

## A. Codage des nombres réels :

**Ce qu'il faut savoir :** Le codage sur 64 bits (8 octets) en binaire des nombres réels pour les représenter en machine est défini par **la norme IEEE 754**.

Pour représenter un nombre à virgule en binaire, on va utiliser une représentation similaire à la notation scientifique en base 10, mais ici en base 2 : il s'agit de **la représentation en virgule flottante**.

Cette représentation est  $S m \times 2^n$  où  $S$  est **le signe** (+ ou -),  $m$  est **la mantisse** ( $1 \leq m < 2$ ) et  $n$  **l'exposant** qui est un nombre entier relatif.

On codera en binaire ce nombre ainsi représenté sur 64 bits ainsi :

64è bit	63è bit	61è bit	.....	53è bit	52è bit	51è bit	50è bit	49è bit	.....	3è bit	2è bit	1er bit
signe		exposant décalé				m - 1 = partie décimale de la matisse						

**1 bit**  
 (0 code + et 1 code -)

**11 bits** codent un nombre entier naturel  
 $d = \text{exposant décalé} = n + 1\ 023$   
 ( $0 \leq d \leq 2\ 047$ ).

Si  $d=0$  ou  $2\ 047$ , cela est réservé à des situations exceptionnelles telles que **l'infini** ou « **not a number** ».

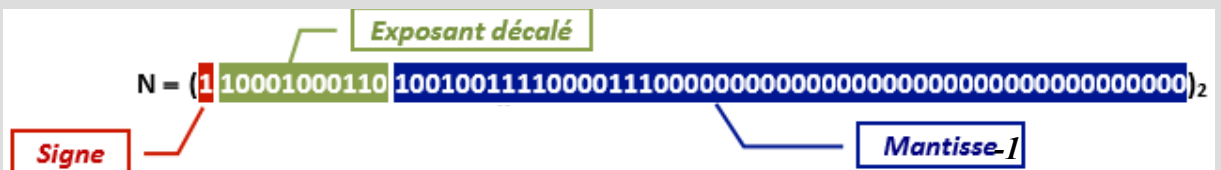
Si  $1 \leq d \leq 2\ 046$ , on soustrait **1 023** pour obtenir l'exposant  $n$  qui est un entier relatif.

**52 bits** codent  $m-1$  avec  $1 \leq m < 2$ .

La partie entière de  $m$  est toujours égale à **1**, il est inutile de coder ce chiffre et de perdre 1 bit.

$m-1$ , la partie décimale de  $m$  est une somme de puissances négatives de 2.

Exemple :



**Signe = 1 :** le nombre **est négatif**.

**Exposant décalé  $d$**  =  $(10001000110)_2 = 2^{10} + 2^6 + 2^2 + 2^1 = 1024 + 64 + 4 + 2 = 1\ 094$ .

Donc **l'exposant  $n$**  =  $1\ 094 - 1\ 023 = 71$ .

**Mantisse** =  $(1,10010011110000111000\dots)_2 = 1 + 2^{-1} + 2^{-4} + 2^{-7} + 2^{-8} + 2^{-9} + 2^{-10} + 2^{-15} + 2^{-16} + 2^{-17} \approx 1,5772$ .

La partie entière de la mantisse est toujours 1.

DONC ce code représente le réel :  $-1,5772 \times 2^{71} \approx -3,724 \times 10^{21}$ .

Cas particulier : On convient que le réel **0 est codé par**  $(0\ 0000000000\ 0000\dots 00)_2$ .



g. Déduis-en le code au format IEEE 754 de -118,375.

**Donc le code de -118,375 au format IEEE 754 est  $(1\ 10000000101\ 110110011000\dots)_2$ .**

3. En reproduisant la méthode de la question 2., détermine le code au format IEEE 754 de ces nombres réels.

a. 75,281 25.

$$75,281\ 25 = 1,1762695313 \times 2^6 \\ (0\ 10000000101\ 00101101001000\dots)_2$$

b. -0,687 5.

$$-0,6875 = -1,375 \times 2^{-1} \\ (1\ 01111111110\ 011000\dots)_2$$

## Activité 2 :

1. a. Comment code-t-on en binaire le nombre entier -7 (sur 8 bits) ? Et le nombre réel -7,0 (sur 64 bits) ?

**Le nombre entier relatif -7 se code sur 8 bits (avec la méthode du complément à 2) par l'entier naturel**

**$-7 + 2^8 = 249$  . On obtient donc  $(11111001)_2$  .**

**$-7 = -1,75 \times 2^2$  , ainsi le code au format IEEE 754 est :  $(1\ 10000000001\ 110000000000000000000000000000\dots)_2$ .**

b. Pourquoi le même nombre n'est-il pas codé de la même façon ? **Parce que les types sont différents.**

Quelle précaution faudra-t-il prendre lorsqu'on écrit un programme (notamment avec Python) ?

**Lorsqu'on écrit un programme il est donc essentiel, lorsqu'on définit une variable, de savoir de quel type elle est (ici l'entier -7 est de type « integer » et le réel -7,0 est de type « float »).**

2. a. Applique la méthode de l'activité 1 pour tenter de convertir le nombre réel 0,2 au format IEEE 754.

Que se passe-t-il ?

**$0,2 \times 2 = 0,4$  ;  $0,4 \times 2 = 0,8$  ;  $0,8 \times 2 = 1,6$  . Donc  $0,2 = 1,6 \times 2^{-3}$  .**

**Le signe est + donc le premier bit est 0.**

**L'exposant décalé est  $-3 + 1023 = 1020$  qui se code  $(0111111100)_2$ .**

**La partie décimale de la mantisse est 0,6.**

**$0,6 \times 2 = 1,2$  donc le premier bit est 1.**

**$0,2 \times 2 = 0,4$  donc le 2<sup>è</sup> bit est 0.**

**$0,4 \times 2 = 0,8$  donc le 3<sup>è</sup> bit est 0.**

**$0,8 \times 2 = 1,6$  donc le 4<sup>è</sup> bit est 1.**

**Comme  $1,6 - 1 = 0,6$  c'est cyclique, et les bits « 1001 » vont se répéter indéfiniment =  $(10011001\dots1001\dots)_2$ .**

b. Quand les 52 bits de la mantisse ne suffisent pas, si le 53<sup>è</sup> bit est 1, on arrondit le bit de la mantisse qui est le plus à droite par excès. Mais si le 53<sup>è</sup> bit est 0, on arrondit le bit de la mantisse qui est le plus à droite par défaut. Quel sera finalement le code de 0,2 au format IEEE 754 ?

**La mantisse se code  $(100110011001\dots)_2$  On a donc  $(1001\ 1001\ 1001\dots1001\ 1\dots)_2$  , en arrondissant :**

  
53<sup>è</sup> bit

**$(10011001100\dots1010)_2$ .**

**Ainsi le code au format IEEE 754 de 0,2 est  $(0\ 01111111100\ 10011001100\dots1010)_2$ .**

c. Déduis-en le code au format IEEE 754 du réel 0,1. Puis applique la méthode de l'activité 1 pour trouver celui du réel 0,3.

**On a vu que  $0,2 = (1 + 2^{-1} + 2^{-4} + 2^{-5} + 2^{-8} + 2^{-8} + \dots) \times 2^{-3}$  . En divisant par 2 (ou en multipliant par  $2^{-1}$ ) on obtient  $0,1 = (1 + 2^{-1} + 2^{-4} + 2^{-5} + 2^{-8} + 2^{-8} + \dots) \times 2^{-4}$  . Donc seul l'exposant change (c'est -4).**

**Ainsi l'exposant décalé est  $-4 + 1023 = 1019$  qui se code  $(01111111011)_2$  .**

**Le code au format IEEE 754 de 0,1 est  $(0\ 01111111011\ 10011001100\dots1010)_2$ .**

On a  $0,3 \times 2 = 0,6$  ;  $0,6 \times 2 = 1,2$  donc  $0,3 = 1,2 \times 2^{-2}$ .

Le bit correspondant au signe + est 0.

L'exposant décalé est  $-2+1\ 023=1\ 021$  qui se code  $(0111111101)_2$  sur 11 bits.

La partie décimale de la mantisse est 0,2.

$0,2 \times 2 = 0,4$  donc le premier bit est 0.

$0,4 \times 2 = 0,8$  donc le 2<sup>è</sup> bit est 0.

$0,8 \times 2 = 1,6$  donc le 3<sup>è</sup> bit est 1.

$0,6 \times 2 = 1,2$  donc le 4<sup>è</sup> bit est 1.

Comme  $1,2-1=0,2$  c'est cyclique, et les bits « 0011 » vont se répéter indéfiniment =  $(00110011...0011...)_2$ .

La mantisse se code  $(0011\ 0011\ 0011...0011\ 0)_2$ , en arrondissant :  $(0011\ 0011\ 0011...0011)_2$ .

53<sup>è</sup> bit

Ainsi le code au format IEEE 754 de 0,3 est  $(0\ 0111111101\ 001100110011...0011)_2$ .

d. Ouvre Python et tape dans la console : `>>>(0.1+0.2)-0.3` puis « Entrée ».

Pourquoi le résultat est-il surprenant ? Comment l'expliquer ?

```
>>> (0.1+0.2)-0.3
```

```
5.551115123125783e-17
```

Ce résultat différent de 0 est surprenant car en théorie le résultat devrait être 0.

Ceci s'explique par les arrondis réalisés sur les mantisses au 52<sup>è</sup> chiffre.

### Ce qu'il faut savoir :

- Il est impératif de bien distinguer le type « **integer** » du type « **float** » dans un programme. En particulier l'entier 7 et le flottant 7.0 ne se codent pas de la même manière en binaire dans une machine.

- En binaire, sur 8 octets (64 bits), avec **la méthode du complément à 2**, les codes représentent **exactement** les nombres entiers compris entre -9 223 372 036 854 775 808 et 9 223 372 036 854 775 807.

En revanche avec **la norme IEEE 754** (sur 64 bits aussi), les codes représentent **parfois des approximations** des nombres réels compris entre  $-1,7976931348623157 \times 10^{308}$  et  $1,7976931348623157 \times 10^{308}$ .

- Sur machine les calculs avec les nombres entiers sont exacts, mais les calculs avec les nombres flottants peuvent présenter d'infimes erreurs d'arrondi.

## B. Exercices d'application :

**Exercice 1 :** Détermine à quels nombres réels correspondent ces codes au format IEEE 754.

a.  $(1\ 0111111110\ 1000...0)_2$

signe = -

exposant décalé = 1 022

exposant =  $1\ 022 - 1\ 023 = -1$

mantisse =  $1 + 2^{-1} = 1,5$

le réel est  $-1,5 \times 2^{-1} = -0,75$ .

b.  $(0\ 10000001000\ 0001100101011100...0)_2$

signe = +

exposant décalé = 1 032

exposant =  $1\ 032 - 1\ 023 = 9$

mantisse =  $1 + 2^{-4} + 2^{-5} + 2^{-8} + 2^{-10} + 2^{-12} + 2^{-13} + 2^{-14} = 1,09906005859375$

le réel est  $1,09906005859375 \times 2^9 = 562,71875$ .

c.  $(0\ 00000000000\ 000...0)_2$

d'après le cours, le réel est 0.

**Exercice 2 :** Convertis les nombres suivants au format IEEE 754.

Pour chaque cas, précise si les codes représentent exactement ou approximativement les nombres réels.

**a.** 3,625.

**b.** -4,5.

**c.** -50.

**d.**  $-\frac{11}{15} = -0,7333 \dots$

**e.** 1.

$3,625 \div 2^1 = 1,8125 \quad (1\ 10000000001\ 00100\dots00)_2$ .

$(1\ 01111111110\ 01110111\dots0111)_2$ .

$3,625 = +1,8125 \times 2^1$

$(1\ 10000000100\ 100100\dots00)_2$ .

$(0\ 01111111111\ 00\dots)_2$ .

*signe + : 0*

*expo décalé 1024 : 10...0*

*partie décimale mantisse 0,8125 :*

$0,8125 \times 2 = 1,625$  donc bit 1

$0,625 \times 2 = 1,25$  donc bit 1

$0,25 \times 2 = 0,5$  donc bit 0

$0,5 \times 2 = 1$  donc bit 1 et fin

*partie décimale de la mantisse : 110100...00*

*Donc 3,625 se code  $(0\ 10000000000\ 110100\dots00)_2$ .*