

T.D. 2 – Corrigé

Systèmes de numération flottante

Exercice 1

Donnez la représentation flottante, en **simple précision**, des nombres suivants :

1. 128

- **S = 0**
- $|128| = 128 = 1000\ 0000_2$
- $128 = (1,0)_2 \times 2^7$
M = 00...0₂ et **e = 7**
- $E = e + \text{biais} = 7 + 127 = 6 + 128$
E = 1000 0110₂
- **128 → 0 10000110 000000000000000000000000**

2. -32,75

- **S = 1**
- $0,75 \times 2 = 1,5$
 $0,5 \times 2 = 1$
- $|-32,75| = 32,75 = 10\ 0000,11_2$
- $32,75 = (1,0000011)_2 \times 2^5$
M = 00000110...0₂ et **e = 5**
- $E = e + \text{biais} = 5 + 127 = 4 + 128$
E = 1000 0100₂
- **-32,75 → 1 10000100 000001100000000000000000**

3. 18,125

- **S = 0**
- $0,125 \times 2 = 0,25$
 $0,25 \times 2 = 0,5$
 $0,5 \times 2 = 1$
- $|18,125| = 18,125 = 1\ 0010,001_2$
- $18,125 = (1,0010001)_2 \times 2^4$
M = 00100010...0₂ et **e = 4**
- $E = e + \text{biais} = 4 + 127 = 3 + 128$
E = 1000 0011₂
- **18,125 → 0 10000011 001000100000000000000000**

4. 0,0625

- $S = 0$
- $0,0625 \times 2 = 0,125$
 $0,125 \times 2 = 0,25$
 $0,25 \times 2 = 0,5$
 $0,5 \times 2 = 1$
- $|0,0625| = 0,0625 = 0,0001_2$
- $0,0625 = (1,0)_2 \times 2^{-4}$
 $M = 00...0_2$ et $e = -4$
- $E = e + \text{biais} = -4 + 127$
 $E = 0111\ 1011_2$
- $0,0625 \rightarrow 0\ 01111011\ 000000000000000000000000$

Exercice 2Donnez la représentation flottante, en **double précision**, des nombres suivants :

1. 1

- $S = 0$
- $|1| = 1 = 1_2$
- $1 = (1,0)_2 \times 2^0$
 $M = 00...0_2$ et $e = 0$
- $E = e + \text{biais} = 0 + 1023$
 $E = 011\ 1111\ 1111_2$
- $1 \rightarrow 0\ 011111111111\ 00.....0$

2. -64

- $S = 1$
- $|-64| = 64 = 100\ 0000_2$
- $64 = (1,0)_2 \times 2^6$
 $M = 00...0_2$ et $e = 6$
- $E = e + \text{biais} = 6 + 1023 = 5 + 1024$
 $E = 100\ 0000\ 0101_2$
- $-64 \rightarrow 1\ 10000000101\ 00.....0$

3. 12,06640625

- **S = 0**
- $0,06640625 \times 2 = 0,1328125$
 $0,1328125 \times 2 = 0,265625$
 $0,265625 \times 2 = 0,53125$
 $0,53125 \times 2 = 1,0625$
 $0,0625 \times 2 = 0,125$
 $0,125 \times 2 = 0,25$
 $0,25 \times 2 = 0,5$
 $0,5 \times 2 = 1$
- $|12,06640625| = 12,06640625 = 1100,00010001_2$
- $12,06640625 = (1,10000010001)_2 \times 2^3$
M = 100000100010...0₂ et $e = 3$
- $E = e + \text{biais} = 3 + 1023 = 2 + 1024$
E = 100 0000 0010₂
- **12,06640625 → 0 10000000010 100000100010.....0**

4. 0,2734375

- **S = 0**
- $0,2734375 \times 2 = 0,546875$
 $0,546875 \times 2 = 1,09375$
 $0,09375 \times 2 = 0,1875$
 $0,1875 \times 2 = 0,375$
 $0,375 \times 2 = 0,75$
 $0,75 \times 2 = 1,5$
 $0,5 \times 2 = 1$
- $|0,2734375| = 0,2734375 = 0,0100011_2$
- $0,2734375 = (1,00011)_2 \times 2^{-2}$
M = 000110...0₂ et $e = -2$
- $E = e + \text{biais} = -2 + 1023$
E = 011 1111 1101₂
- **0,2734375 → 0 01111111101 000110.....0**

Exercice 3

Donnez la représentation décimale des nombres codés en **simple précision** suivants :

1. $1011\ 1101\ 0100\ 0000\ 0000\ 0000\ 0000\ 0000_2$

- $S = 1 \rightarrow$ **négatif**
- $e = E - \text{biais} = 0111\ 1010_2 - 127$
 $e = 122 - 127$
 $e = -5$
- $m = (1,M)_2 = (1,1)_2$
- $-m \times 2^e = -(1,1)_2 \times 2^{-5}$
- $= -(11)_2 \times 2^{-6}$
 $= -3 \times 2^{-6} = -0,046875$

2. $0101\ 0101\ 0110\ 0000\ 0000\ 0000\ 0000\ 0000_2$

- $S = 0 \rightarrow$ **positif**
- $e = E - \text{biais} = 1010\ 1010_2 - 127$
 $e = 170 - 127$
 $e = 43$
- $m = (1,M)_2 = (1,11)_2$
- $+m \times 2^e = (1,11)_2 \times 2^{43}$
- $= +(111)_2 \times 2^{41}$
 $= 7 \times 2^{41} \approx 1,5393 \times 10^{13}$

3. $1100\ 0001\ 1111\ 0000\ 0000\ 0000\ 0000\ 0000_2$

- $S = 1 \rightarrow$ **négatif**
- $e = E - \text{biais} = 1000\ 0011_2 - 127$
 $e = 131 - 127$
 $e = 4$
- $m = (1,M)_2 = (1,111)_2$
- $-m \times 2^e = -(1,111)_2 \times 2^4$
- $= -(11110)_2 \times 2^0$
 $= -30$

4. $1111\ 1111\ 1000\ 0000\ 0000\ 0000\ 0000\ 0000_2$

- $S = 1, E = 1\dots 1$ et $M = 0\dots 0 \rightarrow -\infty$

5. $0000\ 0000\ 0100\ 0000\ 0000\ 0000\ 0000\ 0000_2$

- $E = 0\dots 0$ et $M \neq 0\dots 0 \rightarrow$ **représentation dénormalisée**
- $S = 0 \rightarrow$ **positif**
- $m = (0,M)_2 = (0,1)_2$
- $+m \times 2^{1-\text{biais}} = (0,1)_2 \times 2^{-126}$
- $= (1)_2 \times 2^{-127}$
 $= 2^{-127} \approx 5,877 \times 10^{-39}$

Exercice 4

Donnez la représentation décimale des nombres codés en **double précision** suivants :

1. $403D\ 4800\ 0000\ 0000_{16}$
 $= 0100\ 0000\ 0011\ 1101\ 0100\ 1000\ 0000\dots0$
 - $S = 0 \rightarrow$ **positif**
 - $e = E - \text{biais} = 100\ 0000\ 0011_2 - 1023 = 1027 - 1023$
 $e = 4$
 - $m = (1,M)_2 = (1,110101001)_2$
 - $+m \times 2^e = (1,110101001)_2 \times 2^4$
 - $= (1110101001)_2 \times 2^{-5}$
 $= 937 \times 2^{-5} \approx 29,28125$

2. $C040\ 0000\ 0000\ 0000_{16}$
 $= 1100\ 0000\ 0100\ 0000\dots0$
 - $S = 1 \rightarrow$ **négatif**
 - $e = E - \text{biais} = 100\ 0000\ 0100_2 - 1023 = 1028 - 1023$
 $e = 5$
 - $m = (1,M)_2 = (1,0)_2$
 - $-m \times 2^e = -(1,0)_2 \times 2^5$
 - $= -2^5 = -32$

3. $BFC0\ 0000\ 0000\ 0000_{16}$
 $= 1011\ 1111\ 1100\ 0000\dots0$
 - $S = 1 \rightarrow$ **négatif**
 - $e = E - \text{biais} = 011\ 1111\ 1100_2 - 1023 = 1020 - 1023$
 $e = -3$
 - $m = (1,M)_2 = (1,0)_2$
 - $-m \times 2^e = -(1,0)_2 \times 2^{-3}$
 - $= -2^{-3} = -0,125$

4. $8000\ 0000\ 0000\ 0000_{16}$
 $= 1000\ 0000\ 0000\ 0000\dots0$
 - $S = 1, E = 0\dots0$ et $M = 0\dots0 \rightarrow$ **-0**

5. $FFF0\ 0001\ 0000\ 0000_{16}$
 $= 1111\ 1111\ 1111\ 0000\ 0000\ 0000\ 0001\ 0000\dots0$
 - $E = 1\dots1$ et $M \neq 0\dots0 \rightarrow$ **NaN**

Exercice 5

Pour chaque question, vous traiterez le cas des codages **simples et doubles précisions** du format à **man-tisse normalisée**.

1. Déterminez, en valeur absolue, le plus petit et le plus grand nombre flottant.

- **Simple précision**

- **Minimum**

$$\text{Min}_{\text{simple}} = m_{\text{min}} \times 2^{e_{\text{min}}}$$

$$m_{\text{min}} = (1,0)_2 = 1$$

$$e_{\text{min}} = E_{\text{min}} - \text{biais} \quad \text{avec } E_{\text{min}} = 1$$

$$e_{\text{min}} = 1 - 127 = -126$$

$$\text{Min}_{\text{simple}} = 2^{-126} \approx 1,1755 \times 10^{-38}$$

- **Maximum**

$$\text{Max}_{\text{simple}} = m_{\text{max}} \times 2^{e_{\text{max}}}$$

$$m_{\text{max}} = (1, M_{\text{max}})_2 = 1 + (0, M_{\text{max}})_2 = 1 + M_{\text{max}} \times 2^{-23} \quad \text{avec } M_{\text{max}} = 2^{23} - 1$$

$$m_{\text{max}} = 1 + (2^{23} - 1) \times 2^{-23} = 1 + 1 - 2^{-23} = 2 - 2^{-23} = 2 \times (1 - 2^{-24})$$

$$e_{\text{max}} = E_{\text{max}} - \text{biais} \quad \text{avec } E_{\text{max}} = (2^8 - 1) - 1 = 254$$

$$e_{\text{max}} = 254 - 127 = 127$$

$$\text{Max}_{\text{simple}} = 2 \times (1 - 2^{-24}) \times 2^{127}$$

$$\text{Max}_{\text{simple}} = (1 - 2^{-24}) \times 2^{128} \approx 3,4028 \times 10^{38}$$

- **Double précision**

- **Minimum**

$$\text{Min}_{\text{double}} = m_{\text{min}} \times 2^{e_{\text{min}}}$$

$$m_{\text{min}} = (1,0)_2 = 1$$

$$e_{\text{min}} = E_{\text{min}} - \text{biais} \quad \text{avec } E_{\text{min}} = 1$$

$$e_{\text{min}} = 1 - 1023 = -1022$$

$$\text{Min}_{\text{double}} = 2^{-1022} \approx 2,2251 \times 10^{-308}$$

- **Maximum**

$$\text{Max}_{\text{double}} = m_{\text{max}} \times 2^{e_{\text{max}}}$$

$$m_{\text{max}} = (1, M_{\text{max}})_2 = 1 + (0, M_{\text{max}})_2 = 1 + M_{\text{max}} \times 2^{-52} \quad \text{avec } M_{\text{max}} = 2^{52} - 1$$

$$m_{\text{max}} = 1 + (2^{52} - 1) \times 2^{-52} = 1 + 1 - 2^{-52} = 2 - 2^{-52} = 2 \times (1 - 2^{-53})$$

$$e_{\text{max}} = E_{\text{max}} - \text{biais} \quad \text{avec } E_{\text{max}} = (2^{11} - 1) - 1 = 2046$$

$$e_{\text{max}} = 2046 - 1023 = 1023$$

$$\text{Max}_{\text{double}} = 2 \times (1 - 2^{-53}) \times 2^{1023}$$

$$\text{Max}_{\text{double}} = (1 - 2^{-53}) \times 2^{1024} \approx 1,7977 \times 10^{308}$$

- Ce qui donne :

$$f1 < 2^{28}$$

$$10^n < 2^{28}$$

$$n < \text{Log}(2^{28})$$

$$n < 8,42$$

$$\mathbf{n_{\max} = 8}$$

3. Même question si les variables $f1, f2, f3$ et r sont déclarées en double précision.

Avec un raisonnement identique à celui du codage en simple précision, on obtient :

$$f1 < 2^{5+52}$$

$$10^n < 2^{57}$$

$$n < \text{Log}(2^{57})$$

$$n < 17,15$$

$$\mathbf{n_{\max} = 17}$$

Il ne faut jamais sous-estimer les risques d'erreur liés à la manipulation de variables entières ou flottantes. Des débordements ou des problèmes de précision peuvent survenir à tout moment et déjouer la vigilance de n'importe quel développeur, même expérimenté.

À propos de la fusée Ariane

Voici un exemple célèbre d'erreur de programmation minime aux conséquences énormes. Lors de son tout premier vol le 4 juin 1996, la fusée Ariane 5 a explosé quarante secondes seulement après son décollage de la base de Kourou en Guyane. La perte financière fut estimée à environ 500 millions de dollars. Le CNES (Centre National d'Études Spatiales) et l'ESA (*European Space Agency*) ont immédiatement lancé une enquête. Un comité d'experts internationaux fut réuni et un rapport sans équivoque fut livré un mois plus tard : l'explosion était due à une erreur de logiciel.

En cause, les SRI ou Systèmes de Référence Inertiels. Cette partie du logiciel, qui provenait du lanceur Ariane 4, n'avait pas été adaptée à la plus grande vitesse horizontale d'Ariane 5. Du coup, lors d'une conversion d'un nombre flottant sur 64 bits, contenant la vitesse horizontale, en un entier sur 16 bits, l'opération a provoqué un débordement et une exception a été générée. Malheureusement, aucune routine de traitement de cette exception n'ayant été prévue, c'est le traitement d'exception par défaut qui fut exécuté et le programme tout entier termina son exécution.

Depuis, les concepteurs du logiciel d'Ariane ont mis en place un plan de programmation défensive qui est reconnu comme une référence en la matière.

Jean-Christophe Arnulfo, *Métier Développeur*, Paris, Dunod, 2003

Exercice 7

Sachant que votre compilateur C utilise la norme IEEE 754 pour la gestion des flottants, donnez une fonction en langage C, **de quelques lignes seulement**, permettant de visualiser sous forme hexadécimale la représentation IEEE 754 d'un nombre flottant, simple précision, passé en paramètre.

Le principe de base consiste à trouver l'emplacement où le nombre flottant est stocké en mémoire. Il faut ensuite accéder à cette valeur, en la considérant comme un nombre entier, et l'afficher au format hexadécimal.

Ceci peut être réalisé soit à l'aide de pointeurs, soit à l'aide du mot clé `union` :

- **Solution à l'aide de pointeurs :**

```
void Convert(float fv)
{
    // Déclare un pointeur sur un entier 32 bits.
    long* pl;

    // Convertit le pointeur flottant vers un pointeur entier.
    // Le pointeur entier pointe la même adresse mémoire que le pointeur flottant.
    pl = (long*)&fv;

    // Affiche le résultat au format hexadécimal sur 8 chiffres.
    printf("Code hexadécimal IEEE 754 de %f : %08x\n", fv, *pl);
}
```

- **Solution à l'aide du mot clé `union` :**

```
void Convert(float fv)
{
    // Déclare un flottant et un entier dans le même espace mémoire.
    union
    {
        float f;
        long l;
    };

    // Copie la valeur à convertir dans le flottant de l'union.
    f = fv;

    // Affiche le résultat au format hexadécimal sur 8 chiffres.
    printf("Code hexadécimal IEEE 754 de %f : %08x\n", f, l);
}
```