

TD 3

Exercice 1. Représentation des nombres à virgule flottante

1 Donner l'écriture des nombres 1, 2, 3, 4 et 5 au format IEEE 754 en simple précision.

2 Si l'exposant biaisé est différent de 0 et 255 (ou 2047 en double précision), quels sont les plus petits et les plus grands nombres strictement positifs représentables en simple et double précision ?

3 Même question avec un exposant biaisé valant 0 (représentation dé-normalisée). Quel est le plus petit nombre strictement positif représentable ?

Exercice . Représentation des nombres à virgule flottante

1 Donner l'écriture des nombres 1, 2, 3, 4 et 5 au format IEEE 754 en simple précision.

$1 = + 1,0 \times 2^0$ d'où $\text{signe} = 0$, $\text{mantisse} = 000\dots000$,
 $\text{exposant} = +127_{10} = 01111111$
soit : **0 01111111 000000000000000000000000**

$2 = 10 = 1,0 \times 2^1$ soit : **0 10000000 000000000000000000000000**

$3 = 11 = 1,1 \times 2^1$ soit : **0 10000000 100000000000000000000000**

$4 = 100 = 1,0 \times 2^2$ soit : **0 10000001 000000000000000000000000**

Exercice . Représentation des nombres à virgule flottante

Si l'exposant biaisé est différent de 0 et 255 soit 1 à 254 (ou 2047 en double précision), quels sont les plus petits et les plus grands nombres strictement positifs représentables en simple et double précision ?

Simple	$1.0 \cdot 2^{-126} / 1,17 \cdot 10^{-38}$	$1,1111 \cdot 2^{127} / 3,4 \cdot 10^{38}$
Double	$1.0 \cdot 2^{-1022} / 2,22 \cdot 10^{-308}$	$1,111 \cdot 2^{1023} / 1,79 \cdot 10^{308}$

Exemple de valeurs (*source Wikipédia*)

Age de l'univers : $13,8 \cdot 10^9$ ans =
13,8 milliards d'années

Taille de l'univers : $8,8 \cdot 10^{23}$ Km =
880 000 milliards de milliards de
Km

Poids du soleil : $2 \cdot 10^{30}$ Kg

Volume de l'univers : 10^{81} M³

Taille d'un atome : 10^{-10} à 10^{-15} m

Poids d'un atome : 10^{-25} à 10^{-27} kg

Mais encore

$341 \cdot 10^6$ ordinateurs vendus en 2022 dans le monde

$6,2 \cdot 10^9$ ordinateurs en service dans le monde

$3,1 \cdot 10^6$ ordinateurs vendus en France en 2021

$2,5 \cdot 10^6$ ordinateurs mis au rebus en France en 2021

A méditer

Exercice . Représentation des nombres à virgule flottante

Même question avec un exposant biaisé valant 0 (représentation dé-normalisée). Quel est le plus petit nombre strictement positif représentable ?

avec exp biaisé = 0

	Mini	Maxi
Simple	$0,000 \ 2^{-127} / 1,4 \ 10^{-45}$	$0,111 \ 2^{-127} / 1,17 \ 10^{-45}$
Double	$0,000 \ 2^{-1023} / 5,0 \ 10^{-324}$	$0,111 \ 2^{-1023} / 2,22 \ 10^{-324}$

Donner les représentations en simple précision à virgule flottante de 25 et 2,125.

2⁵

2⁵ = 100000 = 1,0 x 2⁵ → signe = 0 exposant = 127+5=132 = 10000100 Mantisse = 000000000

2⁵ = 0 10000100 000000000000000000000000

2,125

2₁₀ = 10₂

0,125₁₀ 0,125 * 2 = 0,250 → 0

0,250 * 2 = 0,500 → 0

0,500 * 2 = 1,0 → 1

0 * 2 = 0 → 0

d'où 0,125₁₀ = 0,001₂ et 2,125₁₀ = 10,001₂

2,125 = 10,001 x 2⁰ = 1,0001 x 2¹ soit signe = 0 , mantisse = 00010... Exposant = 128 = 10000000

= 0 10000000 000100000000000000000000

Quelle est la représentation en simple précision à virgule flottante de $1/10$? Quelle est l'erreur obtenue ?

$$\begin{aligned} 0,1 &= 0,000110011001100110011\dots \times 2^0 \\ &= 1,100110011001100110011\dots \times 2^{-4} \end{aligned}$$

Soit $s = 0$, $\text{Exp} = 127 + (-4) = 123$

$$0,1 = \mathbf{0\ 01111011\ 10011001100110011001100} \quad \text{en IEEE 754}$$

La régénération donne 0,099998 soit presque 0,1

Même question pour $1/5$

$$1/5 = 0,2 = 0,1 \times 2 \text{ d'où } 1/5 = 0,00110011001100110011...$$

$$1/5 = 0 \text{ 01111100 10011001100110011001100}$$

La régénération donne 0,199998 soit presque 0,2

Perte d'information en arithmétique flottante

(Vous pouvez écrire de petits programmes sur machine pour vous aider)

Soient les nombres

$X = 0\ 10001110\ 000000000000000000000000$

$Y = 0\ 10011010\ 000000000000000000000000$

au format IEEE 754 en simple précision.

1 Donnez les représentations des nombres

a. $A = X + 1$

b. $B = X + Y$

c. $C = Y + A$

Perte d'information en arithmétique flottante

Perte d'information en arithmétique flottante

A = 0 10001110 000000000000000000000000

B = 0 10011010 000000000000000000000000

au format IEEE 754 en simple précision.

$$\text{Soit } A = 2^{142-127} \times 1,0 = 2^{15}$$

$$B = 2^{154-127} \times 1,0 = 2^{27}$$

(a) $C = A + 1$

[illegible]

(b) $D = A + B$

$$\begin{array}{rcl} A = 1,0 \times 2^{15} & = & 0, 0000000000000001 \times 2^{27} \\ B = 1,0 \times 2^{27} & = & 1, 0000000000000000 \times 2^{27} \\ \hline D & = & 1, 0000000000000001 \times 2^{27} \end{array}$$

$$(c) E = B + C$$

$$C = 1,00000000000000001 \times 2^{15}$$

$$B = 1,0 \times 2^{27}$$

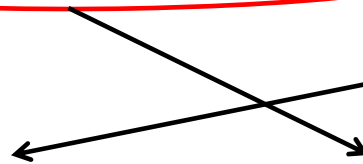
$$C = 0,00000 \ 00000 \ 00100 \ 00000 \ 00000 \ 001 \times 2^{27}$$

$$B = 1,00000 \ 00000 \ 00000 \ 00000 \ 00000 \ 000 \times 2^{27}$$

$$E = 1,00000 \ 00000 \ 00100 \ 00000 \ 00000 \ 001 \times 2^{27}$$

$$E = 1,00000\ 00000\ 00100\ 00000\ 00000\ 001 \times 2^{27}$$

$$1,00000\ 00000\ 00100\ 00000\ 000\ 000\ 001 \times 2^{27}$$



Pour le codage en IEEE 754, on ne garde que 23 bits pour la mantisse , on va donc perdre les 5 derniers bits, soit 00001

D'où : $C+B \rightarrow B$, la valeur de B très grande absorbe C .

Un premier programme

```
A = 2.000000
A = 4.000000
A = 8.000000
A = 16.000000
A = 32.000000
A = 64.000000
A = 128.000000
A = 256.000000
A = 512.000000
A = 1024.000000
A = 2048.000000
A = 4096.000000
...
A = 576460752303423490.000000
A = 1152921504606847000.000000
A = 2305843009213694000.000000
A = 4611686018427387900.000000
A = 9223372036854775800.000000
A = 18446744073709552000.000000
B = 2.000000
```

```
A=1;
while (((A + 1) - A) - 1) == 0)
{ A=2*A;
  printf ("A = %f\n",A);
}
B=1;
while (((A + B) - A) - B) != 0)
{ B=B+1;
  printf ("B = %f\n",B);
}
```

RAPPELS : float a,b;

$$\begin{array}{r}
 1,0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 2^{23}\ A \\
 +\ 0,0000\ 0000\ 0000\ 0000\ 0000\ 0001\ 2^{23}\ 1 \\
 \hline
 1,0000\ 0000\ 0000\ 0000\ 0000\ 0001\ 2^{23}\ A
 \end{array}$$

d'où $((A + 1) - A) - 1 = -1$

$$\begin{array}{r}
 +\ 0,0000\ 0000\ 0000\ 0000\ 0000\ 0010\ 2 \\
 1,0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 2^{23}\ A \\
 \hline
 1,0000\ 0000\ 0000\ 0000\ 0000\ 0011\ A+2
 \end{array}$$

d'où $((A + B) - A) - b = 0$

Un deuxième programme : sommes

n= 10 000

S1 = 9.787613 , S2=9.787604

n = 100 000

S1 = 12.090851 , S2=12.090152

n = 1 000 000

S1 = 14.357358 , S2=14.392652

n = 5 000 000

S1 = 15.403683 , S2=16.007854

n = 6 000 000

S1 = 15.403683 , S2=16.182493

$$S1_n = \sum_{j=1}^n 1/j \quad S2_n = \sum_{j=n}^1 1/j$$

RAPPELS : float s1,s2;