fonction pour qu'elle soit plus satisfaisante ?

Exercice 25**

On aimerait trouver l'écriture dyadique (illimitée) de $\frac{1}{3}$. On note donc

$$\frac{1}{3} = (0, a_1 a_2 a_3 \dots)_2$$

où a_i vaut 1 ou 0.

1. Expliquer pourquoi a_1 vaut *nécessairement* 0.
2. On note $x = \frac{1}{3}$. Montrer que x vérifie $4x = 1 + x$.
3. Quelle est l'écriture dyadique de $4x$?
4. Quelle est celle de $1 + x$?
5. En écrivant que ces 2 écritures représentent le même nombre, en déduire que

$$\frac{1}{3} = (0,0101\,0101\dots)_2$$

Chapitre 4

Représentation du texte

«Peux-tu décoder ce texte?»

1 Le code ASCII

Pour représenter les caractères que nous utilisons pour écrire, on a historiquement choisi d'associer *un numéro* (ou code) à chacun de ces caractères. La correspondance entre chaque caractère et son code était appelée un *Charset*.

Puisqu'à l'origine seul un petit nombre de caractères était utilisé (les caractères de base anglo-saxons), un octet suffisait pour les représenter tous.

Le fait de représenter en machine un jeu de caractères s'appelle réaliser un encodage (*encoding* en Anglais).

Le premier encodage utilisé fut l'ASCII, qui signifie *American Standard Code for Information Interchange*.

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

La table ASCII

Le code ASCII se base sur un tableau contenant les caractères les plus utilisés en langue anglaise : les lettres de l'alphabet en majuscule (de A à Z) et en minuscule (de a à z), les dix chiffres arabes (de 0 à 9), des signes de ponctuation (point, virgule, point-virgule, deux points, points d'exclamation et d'interrogation, apostrophe ou *quote*, guillemet ou *double quotes*, parenthèses, crochets etc.), quelques symboles et certains caractères spéciaux invisibles (espace, retour-chariot, tabulation, retour-arrière, etc.).

Exercice 26

Combien y a-t-il de caractères dans ce catalogue ?
Combien de bits sont nécessaires pour pouvoir représenter tous leurs numéros de code ?

Pour représenter ces symboles ASCII, les ordinateurs utilisaient des cases mémoires de un octet, mais ils réservaient toujours le huitième bit pour le contrôle de parité : c'est un procédé de sécurité pour éviter les erreurs, qui étaient très fréquentes dans les premières mémoires électroniques.

Méthode : contrôle d'erreur par parité

- On dispose d'un mot de 7 bits, par exemple 111 0011
- On compte le nombre de bits à 1, il y en a 5.
- On rajoute le bit de poids fort à 1 pour qu'en tout, il y ait toujours *un nombre pair* de bits à 1.
On obtient 1111 0011, ce bit de poids fort jouant le rôle de *code correcteur*.
- Un autre exemple : 001 1110 est codé 0001 1110.

Exercice 27

Voici un message reçu à l'issue d'une transmission :

53 E1 6C F5 70

Ces 6 octets sont censés représenter 6 caractères ASCII, codées sur 7 bits le 8^e étant réservé au contrôle d'erreur par parité.

1. Décoder ces 6 octets en disant s'il y a des erreurs ou non.
2. Quel était le message initial ?

2 L'insuffisance de l'ASCII

Pour coder les lettres accentuées, inutilisées en Anglais mais très fréquentes dans d'autres langues (notamment le Français), on a décidé d'étendre le codage des caractères au huitième bit (les erreurs-mémoire étant devenues plus rares et les méthodes de contrôle d'erreurs plus efficaces).

Exercice 28

Combien de nouveaux symboles a-t-on pu coder en autorisant le huitième bit dans le codage ?

On a alors pu coder toutes ces lettres et ainsi que de nouveaux caractères typographiques utiles tels que différents tirets.

3 Le problème

Le fait d'utiliser un bit supplémentaire a bien entendu ouvert des possibilités mais malheureusement les caractères de toutes les langues ne pouvaient être pris en charge en même temps. La norme ISO 8859-1 appelée aussi Latin-1 ou Europe occidentale est la première partie d'une norme plus complète appelée **ISO 8859** (qui comprend 16 parties) et qui permet de coder tous les caractères des langues européennes. Cette norme ISO 8859-1 permet de coder 191 caractères de l'alphabet latin qui avaient à l'époque été jugés essentiels dans l'écriture, mais omet quelques caractères fort utiles (ainsi, la ligature œ n'y figure pas).

Dans les pays occidentaux, cette norme est utilisée par de nombreux systèmes d'exploitation, dont Linux et Windows. Elle a donné lieu à quelques extensions et adaptations, dont **Windows-1252** (appelée **ANSI**) et ISO 8859-158 (qui prend en compte le symbole € créé après la norme ISO 8859-1). C'est une source de grande confusion pour les développeurs de programmes informatiques car un même caractère peut être codé différemment suivant la norme utilisée. Voici les tableaux décrivant deux encodages :

Windows-1252 (CP1252)																
	x0	x1	x2	x3	x4	x5	x6	x7	x8	x9	xA	xB	xC	xD	xE	xF
0x	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1x	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2x	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3x	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4x	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5x	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6x	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7x	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL
8x	€		,	ƒ	„	…	†	‡	^	‰	Š	€	CE		Ž	
9x		‘	’	“	”	•	—	—	™	š	›	œ		ž	ÿ	
Ax	NBSP	ı	¢	£	¤	¥	¦	§	¨	©	ª	«	¬	®	¯	
Bx	°	±	²	³	´	µ	¶	·	,	ı	º	»	¼	½	¾	¿
Cx	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
Dx	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
Ex	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
Fx	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

ISO/CEI 8859-15																
	x0	x1	x2	x3	x4	x5	x6	x7	x8	x9	xA	xB	xC	xD	xE	xF
0x																
1x	non utilisé															
2x	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	
3x	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4x	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5x	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6x	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7x	p	q	r	s	t	u	v	w	x	y	z	{		}	~	
8x	non utilisé															
9x	non utilisé															
Ax	ı	¢	£	€	¥	¦	§	¨	©	ª	«	¬	®	¯		
Bx	°	±	²	³	Ž	µ	¶	·	ž	ı	º	»	¼	½	¾	¿
Cx	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
Dx	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
Ex	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
Fx	ð	ñ	ò	ó	ô	õ	÷	ø	ù	ú	û	ü	ý	þ	ÿ	

Exercice 29

Ces deux encodages sont-ils totalement compatibles ? Pourquoi ?

4 La multiplicité des encodages

Au fil du temps une multitude d'encodages sont apparus, multipliant les sources de confusion. Voici pour l'exemple une partie des encodages que PYTHON reconnaît :

Encodage	Alias Python	Langues concernées
ascii	646,us-ascii	English
big5	big5-tw, csbig5	Traditional Chinese
cp424	EBCDIC-CP-HE, IBM424	Hebrew
cp437	437, IBM437	English
cp500	EBCDIC-CP-BE, EBCDIC-CP-CH, IBM500	Western Europe
cp720		Arabic
cp737		Greek
cp856		Hebrew
cp857	857, IBM857	Turkish
cp864	IBM864	Arabic
cp874		Thai
cp932	932, ms932, mskanji, ms-kanji	Japanese
cp1251	windows-1251	Bulgarian, Byelorussian, Macedonian, Russian, Serbian
cp1258	windows-1258	Vietnamese
euc_kr	euckr, korean, ksc5601, ks_c-5601, ks_c-5601-1987	Korean
gbk	936, cp936, ms936	Unified Chinese
latin_1	iso-8859-1, iso8859-1, 8859, cp819, latin, latin1, L1	West Europe
iso8859_14	iso-8859-14, latin8, L8	Celtic languages
koi8_r		Russian
utf_8	U8, UTF, utf8	all languages

Exercice 30

Lequel de ces encodages semble le plus performant ?

5 L'Unicode

La globalisation des échanges culturels et économiques a mis l'accent sur le fait que les langues européennes coexistent avec de nombreuses autres langues aux alphabets spécifiques voire sans alphabet (le Japonais utilise entre autres un syllabaire, chaque symbole représentant une syllabe). La généralisation de l'utilisation d'Internet dans le monde a ainsi nécessité une prise en compte d'un nombre beaucoup plus important de caractères (à titre d'exemple, le mandarin possède à lui tout seul plus de 5000 caractères!).

Une autre motivation pour cette évolution résidait dans les possibles confusions dues au trop faible nombre de caractères pris en compte; ainsi, les symboles monétaires des différents pays n'étaient pas tous représentés dans le système ISO 8859-1, de sorte que les ordres de paiement internationaux transmis par courrier électronique risquaient d'être mal compris. La norme Unicode a donc été créée pour permettre le codage de textes écrits quel que soit le système d'écriture utilisé.

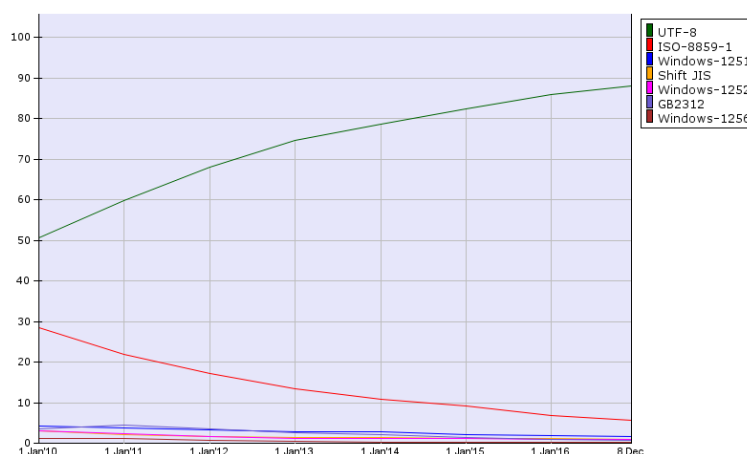
Dans le système UTF-8, on attribue à chaque caractère un nom, une position normative et un bref descriptif qui seront les mêmes quelle que soit la plate-forme informatique ou le logiciel utilisés.

Un consortium composé d'informaticiens, de chercheurs, de linguistes et de personnalités représentant les états ainsi que les entreprises s'occupe d'unifier toutes les pratiques en un seul et même système : *l'Unicode*.

Définition

L'Unicode est une table de correspondance Caractère-Code (Charset), et *l'UTF-8* est l'encodage correspondant (Encoding) le plus répandu.

De nos jours, par défaut, les navigateurs Internet utilisent le codage UTF-8 et les concepteurs de sites pensent de plus en plus à créer leurs pages web en prenant en compte cette même norme ; c'est pourquoi il y a de moins en moins de problèmes de compatibilité : l'UTF-8 est aujourd'hui majoritairement utilisé pour les sites du web, comme le montre ce graphique.



L'UTF-8 est également le codage le plus utilisé dans les systèmes GNU, Linux et compatibles pour gérer le plus simplement possible des textes et leurs traductions dans tous les systèmes d'écritures et tous les alphabets du monde.

6 L'UTF-8 côté technique

La norme Unicode définit entre autres un ensemble (ou répertoire) de caractères. Chaque caractère est repéré dans cet ensemble par un index entier aussi appelé « point de code ».

Par exemple le caractère « € » (euro) est le 8365^{ème} caractère du répertoire Unicode, son index, ou point de code, est donc 8364 (on commence à compter à partir de 0).

Le répertoire Unicode peut contenir plus d'un million de caractères, ce qui est bien trop grand pour être codé par un seul octet (limité à des valeurs entre 0 et 255). La norme Unicode définit donc des méthodes standardisées pour coder et stocker cet index sous forme de séquence d'octets : UTF-8 est la plus utilisée d'entre elles (il y a aussi des variantes comme UTF-16 et

UTF-32). En UTF-8, tout caractère est codé sur 1, 2, 3 ou 4 octets.

La principale caractéristique d'UTF-8 est qu'elle est *rétro-compatible avec la norme ASCII*, c'est-à-dire que tout caractère ASCII se code en UTF-8 sous forme d'un unique octet, identique au code ASCII.

Par exemple «A» (A majuscule) a pour code ASCII 65 et se code en UTF-8 par l'octet 65, il en va de même pour tous les caractères ASCII. Pour les autres, on procède comme ceci :

- Chaque caractère est associé à son index Unicode.
- En général, cet index est exprimé en hexadécimal. Actuellement, presque toutes les valeurs de 0000 à FFFF, c'est-à-dire de 0 à 65535 sont attribuées à des « alphabets » associés à des langues, des plus communes aux plus rares, et à divers symboles, tels que les symboles mathématiques. Au delà de FFFF on trouve des alphabets associés à des langues anciennes (cunéiformes, hiéroglyphes...).
- En fonction du nombre de bits nécessaires pour représenter en binaire cet index, on utilise le codage suivant :

Nombre de bits de l'index	Nombre d'octets pour coder en UTF-8	Schéma de codage
de 0 à 7	1	0xxx xxxx
de 8 à 11	2	110 x xxxx 10 xx xxxx
de 12 à 16	3	1110 xxxx 10 xx xxxx 10 xx xxxx
de 17 à 21	4	1111 0xxx 10 xx xxxx 10 xx xxxx 10 xx xxxx

Par exemple le symbole € a un index Unicode qui vaut 8364.

- 8364 s'écrit 20AC en hexa, ce qui fait 10 0000 1010 1100 en binaire, soit 14 bits.
On va donc utiliser 3 octets pour coder, conformément au schéma de codage.
- On commence par écrire le mot de 16 bits correspondant : 0010 0000 1010 1100.
- On formate comme à la 3ème ligne du tableau :

1110 0010 **1000** 0010 **1010** 1100

Ce qui fait 3 octets : E2 82 AC en hexadécimal.

Exercice 31

Si un ordinateur lit cet encodage UTF-8 du symbole € selon l'encodage ISO8859-15, qu'affichera-t-il ?

7 Conclusion

Même si l'encodage UTF-8 devient le standard international, certains développeurs, sites, ou applications en utilisent malgré tout encore d'autres.

Propriété

La notion de texte brut n'existe pas en informatique : lorsqu'un ordinateur lit un fichier texte il n'a *a priori* aucun moyen de savoir quel est son encodage.

Beaucoup de documents indiquent donc en en tête leur format d'encodage : en HTML, on écrira dans l'en-tête d'une page :

HTML

```
<meta charset="utf8"/>
```

pour préciser qu'elle est encodée en UTF-8.

8 Et Python dans tout ça ?

En PYTHON, on pourra aussi écrire : `# -*- coding: utf8 -*-` en première ligne de tout script pour signifier la même chose, *et cætera*.

PYTHON gère très bien les encodages. On peut fabriquer un convertisseur très rapidement :

Python

```
fichier = open("nom_fichier", 'rt', encoding="utf8")
texte = fichier.read()
fichier.close()
```

Ceci permet de lire le contenu d'un fichier texte (d'où le `'rt'`, pour 'read text') d'un fichier texte encodé en UTF-8, et de le stocker dans la variable `texte`, de type `str`.

Python

```
fichier = open("nom_fichier", 'wt', encoding="utf8")  
fichier.write("Salut")  
fichier.close()
```

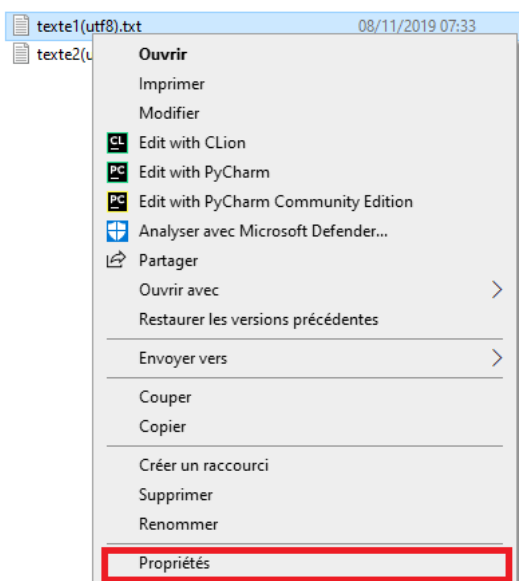
Permet d'écrire un fichier texte (d'où le 'wt', pour write text) en UTF-8.

9 Exercices

Exercice 32

Analyser les deux fichiers `texte1(utf8).txt` et `texte2(utf8).txt` : ils sont tous les deux encodés en UTF-8.

- Ouvrir chaque fichier. Quelles sont les différences de contenu entre ces deux fichiers?
- Faire un clic droit puis **Propriétés** comme ceci :



Quelles sont les tailles de ces deux fichiers?
Comment, dans le détail, la différence de taille s'explique-t-elle?

Exercice 33

Nous allons examiner des problèmes d’encodage.

- Ouvrir le fichier `index1(utf8).html` avec un navigateur, puis le fichier `index2(utf8).html` avec un navigateur.
Que remarquez-vous?
- Ouvrir ce dernier fichier avec un éditeur de texte (comme le bloc-notes Windows).
D’où vient le problème ? Proposer une correction.
Expliquer en détail pourquoi à la place des « é », il y a des « Ã© ».
- Ouvrir le fichier `index3(IS08859-15).html` avec un navigateur.
D’où vient le problème ? Proposer une correction.
Expliquer en détail pourquoi il y a des « points d’interrogation dans des losanges noirs ».

Exercice 34

Écrire un script qui corrige le problème de `index3(IS08859-15).html` en le convertissant en UTF-8 (utiliser les scripts PYTHON fournis à la section précédente).

Table des matières

1	Bases de numération	3
1	Écriture binaire d'un entier naturel	3
1.1	Pourquoi le binaire ?	4
1.2	Comprendre l'écriture en base 2	4
1.3	Un algorithme pour déterminer l'écriture binaire d'un entier naturel . . .	6
1.4	Vocabulaire	7
2	Écriture hexadécimale d'un entier naturel	7
3	Hexadécimal et binaire : un mariage heureux	9
4	Additions	10
5	Exercices	11
2	Représentation des entiers	15
1	Représentation des entiers naturels : l'exemple du unsigned char	15
2	L'exemple du type char	16
2.1	Une première idée... qui n'est pas si bonne	16
2.2	La bonne idée : le complément à deux	16
3	Les principaux formats	19
4	Quelques ordres de grandeur	19
5	Exercices	20
3	Représentation des réels	23
1	De la calculatrice à l'ordinateur	23
2	Le format IEEE 754	24
3	Les limitations du format IEEE 754	25
4	Exercices	26
4	Représentation du texte	31
1	Le code ASCII	31
2	L'insuffisance de l'ASCII	33
3	Le problème	33
4	La multiplicité des encodages	34

5	L'Unicode	35
6	L'UTF-8 côté technique	36
7	Conclusion	38
8	Et Python dans tout ça?	38
9	Exercices	39