

# IEEE 754 : au-delà du réel

Journées académiques 2014

Guillaume CONNAN

IREM de Nantes

Dernière mise à jour : 23 avril 2014 à 21:47



# Sommaire

- 1 Préambule
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



# Sommaire

## 1 Préambule

### 2 La norme IEEE 754

- Différents formats
- Les normaux, les sous-normaux et les paranormaux

### 3 Algèbre des nombres VF

### 4 Réels, arrondis et flottants

### 5 Que la force de l'erreur soit avec vous

- Atelier Padawan # 1 : majorer l'erreur
- Atelier Padawan # 2 : somme de flottants
- Atelier Padawan # 3 : somme compensée de flottants quelconques

### 6 Des sujets de Bac maladroits...

### 7 Lectures recommandées pour aller plus loin













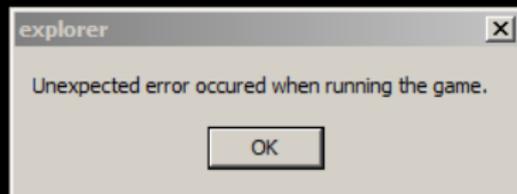
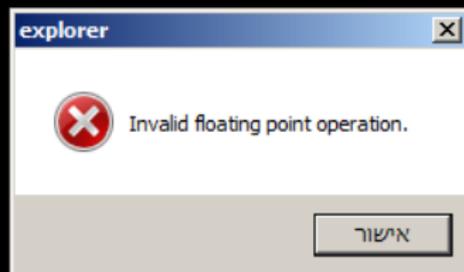
Microsoft Excel - Classeur1

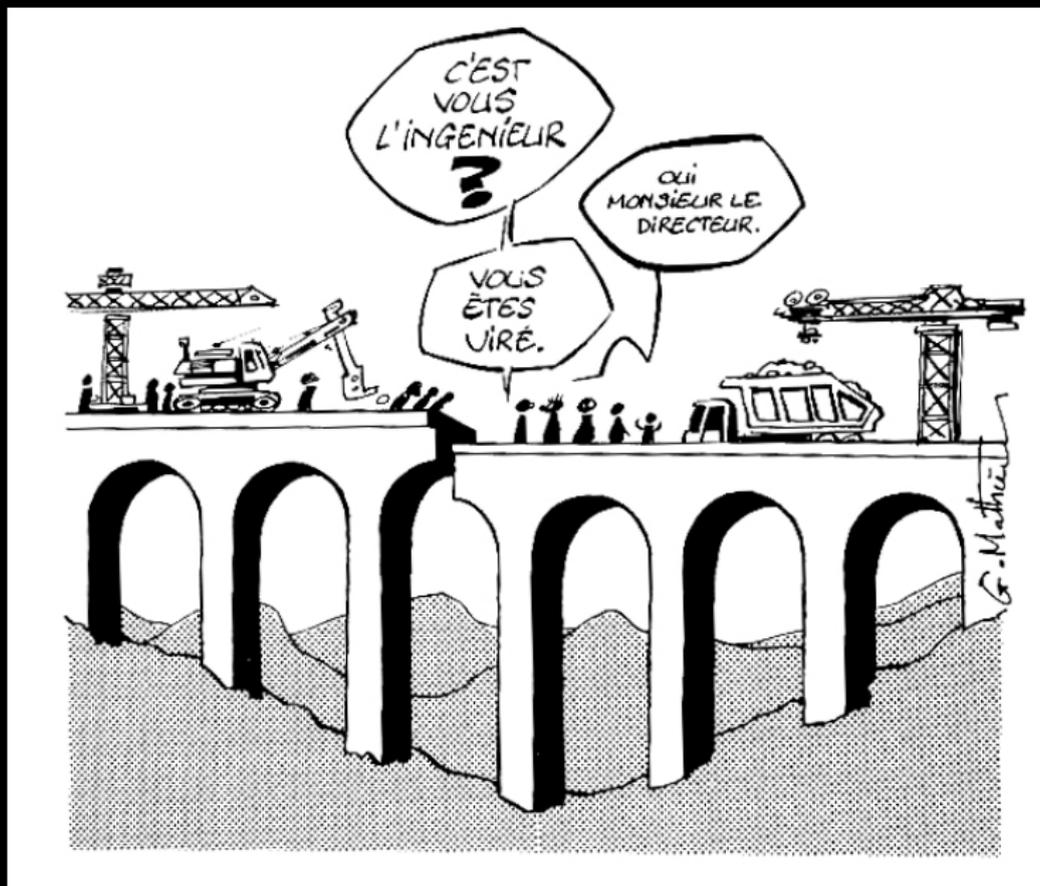
Fichier Edition Affichage Insertion Format Outils Do

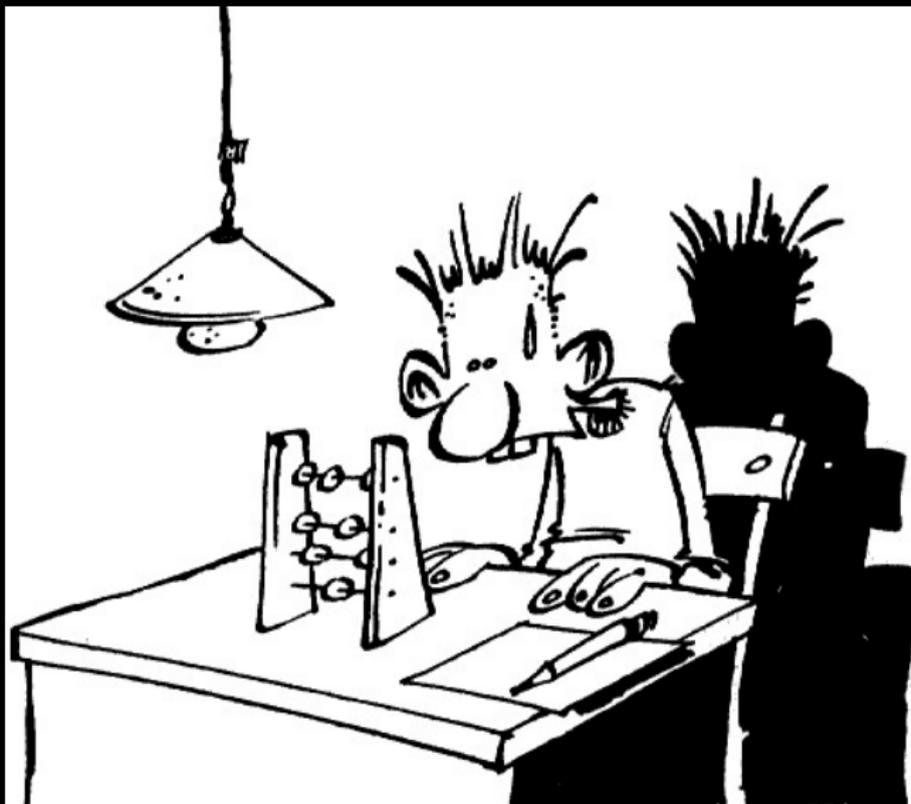
Arial 10 G I S

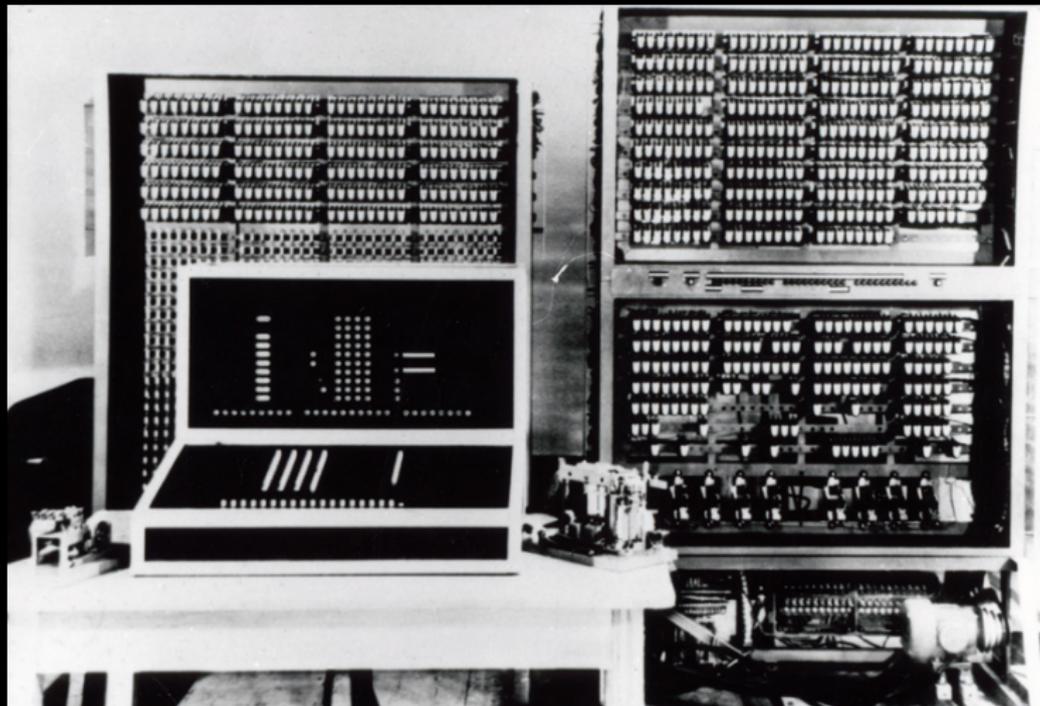
C15 fx

	A	B
1	$B1 = 4/3$	1,33333333333333000000
2	$B2 = B1 - 1$	0,33333333333333300000
3	$B3 = B2 * 3$	1,00000000000000000000
4	$B4 = B3 - 1$	0,00000000000000000000
5	$B5 = B4 * 2^{52}$	0,00000000000000000000
6	$(4/3 - 1) * 3 - 1$	0,00000000000000000000
7	$((4/3 - 1) * 3 - 1)$	-0,000000000000000022204
8	$((4/3 - 1) * 3 - 1) * 2^{52}$	-1,00000000000000000000











- char
- int



- char
- int  $2147483647 + 1 = -2147483648...$



- char
- int  $2147483647 + 1 = -2147483648...$



- char
- `int 2147483647 + 1 = -2147483648...`
- short



- char
- int  $2147483647 + 1 = -2147483648...$
- short long



- char
- int  $2147483647 + 1 = -2147483648...$
- short long



- char
- `int 2147483647 + 1 = -2147483648...`
- short long

• float

• double



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- $\pi$



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1



- char
- int  $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1
- ...



```
In [1]: 3 * 0.1
```

- 
- 
- 



```
In [1]: 3 * 0.1
```

```
Out[1]: 0.30000000000000004
```

- 
- 



```
In [1]: 3 * 0.1
```

```
Out[1]: 0.30000000000000004
```

```
In [2]: sum([0.1 for k in range(10000000)])
```

```
.
```



```
In [1]: 3 * 0.1
```

```
Out[1]: 0.30000000000000004
```

```
In [2]: sum([0.1 for k in range(10000000)])
```

```
Out[2]: 999999.9998389754
```



# LES NOMBRES RÉELS N'EXISTENT PAS.

TOUT CE QUE VOUS AVEZ VU  
AU LYCÉE N'EST  
QU'ILLUSION !



LES NOMBRES RÉELS  
N'EXISTENT PAS.

TOUT CE QUE VOUS AVEZ VU  
AU LYCÉE N'EST  
QU'ILLUSION !



# Sommaire

1 Préambule

2 **La norme IEEE 754**

- Différents formats
- Les normaux, les sous-normaux et les paranormaux

3 Algèbre des nombres VF

4 Réels, arrondis et flottants

5 Que la force de l'erreur soit avec vous

- Atelier Padawan # 1 : majorer l'erreur
- Atelier Padawan # 2 : somme de flottants
- Atelier Padawan # 3 : somme compensée de flottants quelconques

6 Des sujets de Bac maladroits...

7 Lectures recommandées pour aller plus loin



# Sommaire

- 1 Préambule
- 2 **La norme IEEE 754**
  - Différents formats
    - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin





$$v = (-1)^s \times m \times 2^E$$

• binariser



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)

► *binary64*



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)
- *binary64* ( $\#E, \#m$ ) = (11, 53).



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)
- *binary64* ( $\#E, \#m$ ) = (11, 53).



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)
- *binary64* ( $\#E, \#m$ ) = (11, 53).



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)
- *binary64* ( $\#E, \#m$ ) = (11, 53).
- *toy7* ( $\#E, \#m$ ) = (3, 4).



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)
- *binary64* ( $\#E, \#m$ ) = (11, 53).
- *toy7* ( $\#E, \#m$ ) = (3, 4).



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)
- *binary64* ( $\#E, \#m$ ) = (11, 53).
- *toy7* ( $\#E, \#m$ ) = (3, 4).



$$v = (-1)^s \times m \times 2^E$$

- *binary32* ( $\#E, \#m$ ) = (8, 24)
- *binary64* ( $\#E, \#m$ ) = (11, 53).
- *toy7* ( $\#E, \#m$ ) = (3, 4).



On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur  $\#E$  bits on peut coder  $2^{\#E}$  nombres.

La première moitié va donc de 0 à  $2^{\#E-1} - 1$ .

Elle correspond aux exposants réels de  $E_{\min}$  jusqu'à 0.

La deuxième de  $2^{\#E-1}$  à  $2^{\#E} - 1$ .

Elle correspond aux exposants de 1 jusqu'à  $E_{\max}$ .

Il suffit donc de translater les exposants réels de  $2^{\#E-1}$  à 0.

$$\text{Exposant} = \text{Exposant} - 2^{\#E-1}$$



On peut également gagner de la place en ne stockant pas le signe de l'exposant.  
Sur  $\#E$  bits on peut coder  $2^{\#E}$  nombres.

La première moitié va donc de 0 à  $2^{\#E-1} - 1$ .

Elle correspond aux exposants réels de  $E_{\min}$  jusqu'à 0.

La deuxième de  $2^{\#E-1}$  à  $2^{\#E} - 1$ .

Elle correspond aux exposants de 1 jusqu'à  $E_{\max}$ .

Il suffit donc de translater les exposants réels de  $E_{\min}$  à 0.

$$E_{\text{stocké}} = E_{\text{réel}} + E_{\text{biais}}$$



On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur  $\#E$  bits on peut coder  $2^{\#E}$  nombres.

La première moitié va donc de 0 à  $2^{\#E-1} - 1$ .

Elle correspond aux exposants réels de  $E_{\min}$  jusqu'à 0.

La deuxième de  $2^{\#E-1}$  à  $2^{\#E} - 1$ .

Elle correspond aux exposants de 1 jusqu'à  $E_{\max}$ .

Il suffit donc de translater les exposants réels de  $E_{\min}$  à 0.

Il faut donc coder  $E_{\min} + 2^{\#E-1}$  pour l'exposant 0.

Il faut donc coder  $E_{\max} + 2^{\#E-1} - 1$  pour l'exposant 1.



On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur  $\#E$  bits on peut coder  $2^{\#E}$  nombres.

La première moitié va donc de 0 à  $2^{\#E-1} - 1$ .

Elle correspond aux exposants réels de  $E_{\min}$  jusqu'à 0.

La deuxième de  $2^{\#E-1}$  à  $2^{\#E} - 1$ .

Elle correspond aux exposants de 1 jusqu'à  $E_{\max}$ .

Il suffit donc de translater les exposants réels de  $2^{\#E-1} - 1$ ...

$$E_{\text{min}} + 2^{\#E-1} - 1$$



On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur  $\#E$  bits on peut coder  $2^{\#E}$  nombres.

La première moitié va donc de 0 à  $2^{\#E-1} - 1$ .

Elle correspond aux exposants réels de  $E_{\min}$  jusqu'à 0.

La deuxième de  $2^{\#E-1}$  à  $2^{\#E} - 1$ .

Elle correspond aux exposants de 1 jusqu'à  $E_{\max}$ .

Il suffit donc de tradater les exposants réels de  $2^{\#E-1} - 1 \dots$

$$e_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$



On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur  $\#E$  bits on peut coder  $2^{\#E}$  nombres.

La première moitié va donc de 0 à  $2^{\#E-1} - 1$ .

Elle correspond aux exposants réels de  $E_{\min}$  jusqu'à 0.

La deuxième de  $2^{\#E-1}$  à  $2^{\#E} - 1$ .

Elle correspond aux exposants de 1 jusqu'à  $E_{\max}$ .

Il suffit donc de tradater les exposants réels de  $2^{\#E-1} - 1 \dots$

$$e_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$



On peut également gagner de la place en ne stockant pas le signe de l'exposant.  
Sur  $\#E$  bits on peut coder  $2^{\#E}$  nombres.  
La première moitié va donc de 0 à  $2^{\#E-1} - 1$ .  
Elle correspond aux exposants réels de  $E_{\min}$  jusqu'à 0.  
La deuxième de  $2^{\#E-1}$  à  $2^{\#E} - 1$ .  
Elle correspond aux exposants de 1 jusqu'à  $E_{\max}$ .  
Il suffit donc de tradater les exposants réels de  $2^{\#E-1} - 1 \dots$

$$e_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$



On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur  $\#E$  bits on peut coder  $2^{\#E}$  nombres.

La première moitié va donc de 0 à  $2^{\#E-1} - 1$ .

Elle correspond aux exposants réels de  $E_{\min}$  jusqu'à 0.

La deuxième de  $2^{\#E-1}$  à  $2^{\#E} - 1$ .

Elle correspond aux exposants de 1 jusqu'à  $E_{\max}$ .

Il suffit donc de tradater les exposants réels de  $2^{\#E-1} - 1 \dots$

$$e_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$



# 0,75 en toy7

0,75<sub>10</sub> en toy7.

$$0,75_{10} = \frac{3_{10}}{2^2_{10}} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 - e = -1 + 2 = 10_2$$

$$\left[ \begin{array}{c|c|c} s(1 \text{ bit}) & e(3 \text{ bits}) & f(3 \text{ bits}) \\ \hline 0 & 100 & 100 \end{array} \right]$$



# 0,75 en toy7

0,75<sub>10</sub> en toy7.

$$0,75_{10} = \frac{3_{10}}{2^2_{10}} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \rightarrow e = -1 + 3 = 2 = 10_2$$

signe	exposant	mantisse
0	10	1100



## 0,75 en toy7

0,75<sub>10</sub> en toy7.

$$0,75_{10} = \frac{3_{10}}{2^2_{10}} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \rightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0



## 0,75 en toy7

0,75<sub>10</sub> en toy7.

$$0,75_{10} = \frac{3_{10}}{2^2_{10}} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \rightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0



## 0,75 en toy7

0,75<sub>10</sub> en toy7.

$$0,75_{10} = \frac{3_{10}}{2^2_{10}} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0



## 0,75 en toy7

0,75<sub>10</sub> en toy7.

$$0,75_{10} = \frac{3_{10}}{2^2_{10}} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0



## 0,75 en toy7

$0,75_{10}$  en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0



## 0,75 en toy7

$0,75_{10}$  en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0



# 0,75 en $\text{toy7}$ : méthode générale

Comment écrire en base 2 la partie fractionnaire d'un nombre en base 10 ?

$$0,75 = b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$2 \times 0,75 = b_1 + b_2 \times 2^{-1} + \dots$$

On voit poindre un bel algo...



# 0,75 en $\text{toy7}$ : méthode générale

Comment écrire en base 2 la partie fractionnaire d'un nombre en base 10 ?

$$0,75 = b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$2 \times 0,75 = b_1 + b_2 \times 2^{-1} + \dots$$

On voit poindre un bel algo...



# 0,75 en $\text{toy7}$ : méthode générale

Comment écrire en base 2 la partie fractionnaire d'un nombre en base 10 ?

$$0,75 = b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$2 \times 0,75 = b_1 + b_2 \times 2^{-1} + \dots$$

On voit poindre un bel algo...



# 0,75 en $\text{toy7}$ : méthode générale

Comment écrire en base 2 la partie fractionnaire d'un nombre en base 10 ?

$$0,75 = b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$2 \times 0,75 = b_1 + b_2 \times 2^{-1} + \dots$$

On voit poindre un bel algo...



## Exercice 1

*Pourquoi le Patriot a raté son Scud de 600m ?*

*Représentez  $1/10$  avec une mantisse de 24 bits et tronquez le résultat.*

*Calculez l'erreur en seconde puis après 100 heures d'utilisation. Sachant qu'un Scud vole à  $1676\text{ms}^{-1}$ , quelle est environ l'erreur commise en mètres ?*

*Et si on avait arrondi au plus proche ?*



# Sommaire

- 1 Préambule
- 2 **La norme IEEE 754**
  - Différents formats
  - **Les normaux, les sous-normaux et les paranormaux**
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



## QUELQUES HUMANOÏDES RARES

**CENTRALIEN****NORMALIEN****ALIEN**

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7: } e = E + 1, e_{\min} = 000, e_{\max} = 111 = 000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1}) - 1 = 0 - 2^{(\#E)-1} - 1 = -2^{(\#E)-1} - 1$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1}) - 1 = (2^{\#E} - 1) - 2^{(\#E)-1} - 1 = 2^{(\#E)-1} - 2$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

toy7,  $\#E = 3$ ,  $e_{\min} = 000$ ,  $e_{\max} = 111 = 1000 - 1$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1}) - 1 = 0 - 2^{3-1} - 1 = -3$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1}) - 1 = (1000 - 1) - 2^{3-1} - 1 = 7 - 4 - 1 = 2$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

toy7,  $\#E = 3$ ,  $e_{\min} = 000$ ,  $e_{\max} = 111 = 2000 - 1$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1}) - 1 = 0 - 1 = -1$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1}) - 1 = (2000 - 1) - 1 = 1998$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1}) - 1 = 0 - 1 = -1$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1}) - 1 = (111 - 1000) - 1 = -1000 - 1 = -1001$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, \#E} = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = e_{\min} - 1 = (2^{\#E} - 1) - 1 = 2^{\#E} - 2$$

$$E_{\max} = e_{\max} - 1 = (2^{\#E} - 1) - 1 = 2^{\#E} - 2 = 2^{\#E} - 1 - 1 = e_{\max} - 1$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{(\#E)-1} + 1 = 1 - 2^{(\#E)-1}$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^{\#E} - 1 - 2^{(\#E)-1} + 1) - 1 = 2^{\#E} - 1 - 2^{(\#E)-1} - 1 = 2^{\#E} - 2^{(\#E)-1} - 2$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (111 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (111 - 1) - 1 = 100 - 1 = 99$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$



$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$



- zéro sous forme normale ?

- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux



- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux...



- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux....



- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux....



- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux....



## Exercice 2

*En utilisant vos carreaux, représentez les nombres à virgule flottante normaux en  $toy7$  et rajoutez les sous-normaux.*



$s$ (1 bit)	$e$ (3 bits)			$f$ (3 bits)		
0	0	0	0	0	0	0

$s$ (1 bit)	$e$ (3 bits)			$f$ (3 bits)		
1	0	0	0	0	0	0



$s$ (1 bit)	$e$ (3 bits)			$f$ (3 bits)		
0	0	0	0	0	0	0

$s$ (1 bit)	$e$ (3 bits)			$f$ (3 bits)		
1	0	0	0	0	0	0



```
In [1]: x = 1e500
```

- 
- 
- 
- 
- 
- 
- 
- 
- 



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

- 
- 
- 
- 
- 
- 
- 



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

- 
- 
- 
- 
- 
- 



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
•  
•  
•  
•  
•
```



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

- 
- 
- 
- 



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

```
In [4]: 1/x
```

- 
- 
- 



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

```
In [4]: 1/x
```

```
Out[4]: 0.0
```

```
.
```

```
.
```



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

```
In [4]: 1/x
```

```
Out[4]: 0.0
```

```
In [5]: 4 - x
```

```
.
```



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

```
In [4]: 1/x
```

```
Out[4]: 0.0
```

```
In [5]: 4 - x
```

```
Out[5]: -inf
```



$$\lim_{x \rightarrow +\infty} x^2 - x$$

```
In [9]: x**2 - x
```



$$\lim_{x \rightarrow +\infty} x^2 - x$$

```
In [9]: x**2 - x
```

```
.
```



$$\lim_{x \rightarrow +\infty} x^2 - x$$

```
In [9]: x**2 - x
```

```
Out[9]: nan
```



$$\lim_{x \rightarrow +\infty} x^2 - x$$

```
In [9]: x**2 - x
```

```
Out[9]: nan
```





```
In [10]: x
```

```
Out[10]: inf
```

```
In [11]: x - x
```

```
Out[11]: nan
```

```
In [12]: x / x
```

```
Out[12]: nan
```

```
In [13]: x - x == x - x
```

```
Out[13]: False
```

```
In [14]: x**2 - x
```

```
Out[14]: nan
```

```
In [15]: x * (x - 1)
```

```
Out[15]: inf
```

```
In [16]: x * (x - 1) == x**2 - x
```

```
Out[16]: False
```

## Exercice 3

- 1 Comment expliquer les résultats suivants :

```
*Main> let f(x) = x^2 / sqrt(x^3 + 1)
*Main> f(1e100)
1.0e50
*Main> f(1e150)
0.0
*Main> f(1e200)
NaN
```

- 2 Comment éviter le dernier NaN ?



```
In [16]: f = lambda x: x**2 / sqrt(x**3 + 1)
```

```
In [17]: f(1e100)
```

```
Out[17]: 1e+50
```

```
In [18]: f(1e150)
```

```
-----  
OverflowError                                Traceback (most  
  recent call last)
```

```
<ipython-input-18-79775d85f31a> in <module>()  
----> 1 f(1e150)
```

```
<ipython-input-16-eed0b81ef932> in <lambda>(x)
```

```
----> 1 f = lambda x: x**2 / sqrt(x**3 + 1)
```

```
OverflowError: (34, 'Numerical result out of range')
```

```
In [19]: f(1e200)
```

```
In [53]: x = 1e-500
```

```
In [54]: x
```

```
Out[54]: 0.0
```

```
In [55]: x / x
```

```
-----  
ZeroDivisionError                                Traceback (most  
  recent call last)
```

```
<ipython-input-55-fd52a7f8b5f1> in <module>()  
----> 1 x / x
```

```
ZeroDivisionError: float division by zero
```



```
In [56]: sqrt(-1.0)
```

```
-----  
ValueError                                Traceback (most  
  recent call last)  
<ipython-input-56-d1c09f21b443> in <module>()  
----> 1 sqrt(-1.0)  
  
ValueError: math domain error
```

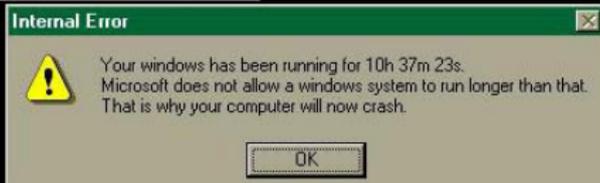
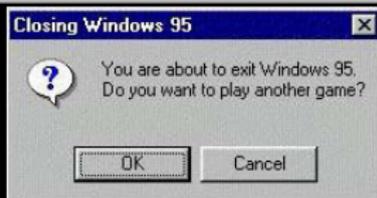
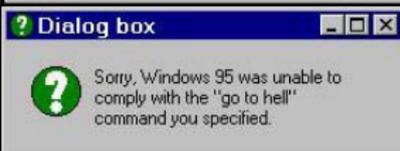
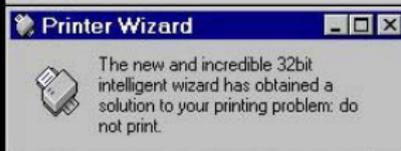
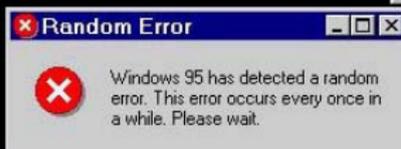
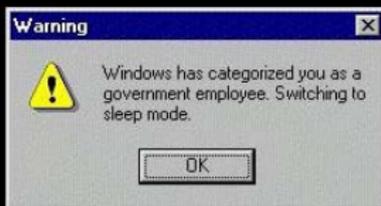


```
*Main> let x = 1e500
*Main> x - x
NaN
*Main> x / x
NaN
*Main> sqrt(-1)
NaN
*Main> 0 / 0
NaN
```



# USS Yorktown 1998





# Air France 447 - 1<sup>er</sup> juin 2009



$s$ (1 bit)	$e$ (3 bits)			$f$ (3 bits)		
1	1	1	1	0	1	0



# Arbre de lecture

$s$	$e$	$f$
-----	-----	-----

$e = 0 \dots 0$	$f = 0$	$\rightarrow v = (-1)^s \times 0.0$
$e = 0 \dots 0$	$f \neq 0$	$\rightarrow v = (-1)^s \times 0.f \times 2^{1-E_{\max}}$
$e = 1 \dots 1$	$f = 0$	$\rightarrow v = (-1)^s \times \infty$
$e = 1 \dots 1$	$f \neq 0$	$\rightarrow \text{NaN}$
$0 \dots 0 < e < 1 \dots 1$	$f \neq 0$	$\rightarrow v = (-1)^s \times 1.f \times 2^{e-E_{\max}}$



# Arbre de lecture



$$e = 0 \dots 0 \quad f = 0 \quad \rightarrow \quad v = (-1)^s \times 0.0$$

$$e = 0 \dots 0 \quad f \neq 0 \quad \rightarrow \quad v = (-1)^s \times 0.f \times 2^{1-E_{\max}}$$

$$e = 1 \dots 1 \quad f = 0 \quad \rightarrow \quad v = (-1)^s \times \infty$$

$$e = 1 \dots 1 \quad f \neq 0 \quad \rightarrow \quad \text{NaN}$$

$$0 \dots 0 < e < 1 \dots 1 \quad f \neq 0 \quad \rightarrow \quad v = (-1)^s \times 1.f \times 2^{e-E_{\max}}$$



# SuccessEUR

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 53$

*float16*

SuccessEUR: IEEE 754-2008, *float16* and *float16b*



# SuccessEUR

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 53$

*single*,  $n = 24$

*double*,  $n = 53$



# Successesseur

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 52$

*nextUp*( $x$ )

successor( $x$ ) =  $x + \text{nextUp}(x)$



# SuccessEUR

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 52$



# SuccessEUR

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{norm}(v) = \lfloor \log_2(|M(v)|) \rfloor + 1 - 2^{E(v)-n}$$



# SuccessEUR

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$



# SuccessEUR

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$



# SuccessEUR

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$



# SuccessEUR

$$n = \#f$$

*toy7*,  $n = 3$

*binary32*,  $n = 23$

*binary64*,  $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$



# Tableau récapitulatif

- On notera  $\varepsilon_m$  l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera  $\lambda$  le plus petit VF normal positif
- On notera  $\mu$  le plus petit VF sous-normal positif
- On notera  $\Omega$  le plus grand VF normal.

Quelle relation existe-t-il entre  $\mu$ ,  $\lambda$  et  $\Omega$  ?



# Tableau récapitulatif

- On notera  $\varepsilon_m$  l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera  $\lambda$  le plus petit VF normal positif
- On notera  $\mu$  le plus petit VF sous-normal positif
- On notera  $\Omega$  le plus grand VF normal.

Quelle relation existe-t-il entre  $\mu$  et  $\varepsilon_m$  ?



# Tableau récapitulatif

- On notera  $\varepsilon_m$  l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera  $\lambda$  le plus petit VF normal positif
- On notera  $\mu$  le plus petit VF sous-normal positif
- On notera  $\Omega$  le plus grand VF normal.

Quelle relation existe-t-il entre  $\mu$ ,  $\varepsilon_m$  et  $\lambda$  ?



# Tableau récapitulatif

- On notera  $\varepsilon_m$  l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera  $\lambda$  le plus petit VF normal positif
- On notera  $\mu$  le plus petit VF sous-normal positif
- On notera  $\Omega$  le plus grand VF normal.

Quelle relation existe-t-il entre  $\mu$ ,  $\varepsilon_m$  et  $\lambda$  ?



# Tableau récapitulatif

- On notera  $\varepsilon_m$  l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera  $\lambda$  le plus petit VF normal positif
- On notera  $\mu$  le plus petit VF sous-normal positif
- On notera  $\Omega$  le plus grand VF normal.

Quelle relation existe-t-il entre  $\mu$ ,  $\varepsilon_m$  et  $\lambda$  ?



# Tableau récapitulatif

- On notera  $\varepsilon_m$  l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera  $\lambda$  le plus petit VF normal positif
- On notera  $\mu$  le plus petit VF sous-normal positif
- On notera  $\Omega$  le plus grand VF normal.

Quelle relation existe-t-il entre  $\mu$ ,  $\varepsilon_m$  et  $\lambda$  ?



Format	#E	#f	$E_{\min}$	$E_{\max}$	$\varepsilon_m$	$\lambda$	$\mu$	$\Omega$
<i>toy7</i>	3	3	-2	3	$2^{-3} = 1/8$	$2^{-2} = 1/4$	$2^{-5} = 1/32$	$1,111 \times 2^3 = 15$
<i>bin32</i>	8	23	-126	127	$2^{-23}$ $\approx 1,2 \times 10^{-7}$	$2^{-126}$ $\approx 1,2 \times 10^{-38}$	$2^{-126-23}$ $\approx 1,4 \times 10^{-45}$	$(2^{24} - 1) \times 2^{127-23}$ $\approx 3,4 \times 10^{38}$
<i>bin64</i>	11	52	-1022	1023	$2^{-52}$ $\approx 2,2 \times 10^{-16}$	$2^{-1022}$ $\approx 2,2 \times 10^{-308}$	$2^{-1074}$ $\approx 5 \times 10^{-324}$	$(2^{53} - 1)2^{1023-52}$ $\approx 1,8 \times 10^{308}$



```
In [1]: 1 + 1e-16 == 1
```

- 
- 
- 
- 
- 
- 
- 
- 
- 
- 
- 
- 
- 
- 
- 



```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

- 
- 
- 
- 
- 
- 
- 
- 
- 
- 
- 



```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

- 
- 
- 
- 
- 
- 
- 
- 
- 
- 



```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

- 
- 
- 
- 
- 
- 
- 



```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

- 
- 
- 
- 
- 



```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

```
In [4]: x - 1e-16
```

```
.
```

```
.
```

```
.
```



```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

```
In [4]: x - 1e-16
```

```
Out[4]: 0.9999999999999999
```

```
.
```

```
.
```



```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

```
In [4]: x - 1e-16
```

```
Out[4]: 0.9999999999999999
```

```
In [5]: 1e-16 + 1e-18
```

```
.
```



```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

```
In [4]: x - 1e-16
```

```
Out[4]: 0.9999999999999999
```

```
In [5]: 1e-16 + 1e-18
```

```
Out[5]: 1.01e-16
```



# Sommaire

- 1 Préambule
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



$\forall b$  $\overline{\forall b}$ 

$$\sqrt{b}$$

$$\overline{\sqrt{b}}$$



# Comparaison

Il est très simple et rapide de comparer deux VF : comment la machine procède-t-elle ? Quel est l'avantage de ce stockage des VF ?



# Addition

- 1 on commence par ramener les deux nombres au même exposant, en l'occurrence le plus grand des deux ;
- 2 on ajoute les deux mantisses *complètes* en tenant compte du signe ;
- 3 on renormalise le nombre obtenu.



# Addition

- 1 on commence par ramener les deux nombres au même exposant, en l'occurrence le plus grand des deux ;
- 2 on ajoute les deux mantisses *complètes* en tenant compte du signe ;
- 3 on renormalise le nombre obtenu.



# Addition

- 1 on commence par ramener les deux nombres au même exposant, en l'occurrence le plus grand des deux ;
- 2 on ajoute les deux mantisses *complètes* en tenant compte du signe ;
- 3 on renormalise le nombre obtenu.



# Addition

En *toy7* :  $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 : 

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 : 

0	0	0	1	1	1	0
---	---	---	---	---	---	---

1,100  
0,0011  
-----  
1,1011





# Addition

En *toy7* :  $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 : 

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 : 

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111

```



# Addition

En *toy7* :  $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 : 

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 : 

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111
  
```



# Addition

En *toy7* :  $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 : 

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 : 

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111

```



## Addition

En *toy7* :  $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 : 

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 : 

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111

```



# Addition

En *toy7* :  $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 : 

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 : 

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111
  
```



# Les arrondis

## Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers  $+\infty$  (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers  $-\infty$  (RD) qui arrondit à la valeur inférieure la plus grande ;

Le mode d'arrondi par défaut est le premier.



# Les arrondis

## Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers  $+\infty$  (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers  $-\infty$  (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.



# Les arrondis

## Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers  $+\infty$  (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers  $-\infty$  (RD) qui arrondit à la valeur inférieure la plus grande.



# Les arrondis

## Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers  $+\infty$  (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers  $-\infty$  (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.



# Les arrondis

## Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers  $+\infty$  (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers  $-\infty$  (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.



# Les arrondis

## Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers  $+\infty$  (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers  $-\infty$  (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.



# Les arrondis

## Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers  $+\infty$  (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers  $-\infty$  (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.



$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \hline
 1,10111
 \end{array}$$

$x = 1,10111$  VF en  $\text{toy7}$ ?

$$1,10100 < \underbrace{1,10111}_x < 1,11000$$

$$\frac{1,100 + 1,00 \cdot 10^6}{10^6} < x < \frac{1,110 + 10,00 \cdot 10^6}{10^6}$$

soit

$(x - 1) \cdot 10^6 \in ]0,100[$  et  $(x - 1) \cdot 10^6 \in ]0,1100[$  donc  $(x - 1) \cdot 10^6 \in ]0,100[$

$$1,1 + 0,00111 = 1,111$$



$$\begin{array}{r} 1,100 \\ 0,00111 \\ \hline 1,10111 \end{array}$$

$x = 1,10111$  VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ ulp}(x)}_{\text{succ}(v)}$$

$|x - v| = 0,11 \text{ ulp}(x)$   $|x - \text{succ}(v)| = 0,01 \text{ ulp}(x)$  donc  $\text{RN}(x) = \text{succ}(v)$

$$1,1 + 0,00111 = 1,111$$



$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \hline
 1,10111
 \end{array}$$

$x = 1,10111$  VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ ulp}(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11 \text{ ulp}(x) \quad |x - \text{succ}(v)| = 0,01 \text{ ulp}(x) \quad \text{donc } \text{round}(x) = \text{succ}(v)$$

$$1,1 + 0,00111 = 1,1111$$



$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \text{---} \\
 1,10111
 \end{array}$$

$x = 1,10111$  VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ ulp}(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11 \text{ ulp}(x) \quad |x - \text{succ}(v)| = 0,01 \text{ ulp}(x) \quad \text{donc } RN(x) = \text{succ}(v)$$

$$1,1 + 0,00111 = 1,111$$



$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \text{---} \\
 1,10111
 \end{array}$$

$x = 1,10111$  VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ ulp}(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11 \text{ ulp}(x) \quad |x - \text{succ}(v)| = 0,01 \text{ ulp}(x) \quad \text{donc } RN(x) = \text{succ}(v)$$

$$1,1 \oplus 0,00111 = 1,110$$



$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \text{---} \\
 1,10111
 \end{array}$$

$x = 1,10111$  VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ ulp}(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11 \text{ ulp}(x) \quad |x - \text{succ}(v)| = 0,01 \text{ ulp}(x) \quad \text{donc } RN(x) = \text{succ}(v)$$

$$1,1 \oplus 0,00111 = 1,110$$



$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \text{---} \\
 1,10111
 \end{array}$$

$x = 1,10111$  VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ ulp}(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11 \text{ ulp}(x) \quad |x - \text{succ}(v)| = 0,01 \text{ ulp}(x) \quad \text{donc } RN(x) = \text{succ}(v)$$

$$1,1 \oplus 0,00111 = 1,110$$



$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \text{---} \\
 1,10111
 \end{array}$$

$x = 1,10111$  VF en *toy7* ?

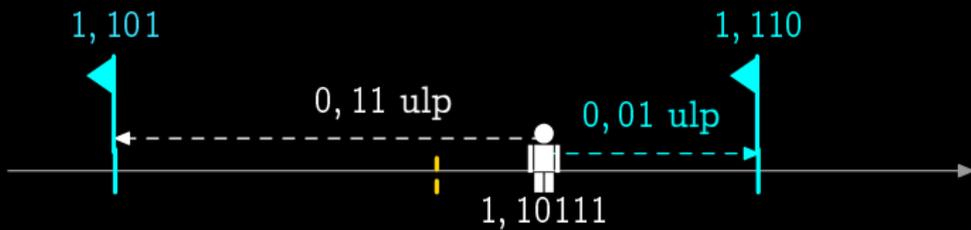
$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ ulp}(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11 \text{ ulp}(x) \quad |x - \text{succ}(v)| = 0,01 \text{ ulp}(x) \quad \text{donc } RN(x) = \text{succ}(v)$$

$$1,1 \oplus 0,00111 = 1,110$$





- 1 Est-ce que l'addition des VF est associative ?

```
In [14]: (1 + 1e-16) + 1e-16
```

```
Out[14]: 1.0
```

```
In [15]: 1 + (1e-16 + 1e-16)
```

```
Out[15]: 1.00000000000000002
```

- 2 Est-on sûr de l'ordre dans le quel un compilateur calcule  $a+b+c+d$  ?



- 1 Est-ce que l'addition des VF est associative ?

```
In [14]: (1 + 1e-16) + 1e-16
```

```
Out[14]: 1.0
```

```
In [15]: 1 + (1e-16 + 1e-16)
```

```
Out[15]: 1.0000000000000002
```

- 2 Est-on sûr de l'ordre dans le quel un compilateur calcule  $a+b+c+d$  ?



- 1 Est-ce que l'addition des VF est associative ?

```
In [14]: (1 + 1e-16) + 1e-16
```

```
Out[14]: 1.0
```

```
In [15]: 1 + (1e-16 + 1e-16)
```

```
Out[15]: 1.0000000000000002
```

- 2 Est-on sûr de l'ordre dans le quel un compilateur calcule  $a + b + c + d$  ?



# Multiplication

- 1 on « xore » les bits de signe ;
- 2 on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
- 3 on multiplie les mantisses ;
- 4 on normalise.



# Multiplication

- 1 on « xore » les bits de signe ;
- 2 on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
- 3 on multiplie les mantisses ;
- 4 on normalise.



# Multiplication

- 1 on « xore » les bits de signe ;
- 2 on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
- 3 on multiplie les mantisses ;
- 4 on normalise.



# Multiplication

- 1 on « xore » les bits de signe ;
- 2 on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
- 3 on multiplie les mantisses ;
- 4 on normalise.



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

⊛ 0 sur 1 donne 1, le bit de signe est 1

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$\begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array}$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec  $|x - 1,100| = 0,011011$  et  $|x - 1,101| = 0,010111$  donc

$$|x - 1,101| < |x - 1,100|$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec  $|x - 1,100| = 0,011ulp(x)$  et  $|x - 1,101| = 0,101ulp(x)$  donc

$$1,100 < x < 1,100 + 0,011ulp(x) = 1,10111$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec  $|x - 1,100| = 0,011 \text{ulp}(x)$  et  $|x - 1,101| = 0,101 \text{ulp}(x)$  donc

$$1,100 < x < 1,100 + 0,011 \text{ulp}(x)$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec  $|x - 1,100| = 0,011ulp(x)$  et  $|x - 1,101| = 0,101ulp(x)$  donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec  $|x - 1,100| = 0,011 \text{ulp}(x)$  et  $|x - 1,101| = 0,101 \text{ulp}(x)$  donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec  $|x - 1,100| = 0,011 \text{ulp}(x)$  et  $|x - 1,101| = 0,101 \text{ulp}(x)$  donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$



$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ②  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
- ③  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
- ④ le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec  $|x - 1,100| = 0,011 \text{ulp}(x)$  et  $|x - 1,101| = 0,101 \text{ulp}(x)$  donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$



- ① Est-ce que la multiplication des VF est associative ?

```
In [1]: (1.00000001 * 1e-15) * 1e15
```

```
Out[1]: 1.0000000100000002
```

```
In [2]: 1.00000001 * (1e-15 * 1e15)
```

```
Out[2]: 1.00000001
```

- ② Est-ce que la multiplication des VF est distributive sur l'addition ?

```
In [3]: 1e15 * (1e-16 + 1)
```

```
Out[3]: 10000000000000000.0
```

```
In [4]: (1e15 * 1e-16) + (1e15 * 1)
```

```
Out[4]: 10000000000000000.1
```



- ① Est-ce que la multiplication des VF est associative ?

```
In [1]: (1.00000001 * 1e-15) * 1e15
```

```
Out[1]: 1.0000000100000002
```

```
In [2]: 1.00000001 * (1e-15 * 1e15)
```

```
Out[2]: 1.00000001
```

- ② Est-ce que la multiplication des VF est distributive sur l'addition ?

```
In [3]: 1e15 * (1e-16 + 1)
```

```
Out[3]: 10000000000000000.0
```

```
In [4]: (1e15 * 1e-16) + (1e15 * 1)
```

```
Out[4]: 10000000000000000.1
```



- ① Est-ce que la multiplication des VF est associative ?

```
In [1]: (1.00000001 * 1e-15) * 1e15
```

```
Out[1]: 1.0000000100000002
```

```
In [2]: 1.00000001 * (1e-15 * 1e15)
```

```
Out[2]: 1.00000001
```

- ② Est-ce que la multiplication des VF est distributive sur l'addition ?

```
In [3]: 1e15 * (1e-16 + 1)
```

```
Out[3]: 10000000000000000.0
```

```
In [4]: (1e15 * 1e-16) + (1e15 * 1)
```

```
Out[4]: 10000000000000000.1
```



- ① Est-ce que la multiplication des VF est associative ?

```
In [1]: (1.00000001 * 1e-15) * 1e15
```

```
Out[1]: 1.0000000100000002
```

```
In [2]: 1.00000001 * (1e-15 * 1e15)
```

```
Out[2]: 1.00000001
```

- ② Est-ce que la multiplication des VF est distributive sur l'addition ?

```
In [3]: 1e15 * (1e-16 + 1)
```

```
Out[3]: 10000000000000000.0
```

```
In [4]: (1e15 * 1e-16) + (1e15 * 1)
```

```
Out[4]: 10000000000000000.1
```



# Sommaire

- 1 Préambule
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



```
def base1(a):  
    if (a + 1.0) - a != 1.0:  
        return a  
    else:  
        return base1(2.0 * a)  
  
def base2(a, b):  
    if (a + b) - a == b:  
        return b  
    else:  
        return base2(a, b + 1.0)
```

????!!!!?????!!???



```

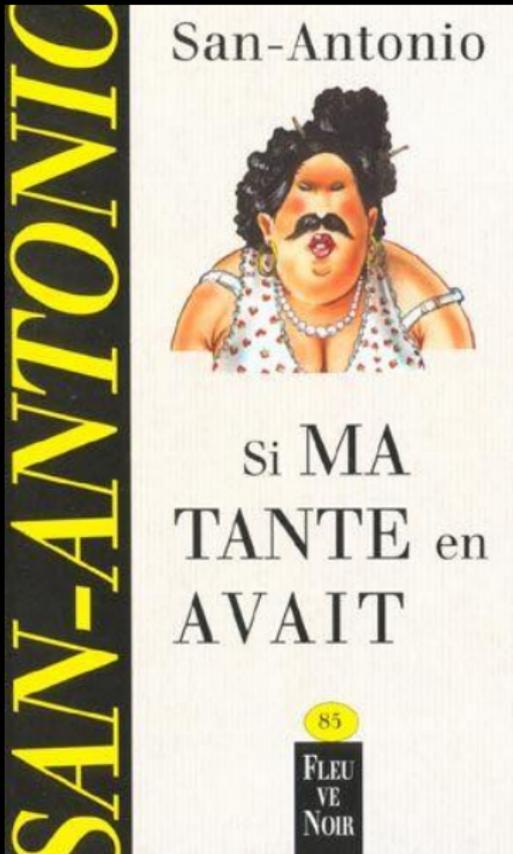
def base1(a):
    if (a + 1.0) - a != 1.0:
        return a
    else:
        return base1(2.0 * a)

def base2(a, b):
    if (a + b) - a == b:
        return b
    else:
        return base2(a, b + 1.0)

```

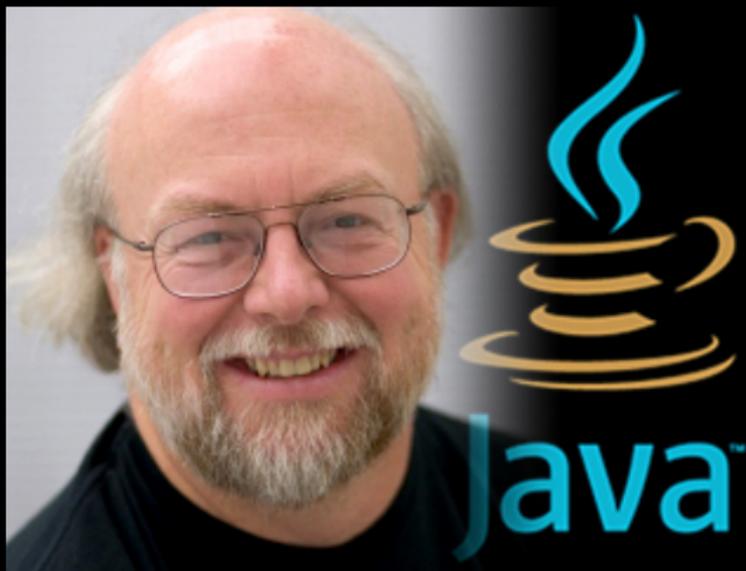
????!!!!?????!!???

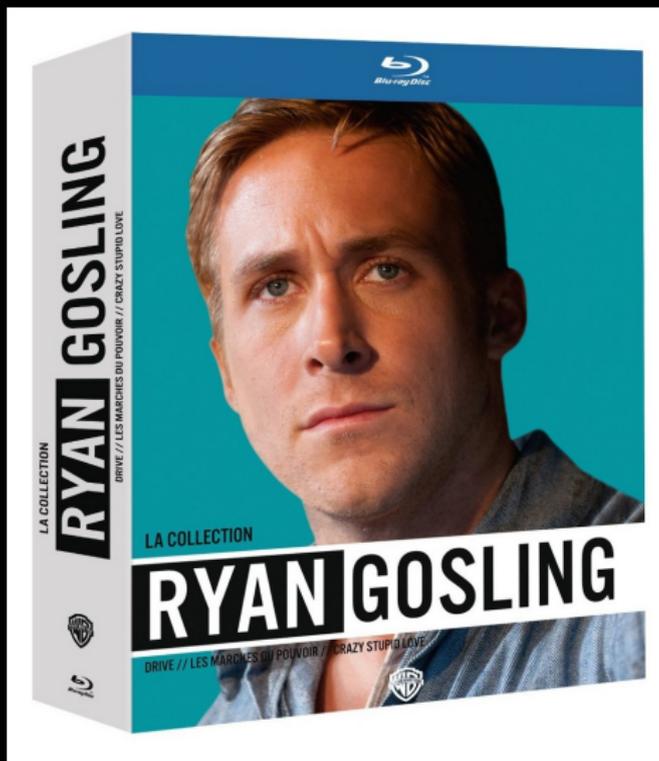




*95 % of the folks out there are completely clueless about floating-point*

*James Gosling (M. Java) - 28 février 1998*





**Précision** (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

3.1777777777777777 est une approximation plutôt précise (16 décimales) mais pas celle de  $\frac{1}{3}$  (décimales à l'infini).



**Précision** (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

3.1777777777777777 est une approximation plutôt précise à 16 décimales, mais pas le nombre décimal exact.



**Précision** (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

L'exactitude (*accuracy*) est ce qui relie un nombre au contexte dans lequel il est employé.

3,1777777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de  $\pi$ .



**Précision** (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

**Exactitude** (*accuracy*) : c'est ce qui relie un nombre au contexte dans lequel il est employé.

3,1777777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de  $\pi$ .



**Précision** (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

**Exactitude** (*accuracy*) : c'est ce qui relie un nombre au contexte dans lequel il est employé.

3,17777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de  $\pi$ .



**Précision** (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

**Exactitude** (*accuracy*) : c'est ce qui relie un nombre au contexte dans lequel il est employé.

3,17777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de  $\pi$ .



Soit  $x$  un nombre et  $\hat{x}$  le nombre qui le représente. On distingue :

l'erreur absolue

l'erreur relative

comme on peut en faire la preuve



Soit  $x$  un nombre et  $\hat{x}$  le nombre qui le représente. On distingue :

l'erreur absolue  $|x - \hat{x}|$

l'erreur relative  $\eta = \frac{x - \hat{x}}{x}$

commises en prenant  $x$  à la place de  $\hat{x}$ .



Soit  $x$  un nombre et  $\hat{x}$  le nombre qui le représente. On distingue :

l'erreur absolue  $|x - \hat{x}|$

l'erreur relative  $\eta = \frac{x - \hat{x}}{x}$  alors  $\hat{x} = x(1 + \eta)$

commise en prenant  $\hat{x}$  à la place de  $x$ .



Soit  $x$  un nombre et  $\hat{x}$  le nombre qui le représente. On distingue :

l'erreur absolue  $|x - \hat{x}|$

l'erreur relative  $\eta = \frac{x - \hat{x}}{x}$  alors  $\hat{x} = x(1 + \eta)$

commise en prenant  $\hat{x}$  à la place de  $x$ .



Soit  $x$  un nombre et  $\hat{x}$  le nombre qui le représente. On distingue :

l'erreur absolue  $|x - \hat{x}|$

l'erreur relative  $\eta = \frac{x - \hat{x}}{x}$  alors  $\hat{x} = x(1 + \eta)$

commise en prenant  $\hat{x}$  à la place de  $x$ .



Soit  $x$  un nombre et  $\hat{x}$  le nombre qui le représente. On distingue :

l'erreur absolue  $|x - \hat{x}|$

l'erreur relative  $\eta = \frac{x - \hat{x}}{x}$  alors  $\hat{x} = x(1 + \eta)$

commise en prenant  $\hat{x}$  à la place de  $x$ .



*Feel nervous, but feel in control. It's not dark magic, it's science.*

Florent de Dinechin



# Sommaire

- 1 Préambule
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



# Sommaire

- 1 Preamble
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 **Que la force de l'erreur soit avec vous**
  - **Atelier Padawan # 1 : majorer l'erreur**
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



# Précision machine



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \tilde{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \tilde{x}}{x} \right| \leq \left| \frac{x - \tilde{x}}{m \times 2^E} \right| \leq \frac{2^{E-n}}{2 \times m \times 2^E} = \frac{1}{2 \times m \times 2^n}$$



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2m} 2^{-n}$$



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2m} 2^{-n}$$



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\epsilon M}{m}$$



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\epsilon_M}{m}$$



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

• Si  $x \approx M$ , alors l'erreur relative est majorée par  $\frac{\varepsilon_M}{2m}$ .

• Si  $x \approx 2^E$ , alors l'erreur relative est majorée par  $\frac{\varepsilon_M}{2}$ .



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

- 1 Si  $\hat{x}$  VF alors l'**erreur relative** est majorée  $\left| \frac{x - \hat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$
- 2 Si  $\hat{x}$  est sous-normal, alors l'**erreur absolue** est majorée



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

- 1 Si  $\hat{x}$  VF alors l'**erreur relative** est majorée  $\left| \frac{x - \hat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$
- 2 Si  $\hat{x}$  est sous-normal, alors l'**erreur absolue** est majorée  $|x - \hat{x}| \leq \frac{1}{2} 2^{E_{\min} - 1 - n}$



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

- 1 Si  $\hat{x}$  VF alors l'**erreur relative** est majorée  $\left| \frac{x - \hat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$
- 2 Si  $\hat{x}$  est sous-normal, alors l'**erreur absolue** est majorée  $|x - \hat{x}| \leq \frac{1}{2} 2^{E_{\min} - 1 - n}$



# Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

- 1 Si  $\hat{x}$  VF alors l'**erreur relative** est majorée  $\left| \frac{x - \hat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$
- 2 Si  $\hat{x}$  est sous-normal, alors l'**erreur absolue** est majorée  $|x - \hat{x}| \leq \frac{1}{2} 2^{E_{\min} - 1 - n}$



NE FAITES PAS DES TESTS D'ÉGALITÉ  
MAIS DES TESTS D'APPARTENANCE  
À DES INTERVALLES DE LARGEUR  $L'\epsilon$   
MACHINE!



# Précision machine

$$\hat{x} = x(1 + \eta_1) + \eta_2$$

avec :

- si  $\hat{x}$  est normal  $|\eta_1| \leq \frac{1}{2}\varepsilon_M$  et  $\eta_2 = 0$
- si  $\hat{x}$  est sous-normal  $\eta_1 = 0$  et  $|\eta_2| \leq \frac{1}{2}2^{E_{\min}-1-n}$

Mouais...



# Précision machine

$$\hat{x} = x(1 + \eta_1) + \eta_2$$

avec :

- si  $\hat{x}$  est normal  $|\eta_1| \leq \frac{1}{2}\varepsilon_M$  et  $\eta_2 = 0$
- si  $\hat{x}$  est sous-normal  $\eta_1 = 0$  et  $|\eta_2| \leq \frac{1}{2}2^{E_{\min}-1-n}$

Mouais...



# Précision machine

$$\hat{x} = x(1 + \eta_1) + \eta_2$$

avec :

- si  $\hat{x}$  est normal  $|\eta_1| \leq \frac{1}{2}\varepsilon_M$  et  $\eta_2 = 0$
- si  $\hat{x}$  est sous-normal  $\eta_1 = 0$  et  $|\eta_2| \leq \frac{1}{2}2^{E_{\min}-1-n}$

Mouais...



# Précision machine

$$\hat{x} = x(1 + \eta_1) + \eta_2$$

avec :

- si  $\hat{x}$  est normal  $|\eta_1| \leq \frac{1}{2}\varepsilon_M$  et  $\eta_2 = 0$
- si  $\hat{x}$  est sous-normal  $\eta_1 = 0$  et  $|\eta_2| \leq \frac{1}{2}2^{E_{\min}-1-n}$

Mouais...



# Élimination

$$\hat{a} = a(1 + \eta_a), \hat{b} = b(1 + \eta_b), x = a - b \text{ et } \hat{x} = \widehat{a - b} = \hat{a} - \hat{b}$$

$$\left| \frac{\hat{x} - x}{x} \right| = \left| \frac{a(1 + \eta_a) - b(1 + \eta_b) - (a - b)}{a - b} \right| = \left| \frac{a\eta_a - b\eta_b}{a - b} \right|$$



# Élimination

$$\hat{a} = a(1 + \eta_a), \hat{b} = b(1 + \eta_b), x = a - b \text{ et } \hat{x} = \widehat{a - b} = \hat{a} - \hat{b}.$$

$$\left| \frac{\hat{x} - x}{x} \right| = \left| \frac{a(1 + \eta_a) - b(1 + \eta_b) - (a - b)}{a - b} \right| = \left| \frac{a\eta_a - b\eta_b}{a - b} \right| \leq \frac{\max\{|a|, |b|\}}{|a - b|} \max\{|\eta_a|, |\eta_b|\}$$



# Élimination

$$\hat{a} = a(1 + \eta_a), \hat{b} = b(1 + \eta_b), x = a - b \text{ et } \hat{x} = \widehat{a - b} = \hat{a} - \hat{b}.$$

$$\left| \frac{x - \hat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| = \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$



# Élimination

$$\hat{a} = a(1 + \eta_a), \hat{b} = b(1 + \eta_b), x = a - b \text{ et } \hat{x} = \widehat{a - b} = \hat{a} - \hat{b}.$$

$$\left| \frac{x - \hat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$



# Élimination

$\hat{a} = a(1 + \eta_a), \hat{b} = b(1 + \eta_b), x = a - b$  et  $\hat{x} = \widehat{a - b} = \hat{a} - \hat{b}$ .

$$\left| \frac{x - \hat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$



# Élimination

$\hat{a} = a(1 + \eta_a)$ ,  $\hat{b} = b(1 + \eta_b)$ ,  $x = a - b$  et  $\hat{x} = \widehat{a - b} = \hat{a} - \hat{b}$ .

$$\left| \frac{x - \hat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$



# Élimination

$\hat{a} = a(1 + \eta_a)$ ,  $\hat{b} = b(1 + \eta_b)$ ,  $x = a - b$  et  $\hat{x} = \widehat{a - b} = \hat{a} - \hat{b}$ .

$$\left| \frac{x - \hat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$



# Catastrophe

```
def deriv(f,h,x):  
    return (f(x+h) - f(x)) / h
```

```
def derivn(f,x,h,n):  
    if n == 0:  
        return f(x)  
    else:  
        return derivn(lambda x: deriv(f,h,x),x,h,n-1)
```



# Catastrophe

```
def deriv(f,h,x):  
    return (f(x+h) - f(x)) / h
```

```
def derivn(f,x,h,n):  
    if n == 0:  
        return f(x)  
    else:  
        return derivn(lambda x: deriv(f,h,x),x,h,n-1)
```



# Catastrophe

```
In [1]: [derivn(lambda x: x**4,1,1e-7,k) for k in range(5)]
```

```
Out[1]:
```

```
[1,  
 4.000000601855902,  
12.012613126444194,  
-222044.6049250313,  
2220446049250.313]
```

```
In [2]: [derivn(lambda x: x**4,1,1e-8,k) for k in range(5)]
```

```
Out[2]: [1, 4.000000042303498, 11.102230246251565, 0.0, 2.  
220446049250313e+16]
```



# Catastrophe

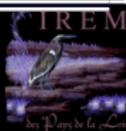
```
In [1]: [derivn(lambda x: x**4,1,1e-7,k) for k in range(5)]
```

```
Out[1]:
```

```
[1,  
 4.000000601855902,  
12.012613126444194,  
-222044.6049250313,  
2220446049250.313]
```

```
In [2]: [derivn(lambda x: x**4,1,1e-8,k) for k in range(5)]
```

```
Out[2]: [1, 4.000000042303498, 11.102230246251565, 0.0, 2.  
220446049250313e+16]
```



# Catastrophe

```
In [3]: [derivn(lambda x: x**4,1,1e-9,k) for k in range(5)]  
Out[3]: [1, 4.000000330961484, 0.0, 0.0, 0.0]
```

```
In [4]: [derivn(lambda x: x**4,1,1e-3,k) for k in range(6)]  
Out[4]:  
[1,  
 4.006004000999486,  
 12.024014000244776,  
 24.03599941303014,  
 24.001245435556484,  
 -2.4424906541753444]
```



# Catastrophe

```
In [3]: [derivn(lambda x: x**4,1,1e-9,k) for k in range(5)]
```

```
Out[3]: [1, 4.000000330961484, 0.0, 0.0, 0.0]
```

```
In [4]: [derivn(lambda x: x**4,1,1e-3,k) for k in range(6)]
```

```
Out[4]:
```

```
[1,  
 4.006004000999486,  
12.024014000244776,  
24.03599941303014,  
24.001245435556484,  
-2.4424906541753444]
```



# Catastrophe

```
In [5]: [derivn(lambda x: x**4,1,2**-10,k) for k in range(6)]  
Out[5]: [1, 4.005863190628588, 12.02345085144043, 24.  
03515625, 24.0, 0.0]
```

```
In [6]: [derivn(lambda x: x**4,1,1.0/1026.,k) for k in range  
(6)]  
Out[6]:  
[1,  
4.005851753982072,  
12.023405112200917,  
24.035087928668187,  
23.999573957547014,  
0.755876563500351]
```



# Catastrophe

```
In [5]: [derivn(lambda x: x**4,1,2**-10,k) for k in range(6)]  
Out[5]: [1, 4.005863190628588, 12.02345085144043, 24.  
03515625, 24.0, 0.0]
```

```
In [6]: [derivn(lambda x: x**4,1,1.0/1026.,k) for k in range  
(6)]  
Out[6]:  
[1,  
4.005851753982072,  
12.023405112200917,  
24.035087928668187,  
23.999573957547014,  
0.755876563500351]
```



# Catastrophe

```
In [7]: [derivn(lambda x: x**4,1,2**-20,k) for k in range(6)]  
Out[7]: [1, 4.000005722045898, 12.0, 0.0, 268435456.0,  
-844424930131968.0]
```



# Catastrophe

```
In [7]: [derivn(exp,0,2**k,k) for k in range(10)]
```

```
Out[7]:
```

```
[1.0,  
 1.0157890399712883,  
 1.0318273737254913,  
 1.0481189373895177,  
 1.0646677284967154,  
 1.081477828323841,  
 1.0985534191131592,  
 1.1158447265625,  
 1.1376953125,  
 0.953125]
```



# Catastrophe

```
In [8]: [derivn(exp,0,2**(-20),k) for k in range(10)]  
Out[8]: [1.0, 1.0000004768371582, 1.0, 0.0, 0.0, 0.0, 0.0, 0.  
0, 0.0, 0.0]
```



# TROP DE PRÉCISION TUE LA PRÉCISION !

DIVISEZ PAR DES PUISSANCES DE 2...PAS DES PUISSANCES DE 10 !



TROP DE PRÉCISION TUE LA  
PRÉCISION !  
DIVISEZ PAR DES PUISSANCES DE  
2...PAS DES PUISSANCES DE 10 !



TROP DE PRÉCISION TUE LA  
PRÉCISION !  
DIVISEZ PAR DES PUISSANCES DE  
2...PAS DES PUISSANCES DE 10 !



## Exercice 4

*Vous savez résoudre une équation du type  $ax^2 + bx + c = 0$  avec  $a \neq 0$ ...*

*Les racines, si elles existent, sont données par une formule bien connue dépendant de  $a$ ,  $b$  et  $c$  :*

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

*Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».*

- ❶ *Que se passe-t-il lorsque  $b^2 \gg |4ac|$  ? En quoi la formule  $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$  peut aider ?*
- ❷ *Que se passe-t-il lorsque  $b^2 \approx 4ac$  ? Peut-on y remédier ? Que peut-on dire de  $\Delta$  par rapport à  $b^2$  ? Y a-t-il élimination catastrophique ?*
- ❸ *Que se passe-t-il dans le cas de l'équation  $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$  ?*

## Exercice 4

*Vous savez résoudre une équation du type  $ax^2 + bx + c = 0$  avec  $a \neq 0$ ...*

*Les racines, si elles existent, sont données par une formule bien connue dépendant de  $a$ ,  $b$  et  $c$  :*

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

*Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».*

- ① *Que se passe-t-il lorsque  $b^2 \gg |4ac|$  ? En quoi la formule  $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$  peut aider ?*
- ② *Que se passe-t-il lorsque  $b^2 \approx 4ac$  ? Peut-on y remédier ? Que peut-on dire de  $\Delta$  par rapport à  $b^2$  ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation  $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$  ?*
- ④ *Et dans le cas de  $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$  ?*

## Exercice 4

*Vous savez résoudre une équation du type  $ax^2 + bx + c = 0$  avec  $a \neq 0$ ...*

*Les racines, si elles existent, sont données par une formule bien connue dépendant de  $a$ ,  $b$  et  $c$  :*

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

*Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».*

- ① *Que se passe-t-il lorsque  $b^2 \gg |4ac|$  ? En quoi la formule  $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$  peut aider ?*
- ② *Que se passe-t-il lorsque  $b^2 \approx 4ac$  ? Peut-on y remédier ? Que peut-on dire de  $\Delta$  par rapport à  $b^2$  ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation  $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$  ?*
- ④ *Et dans le cas de  $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$  ?*
- ⑤ *Moralité ?*

## Exercice 4

*Vous savez résoudre une équation du type  $ax^2 + bx + c = 0$  avec  $a \neq 0$ ...*

*Les racines, si elles existent, sont données par une formule bien connue dépendant de  $a$ ,  $b$  et  $c$  :*

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

*Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».*

- ① *Que se passe-t-il lorsque  $b^2 \gg |4ac|$  ? En quoi la formule  $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$  peut aider ?*
- ② *Que se passe-t-il lorsque  $b^2 \approx 4ac$  ? Peut-on y remédier ? Que peut-on dire de  $\Delta$  par rapport à  $b^2$  ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation  $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$  ?*
- ④ *Et dans le cas de  $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$  ?*
- ⑤ *Moralité ?*

## Exercice 4

*Vous savez résoudre une équation du type  $ax^2 + bx + c = 0$  avec  $a \neq 0$ ...*

*Les racines, si elles existent, sont données par une formule bien connue dépendant de  $a$ ,  $b$  et  $c$  :*

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

*Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».*

- ① *Que se passe-t-il lorsque  $b^2 \gg |4ac|$  ? En quoi la formule  $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$  peut aider ?*
- ② *Que se passe-t-il lorsque  $b^2 \approx 4ac$  ? Peut-on y remédier ? Que peut-on dire de  $\Delta$  par rapport à  $b^2$  ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation  $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$  ?*
- ④ *Et dans le cas de  $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$  ?*
- ⑤ *Moralité ?*

## Exercice 4

*Vous savez résoudre une équation du type  $ax^2 + bx + c = 0$  avec  $a \neq 0$ ...*

*Les racines, si elles existent, sont données par une formule bien connue dépendant de  $a$ ,  $b$  et  $c$  :*

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

*Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».*

- ① *Que se passe-t-il lorsque  $b^2 \gg |4ac|$  ? En quoi la formule  $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$  peut aider ?*
- ② *Que se passe-t-il lorsque  $b^2 \approx 4ac$  ? Peut-on y remédier ? Que peut-on dire de  $\Delta$  par rapport à  $b^2$  ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation  $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$  ?*
- ④ *Et dans le cas de  $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$  ?*
- ⑤ *Moralité ?*

## Exercice 5 (Max IEEE)

**Fonction**  $\text{max}(x, y : \text{flottants}) : \text{flottant}$

**Si**  $x \geq y$  **Alors**

|  $\text{max} = x$

**Sinon**

|  $\text{max} = y$

**FinSi**



```
In [1]: max(1,1+1e-16)
```

```
Out[1]: 1
```

```
In [2]: max(0.1,1e500/1e500)
```

```
Out[2]: 0.1
```

```
In [3]: 0.1 <= 1e500/1e500
```

```
Out[3]: False
```

```
In [4]: 0.1 > 1e500/1e500
```

```
Out[4]: False
```



# Sommaire

- 1 Preamble
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 **Que la force de l'erreur soit avec vous**
  - Atelier Padawan # 1 : majorer l'erreur
  - **Atelier Padawan # 2 : somme de flottants**
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



# Élimination

```
def expn(n):  
    return (1. + 1./n)**n
```



# Élimination

```
In [1]: [(exp(1) - expn(10.  
          0**k)) for k in range(20)  
          ]
```

Out[1]:

```
[0.7182818284590451,  
 0.12453936835904278,  
 0.01346799903751661,  
 0.0013578962234515046,  
 0.000135901634119584,  
 1.359126674760347e-05,  
 1.359363291708604e-06,  
 1.3432696333026684e-07,  
 3.011168736577474e-08,
```

```
-2.2355251516614771e-07,  
-2.2477574246337895e-07,  
-2.248980650598753e-07,  
-0.00024166757819266138,  
 0.002171794372144209,  
 0.0021717943720220845,  
-0.31675337809021675,  
 1.718281828459045,  
 1.718281828459045,  
 1.718281828459045,  
 1.718281828459045]
```



# Élimination

```
In [2]: [(exp(1) - expn(2.  
    0**(3*k))) for k in range  
    (20)]
```

Out[2]:

```
[0.7182818284590451,  
0.1524973145086972,  
0.020936875893946105,  
0.0026498282900537795,  
0.0003317472693793455,  
4.147652875108321e-05,  
5.1846929989274315e-06,  
6.480886076687398e-07,  
8.101110671177025e-08,
```

```
1.0126388616527038e-08,  
1.2657985770658797e-09,  
1.582245445774788e-10,  
1.977795704988239e-11,  
2.4722446312352986e-12,  
3.090860900556436e-13,  
3.863576125695545e-14,  
4.884981308350689e-15,  
4.440892098500626e-16,  
1.718281828459045,  
1.718281828459045]
```



```
In [89]: [2.0**53 + k for k in range(10)]
```

```
Out[89]:
```

```
[9007199254740992.0,  
 9007199254740992.0,  
 9007199254740994.0,  
 9007199254740996.0,  
 9007199254740996.0,  
 9007199254740996.0,  
 9007199254740998.0,  
 9007199254741000.0,  
 9007199254741000.0,  
 9007199254741000.0]
```



```
In [91]: [2.0**55 + k for k in range(10)]
```

```
Out[91]:
```

```
[3.602879701896397e+16,  
 3.602879701896397e+16,  
 3.602879701896397e+16,  
 3.602879701896397e+16,  
 3.602879701896397e+16,  
 3.6028797018963976e+16,  
 3.6028797018963976e+16,  
 3.6028797018963976e+16,  
 3.6028797018963976e+16,  
 3.6028797018963976e+16]
```



## Attention !

Pour éviter de déduire des lemmes suivants des théorèmes totalement faux, n'oubliez pas que **DANS CE QUI SUIT X ET Y SONT DES NOMBRES À VIRGULE FLOTTANTE !**



## Lemme 2 (Majoration de l'erreur d'une somme)

Posons  $x \oplus y = x + y + \text{err}(x \oplus y)$ . Alors, s'il n'y a pas de dépassement de capacité,

$$|\text{err}(x \oplus y)| \leq \min(|x|, |y|)$$

On a bien sûr un résultat analogue pour la différence



- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x$  est un VF situé à la distance  $|y|$  de  $x + y$
- $x \oplus y$  le VF le plus proche de  $x + y$
- $\text{err}(x \oplus y) = |y| \epsilon$
- de même  $\text{err}(x \ominus y) = |y| \epsilon$



- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x$  est un VF situé à la distance  $|y|$  de  $x + y$
- $x \oplus y$  le VF le plus proche de  $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même  $|\text{err}(x \ominus y)| \leq |y|$



- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x$  est un VF situé à la distance  $|y|$  de  $x + y$
- $x \oplus y$  le VF le plus proche de  $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même  $|\text{err}(x \ominus y)| \leq |x|$



- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x$  est un VF situé à la distance  $|y|$  de  $x + y$
- $x \oplus y$  le VF le plus proche de  $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même  $|\text{err}(x \ominus y)| \leq |x|$



- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x$  est un VF situé à la distance  $|y|$  de  $x + y$
- $x \oplus y$  le VF le plus proche de  $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même  $|\text{err}(x \ominus y)| \leq |x|$



- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x$  est un VF situé à la distance  $|y|$  de  $x + y$
- $x \oplus y$  le VF le plus proche de  $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même  $|\text{err}(x \ominus y)| \leq |x|$



- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x$  est un VF situé à la distance  $|y|$  de  $x + y$
- $x \oplus y$  le VF le plus proche de  $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même  $|\text{err}(x \ominus y)| \leq |x|$

De l'importance de disposer avec la IEEE 754 de la meilleure approximation



- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x$  est un VF situé à la distance  $|y|$  de  $x + y$
- $x \oplus y$  le VF le plus proche de  $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même  $|\text{err}(x \ominus y)| \leq |x|$

## À noter

De l'importance de disposer avec la IEEE 754 de la **meilleure approximation** !



## Corollaire 3 (Møller, Knuth, Dekker)

*L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.*



## Corollaire 3 (Møller, Knuth, Dekker)

*L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.*

• Møller

• Knuth

• Dekker



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

- $|x| \geq |y|$
- $x$  et  $y$  sont des VF



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

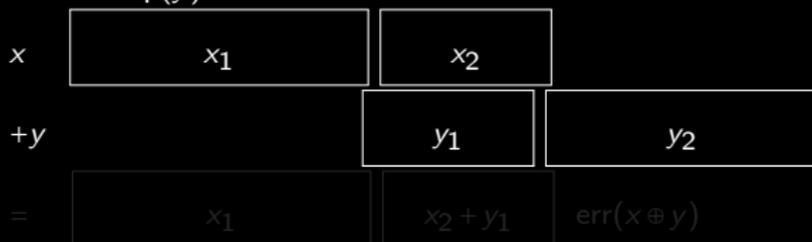
- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

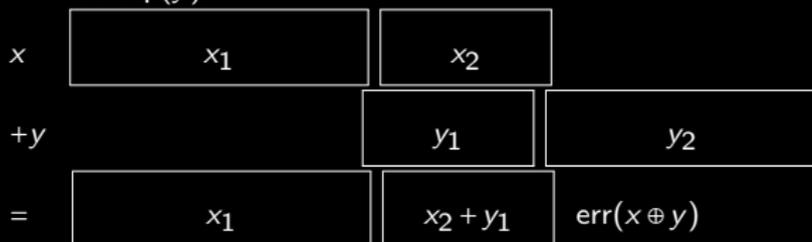
- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

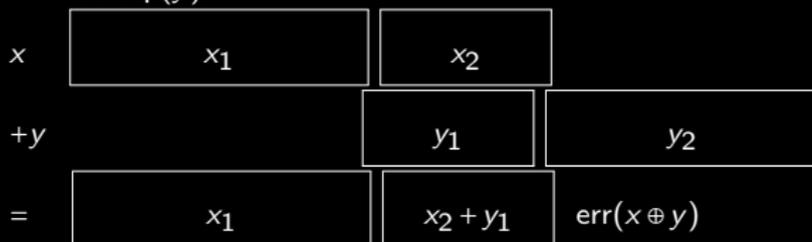
- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



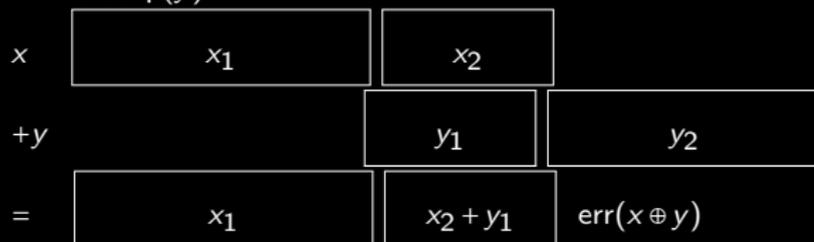
- $|\text{err}(x \oplus y)| \leq |y|$  donc la mantisse entière de  $\text{err}(x \oplus y)$  a une longueur inférieure à  $p$  bits



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



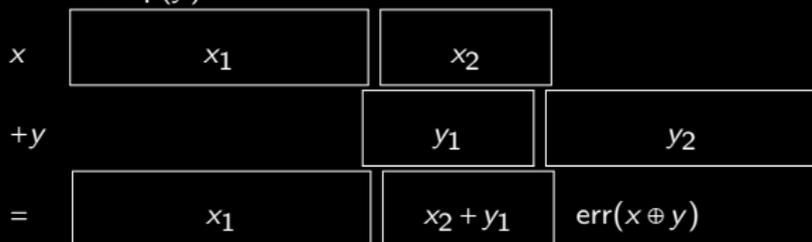
- $|\text{err}(x \oplus y)| \leq |y|$  donc la mantisse entière de  $\text{err}(x \oplus y)$  a une longueur inférieure à  $p$  bits



## Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

- $|x| \geq |y|$
- $x$  et  $y$  sont des VF : le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$



- $|\text{err}(x \oplus y)| \leq |y|$  donc la mantisse entière de  $\text{err}(x \oplus y)$  a une longueur inférieure à  $p$  bits



## Lemme 4

Supposons que  $|x + y| \leq \min(|x|, |y|)$ , alors  $x \oplus y = x + y$ .  
On obtient un résultat analogue pour la soustraction.

- 1 Supposons, sans perdre de généralité, que  $|x| \geq |y|$
- 2 Le plus petit bit significatif de  $x + y$  est au moins de magnitude celle de  $\text{ulp}(y)$
- 3  $|x + y| \leq |y|$  donc la mantisse entière de  $x + y$  a une longueur inférieure à  $p$  bits



## Lemme 4

Supposons que  $|x + y| \leq \min(|x|, |y|)$ , alors  $x \oplus y = x + y$ .  
On obtient un résultat analogue pour la soustraction.

- 1 Supposons, sans perdre de généralité, que  $|x| \geq |y|$
- 2 Le plus petit bit significatif de  $x + y$  est au moins de magnitude celle de  $\text{ulp}(y)$
- 3  $|x + y| \leq |y|$  donc la mantisse entière de  $x + y$  a une longueur inférieure à  $p$  bits



## Lemme 4

Supposons que  $|x + y| \leq \min(|x|, |y|)$ , alors  $x \oplus y = x + y$ .  
On obtient un résultat analogue pour la soustraction.

- 1 Supposons, sans perdre de généralité, que  $|x| \geq |y|$
- 2 Le plus petit bit significatif de  $x + y$  est au moins de magnitude celle de  $\text{ulp}(y)$
- 3  $|x + y| \leq |y|$  donc la mantisse entière de  $x + y$  a une longueur inférieure à  $p$  bits



## Lemme 4

Supposons que  $|x + y| \leq \min(|x|, |y|)$ , alors  $x \oplus y = x + y$ .  
On obtient un résultat analogue pour la soustraction.

- 1 Supposons, sans perdre de généralité, que  $|x| \geq |y|$
- 2 Le plus petit bit significatif de  $x + y$  est au moins de magnitude celle de  $\text{ulp}(y)$
- 3  $|x + y| \leq |y|$  donc la mantisse entière de  $x + y$  a une longueur inférieure à  $p$  bits



## DANGER

Nous avons démontré nos lemmes en considérant des VF  $x$  et  $y$ .  
Est-ce que le résultat suivant contredit notre dernier lemme ?

```
In [86]: 1 - 0.9
```

```
Out[86]: 0.09999999999999998
```



## Lemme 5 (Lemme de Sterbenz (1973))

Soit  $(x, y) \in \mathbb{V}^2$  vérifiant  $\frac{x}{2} \leq y \leq 2x$  alors

$$x \ominus y = x - y$$

*La différence de deux VF suffisamment proches est donc exacte.*

► [Lecture de la preuve de Sterbenz \(1973\) pour valider le lemme de Sterbenz](#)



## Lemme 5 (Lemme de Sterbenz (1973))

Soit  $(x, y) \in \mathbb{V}^2$  vérifiant  $\frac{x}{2} \leq y \leq 2x$  alors

$$x \ominus y = x - y$$

*La différence de deux VF suffisamment proches est donc exacte.*

• Si  $x \leq y$  alors  $x - y \leq -x/2$  donc 0

• Pour la réciproque, on peut utiliser le lemme de Sterbenz



## Lemme 5 (Lemme de Sterbenz (1973))

Soit  $(x, y) \in \mathbb{V}^2$  vérifiant  $\frac{x}{2} \leq y \leq 2x$  alors

$$x \ominus y = x - y$$

*La différence de deux VF suffisamment proches est donc exacte.*

- ①  $x < y$  : alors  $x < y \leq 2x$  donc  $0 < y - x \leq x \leq y$  ;
- ②  $x \geq y$  : alors  $\frac{x}{2} \leq y \leq x$  donc  $-\frac{x}{2} \leq y - x \leq 0$  et par suite  $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$ .

Dans tous les cas  $|x - y| \leq \min(|x|, |y|)$  et on peut utiliser le lemme 4 page 294.



## Lemme 5 (Lemme de Sterbenz (1973))

Soit  $(x, y) \in \mathbb{V}^2$  vérifiant  $\frac{x}{2} \leq y \leq 2x$  alors

$$x \ominus y = x - y$$

*La différence de deux VF suffisamment proches est donc exacte.*

- 1  $x < y$  : alors  $x < y \leq 2x$  donc  $0 < y - x \leq x \leq y$  ;
- 2  $x \geq y$  : alors  $\frac{x}{2} \leq y \leq x$  donc  $-\frac{x}{2} \leq y - x \leq 0$  et par suite  $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$ .

Dans tous les cas  $|x - y| \leq \min(|x|, |y|)$  et on peut utiliser le lemme 4 page 294.



## Lemme 5 (Lemme de Sterbenz (1973))

Soit  $(x, y) \in \mathbb{V}^2$  vérifiant  $\frac{x}{2} \leq y \leq 2x$  alors

$$x \ominus y = x - y$$

*La différence de deux VF suffisamment proches est donc exacte.*

- 1  $x < y$  : alors  $x < y \leq 2x$  donc  $0 < y - x \leq x \leq y$  ;
- 2  $x \geq y$  : alors  $\frac{x}{2} \leq y \leq x$  donc  $-\frac{x}{2} \leq y - x \leq 0$  et par suite  $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$ .

Dans tous les cas  $|x - y| \leq \min(|x|, |y|)$  et on peut utiliser le lemme 4 page 294.



## Lemme 5 (Lemme de Sterbenz (1973))

Soit  $(x, y) \in \mathbb{V}^2$  vérifiant  $\frac{x}{2} \leq y \leq 2x$  alors

$$x \ominus y = x - y$$

*La différence de deux VF suffisamment proches est donc exacte.*

- 1  $x < y$  : alors  $x < y \leq 2x$  donc  $0 < y - x \leq x \leq y$  ;
- 2  $x \geq y$  : alors  $\frac{x}{2} \leq y \leq x$  donc  $-\frac{x}{2} \leq y - x \leq 0$  et par suite  $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$ .

Dans tous les cas  $|x - y| \leq \min(|x|, |y|)$  et on peut utiliser le lemme 4 page 294.



- Concrètement, à quoi correspond cette condition  $\frac{x}{2} \leq y \leq 2x$  ?
- Demandez-vous ce que l'on peut dire de l'écart maximum entre les exposants de  $x$  et  $y$ .



- Concrètement, à quoi correspond cette condition  $\frac{x}{2} \leq y \leq 2x$  ?
- Demandez-vous ce que l'on peut dire de l'écart maximum entre les exposants de  $x$  et  $y$ .



# Sommaire

- 1 Preamble
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 **Que la force de l'erreur soit avec vous**
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - **Atelier Padawan # 3 : somme compensée de flottants quelconques**
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin





## Théorème 6 (Fast2Sum - Dekker & Kahan)

On considère deux VF  $x$  et  $y$  tels que  $|x| \geq |y|$  et l'algorithme suivant :

```
1   $s \leftarrow x \oplus y$ 
2   $y_v \leftarrow s \ominus x$ 
3   $d \leftarrow y \ominus y_v$ 
4  Retourner  $(s, d)$ 
```

Alors  $x + y = s + d$  avec  $s = x \oplus y$  et  $d = \text{err}(x \oplus y)$ . De plus  $s$  et  $d$  ne se chevauchent pas.



## Théorème 6 (Fast2Sum - Dekker & Kahan)

On considère deux VF  $x$  et  $y$  tels que  $|x| \geq |y|$  et l'algorithme suivant :

```
1   $s \leftarrow x \oplus y$ 
2   $y_v \leftarrow s \ominus x$ 
3   $d \leftarrow y \ominus y_v$ 
4  Retourner  $(s, d)$ 
```

Alors  $x + y = s + d$  avec  $s = x \oplus y$  et  $d = \text{err}(x \oplus y)$ . De plus  $s$  et  $d$  ne se chevauchent pas.

- si  $x$  et  $y$  sont de même signe (il s'il  $|y| < \epsilon$ ), alors  $d = 0$  et on peut appliquer le lemme.

• sinon



## Théorème 6 (Fast2Sum - Dekker & Kahan)

On considère deux VF  $x$  et  $y$  tels que  $|x| \geq |y|$  et l'algorithme suivant :

```

1  s ← x ⊕ y
2  yv ← s ⊖ x
3  d ← y ⊖ yv
4  Retourner (s, d)
    
```

Alors  $x + y = s + d$  avec  $s = x \oplus y$  et  $d = \text{err}(x \oplus y)$ . De plus  $s$  et  $d$  ne se chevauchent pas.

- si  $x$  et  $y$  sont de même signe OU si  $|y| \leq \frac{|x|}{2}$  : alors  $\frac{x}{2} \leq s \leq 2x$  et on peut appliquer le lemme ;
- sinon : on a  $x$  et  $y$  de signes opposés ET  $y > \frac{|x|}{2}$  alors si  $y$  est négatif  $x/2 < -y < x$  et sinon  $-x/2 < y < -x$ . Dans les deux cas, d'après le lemme de Sterbenz,  $s$  est calculée exactement et alors  $y_v = y$ .



## Théorème 6 (Fast2Sum - Dekker & Kahan)

On considère deux VF  $x$  et  $y$  tels que  $|x| \geq |y|$  et l'algorithme suivant :

```

1  s ← x ⊕ y
2  yv ← s ⊖ x
3  d ← y ⊖ yv
4  Retourner (s, d)
    
```

Alors  $x + y = s + d$  avec  $s = x \oplus y$  et  $d = \text{err}(x \oplus y)$ . De plus  $s$  et  $d$  ne se chevauchent pas.

- si  $x$  et  $y$  sont de même signe OU si  $|y| \leq \frac{|x|}{2}$  : alors  $\frac{x}{2} \leq s \leq 2x$  et on peut appliquer le lemme ;
- **sinon** : on a  $x$  et  $y$  de signes opposés ET  $y > \frac{|x|}{2}$  alors si  $y$  est négatif  $x/2 < -y < x$  et sinon  $-x/2 < y < -x$ . Dans les deux cas, d'après le lemme de Sterbenz,  $s$  est calculée exactement et alors  $y_v = y$ .





```
def fast2sum(x,y):  
    if abs(x) >= abs(y):  
        s = x + y  
        yv = s - x  
        d = y - yv  
        return (s,d)  
    else:  
        return fast2sum(y,x)
```



```
In [100]: fast2sum(2.0**54,1.0)
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)
Out[104]: (1.8014398509481988e+16, 1.0)
```



```
In [100]: fast2sum(2.0**54,1.0)
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)
Out[104]: (1.8014398509481988e+16, 1.0)
```



```
In [100]: fast2sum(2.0**54,1.0)  
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)  
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)  
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)  
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)  
Out[104]: (1.8014398509481988e+16, 1.0)
```



```
In [100]: fast2sum(2.0**54,1.0)  
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)  
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)  
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)  
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)  
Out[104]: (1.8014398509481988e+16, 1.0)
```



```
In [100]: fast2sum(2.0**54,1.0)  
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)  
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)  
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)  
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)  
Out[104]: (1.8014398509481988e+16, 1.0)
```



## Théorème 7 (2Sum - Knuth )

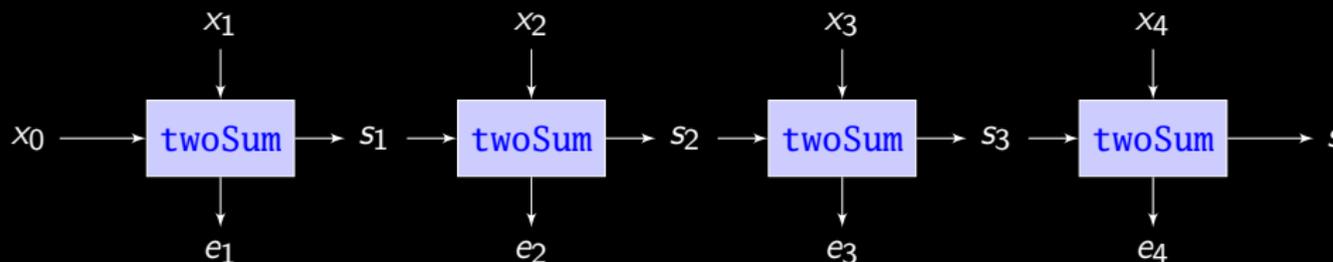
On considère deux VF  $x$  et  $y$  et l'algorithme suivant :

```
1   $s \leftarrow x \oplus y$   
2   $y_v \leftarrow s \ominus x$   
3   $x_v \leftarrow s \ominus y_v$   
4   $y_a \leftarrow y \ominus y_v$   
5   $x_a \leftarrow x \ominus x_v$   
6   $d \leftarrow x_a \oplus y_a$   
7  Retourner  $(s, d)$ 
```

Alors  $x + y = s + d$  avec  $s = x \oplus y$  et  $d = \text{err}(x \oplus y)$ . De plus  $s$  et  $d$  ne se chevauchent pas.



# Somme compensée d'un nombre quelconque de flottants



# Somme compensée d'un nombre quelconque de flottants

```
def sommeKahan(liste):  
    (s,c) = (0.,0.)  
    for x in liste:  
        (s,c) = fast2sum(s, x + c)  
    return s
```

```
def sommePichat(liste):  
    s,e = 0.,0.  
    for x in liste:  
        (s,c) = fast2sum(s, x)  
        e += c  
    return s + e
```



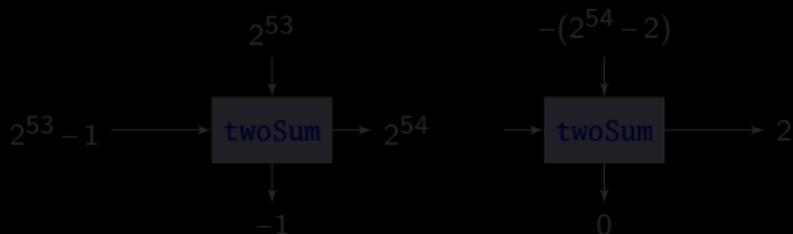
# Somme compensée d'un nombre quelconque de flottants

```
def sommeKahan(liste):  
    (s,c) = (0.,0.)  
    for x in liste:  
        (s,c) = fast2sum(s, x + c)  
    return s
```

```
def sommePichat(liste):  
    s,e = 0.,0.  
    for x in liste:  
        (s,c) = fast2sum(s, x)  
        e += c  
    return s + e
```



$$x_0 = 2^{53} - 1, \quad x_1 = 2^{53} \quad \text{et} \quad x_2 = -(2^{54} - 2)$$

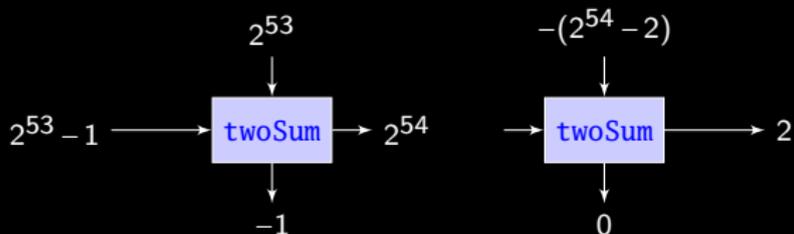


## Exercice

Expliquez les résultats obtenus. Que vous dit l'algorithme de somme compensée ? Et celui de somme en cascade.



$$x_0 = 2^{53} - 1, x_1 = 2^{53} \text{ et } x_2 = -(2^{54} - 2)$$



## Exercice 6

*Expliquez les résultats trouvés. Que va faire l'algorithme de somme compensée ? Et celui de somme en cascade ?*



```
In [105]: sommePichat([1. / n for n in range(1,100001)])
```

```
Out[105]: 12.090146129863427
```

```
In [106]: sommeKahan([1. / n for n in range(1,100001)])
```

```
Out[106]: 12.090146129863427
```

```
In [107]: sum([1. / n for n in range(1,100001)])
```

```
Out[107]: 12.090146129863335
```

```
In [108]: sum([1. / n for n in range(100000,0,-1)])
```

```
Out[108]: 12.090146129863408
```

```
sage: sum(1. / n, n , 1 , 100000).n(64)
```

```
12.0901461298634279
```



```
In [105]: sommePichat([1. / n for n in range(1,100001)])  
Out[105]: 12.090146129863427
```

```
In [106]: sommeKahan([1. / n for n in range(1,100001)])  
Out[106]: 12.090146129863427
```

```
In [107]: sum([1. / n for n in range(1,100001)])  
Out[107]: 12.090146129863335
```

```
In [108]: sum([1. / n for n in range(100000,0,-1)])  
Out[108]: 12.090146129863408
```

```
sage: sum(1. / n, n , 1 , 100000).n(64)  
12.0901461298634279
```



# Sommaire

- 1 Préambule
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



Le sujet du Bac S des Centres Étrangers 2012 fait étudier la suite (vu dans BV APMEP 507) :

$$I_{n+2} = \frac{1}{2}e^{-\frac{n+1}{2}} I_n, \quad I_1 = \frac{1}{2}(e-1)$$

```
def bac(n):
    if n == 1:
        return 0.5 * (exp(1.) - 1.)
    else:
        return (0.5 * (exp(1.) - (n - 1) * bac2012(n - 2)))
```



Le sujet du Bac S des Centres Étrangers 2012 fait étudier la suite (vu dans BV APMEP 507) :

$$I_{n+2} = \frac{1}{2}e - \frac{n+1}{2}I_n, \quad I_1 = \frac{1}{2}(e-1)$$

```
def bac(n):
    if n == 1:
        return 0.5 * (exp(1.) - 1.)
    else:
        return (0.5 * (exp(1.) - (n - 1) * bac2012(n - 2)))
```



```
In [1]: [(2*k + 1, bac(2*k +
1)) for k in range(30)]
```

```
Out[1]:
```

```
[(1, 0.8591409142295225),
 (3, 0.5),
 (5, 0.35914091422952255),
 (7, 0.2817181715409549),
 (9, 0.2322682280657029),
 (11, 0.197799773901008),
 (13, 0.17234227082347453),
 (15, 0.15274501846520083),
 (17, 0.13718076650791589),
 (19, 0.12451401565827958),
 (21, 0.11400075764672679),
 (23, 0.10513258011552784),
 (25, 0.09754995284318846),
 (27, 0.09099152726807258),
```

```
(29, 0.08525953247650642),
 (31, 0.0802479270819263),
 (33, 0.07517408091870181),
 (35, 0.08118153861159172),
 (37, -0.1021267807791284),
 (39, 3.299549749032962),
 (41, -64.63185406642971),
 (43, 1358.6280763092534),
 (45, -29888.458537889346),
 (47, 687435.9055123692),
 (49, -16498460.373155948),
 (51, 412461510.6880396),
 (53, -10723999276.52989),
 (55, 289547980467.66614),
 (57, -8107343453093.293),
 (59, 235112960139706.84)]
```

# Banque chaotique

## Exemple dû à Jean-Michel Muller (in Muller et coll. [2010])

*M. X a récemment été à sa banque (Chaotic Bank Society), pour connaître les nouvelles offres proposées aux meilleurs clients. Son banquier lui propose l'offre suivante : « vous déposez tout d'abord  $e - 1$  euros sur votre compte (où  $e = 2.7182818\dots$  est la base du logarithme népérien). La première année, nous prenons 1 euro sur votre compte de frais de gestion. Par contre, la deuxième année est plus intéressante pour vous, car nous multiplions votre capital restant par 2 et prenons 1 euro de frais de gestion. La troisième année est encore plus intéressante, car nous multiplions votre capital par 3 et prenons 1 euro de frais de gestion. Et ainsi de suite : la  $n$ -ième année, nous multiplions votre capital par  $n$  et prenons 1 euro de frais de gestion. Intéressant, non ? » Pour prendre sa décision, M. X décide de demander l'aide d'un informaticien et se retourne vers vous.*



# Banque chaotique

Exemple dû à Jean-Michel Muller (in Muller et coll. [2010])

*M. X a récemment été à sa banque (Chaotic Bank Society), pour connaître les nouvelles offres proposées aux meilleurs clients. Son banquier lui propose l'offre suivante : « vous déposez tout d'abord  $e - 1$  euros sur votre compte (où  $e = 2.7182818\dots$  est la base du logarithme népérien). La première année, nous prenons 1 euro sur votre compte de frais de gestion. Par contre, la deuxième année est plus intéressante pour vous, car nous multiplions votre capital restant par 2 et prenons 1 euro de frais de gestion. La troisième année est encore plus intéressante, car nous multiplions votre capital par 3 et prenons 1 euro de frais de gestion. Et ainsi de suite : la  $n$ -ième année, nous multiplions votre capital par  $n$  et prenons 1 euro de frais de gestion. Intéressant, non ? » Pour prendre sa décision, M. X décide de demander l'aide d'un informaticien et se retourne vers vous.*



```
def chaotic(n):
    cpt = exp(1) - 1
    for i in range(1,n + 1):
        cpt = i*cpt - 1
    return cpt
```

$$c_n = n! \times \left( c_0 - 1 - \frac{1}{2!} - \frac{1}{3!} - \dots - \frac{1}{n!} \right)$$

$$= n! \times (c_0 - (e-1) + \frac{1}{(n+1)!} + \frac{1}{(n+2)!} + \dots)$$



```
def chaotic(n):
    cpt = exp(1) - 1
    for i in range(1,n + 1):
        cpt = i*cpt - 1
    return cpt
```

$$c_n = n! \times \left( c_0 - 1 - \frac{1}{2!} - \frac{1}{3!} - \dots - \frac{1}{n!} \right)$$

$$= n! \times (c_0 - (e-1) + \frac{1}{(n+1)!} + \frac{1}{(n+2)!} + \dots)$$



# Sommaire

- 1 Préambule
- 2 La norme IEEE 754
  - Différents formats
  - Les normaux, les sous-normaux et les paranormaux
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous
  - Atelier Padawan # 1 : majorer l'erreur
  - Atelier Padawan # 2 : somme de flottants
  - Atelier Padawan # 3 : somme compensée de flottants quelconques
- 6 Des sujets de Bac maladroits...
- 7 Lectures recommandées pour aller plus loin



A. Casamayou, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, N. M. Thiéry et P. Zimmermann,  
Calcul mathématique avec Sage,  
 Amazon, ISBN 9781481191043, adresse : <http://hal.inria.fr/inria-00540485>, 2013,  
 electronic version available under Creative Commons license.

F. de Dinechin,  
 « Page personnelle », adresse : <http://perso.citi-lab.fr/fdedinec/>, 2014.

L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Mpfpr : A multiple-precision binary floating-point library with correct rounding »,  
ACM Trans. Math. Softw., vol. 33, **12**, juin 2007, ISSN 0098-3500, doi : 10.1145/1236463.1236468, adresse :  
<http://doi.acm.org/10.1145/1236463.1236468>.

D. Goldberg,  
 « What every computer scientist should know about floating point arithmetic »,  
ACM Computing Surveys, vol. 23, **11**, p. 5–48, 1991.

F. Goussard,  
 « Page personnelle », adresse : <http://www.goussard.fr/>

W. Ren, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « The GNU Multiple Precision Arithmetic Library »,  
Journal of Supercomputing, vol. 28, **2**, p. 151–178, 2015, doi : 10.1007/s11227-015-0001-0, adresse :  
<http://dx.doi.org/10.1007/s11227-015-0001-0>

F. Goussard,  
 « Page personnelle », adresse : <http://www.goussard.fr/>

M. Muller, G. Hanrot, F. de Dinechin, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Handbook of Floating-Point Arithmetic »,  
 Springer, 2015, ISBN 978-1-4939-9723-7, adresse :  
<http://www.springer.com/9781493997237>



A. Casamayou, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, N. M. Thiéry et P. Zimmermann,  
Calcul mathématique avec Sage,  
 Amazon, ISBN 9781481191043, adresse : <http://hal.inria.fr/inria-00540485>, 2013,  
 electronic version available under Creative Commons license.

F. de Dinechin,  
 « Page personnelle », adresse : <http://perso.citi-lab.fr/fdedinec/>, 2014.

L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Mpfir : A multiple-precision binary floating-point library with correct rounding »,  
ACM Trans. Math. Softw., vol. 33, **12**, juin 2007, ISSN 0098-3500, doi : 10.1145/1236463.1236468, adresse :  
<http://doi.acm.org/10.1145/1236463.1236468>.

D. Goldberg,  
 « What every computer scientist should know about floating point arithmetic »,  
ACM Computing Surveys, vol. 23, **1**, p. 5–48, 1991.

F. Goualard,  
 « Page personnelle », adresse : <http://goualard.frederic.free.fr/>, 2014.

« [www.inria.fr/Equipe/Equipe%20MATH-EDS](http://www.inria.fr/Equipe/Equipe%20MATH-EDS) »  
 « [www.inria.fr/Equipe/Equipe%20MATH-EDS](http://www.inria.fr/Equipe/Equipe%20MATH-EDS) »  
 « [www.inria.fr/Equipe/Equipe%20MATH-EDS](http://www.inria.fr/Equipe/Equipe%20MATH-EDS) »

« [www.inria.fr/Equipe/Equipe%20MATH-EDS](http://www.inria.fr/Equipe/Equipe%20MATH-EDS) »

« [www.inria.fr/Equipe/Equipe%20MATH-EDS](http://www.inria.fr/Equipe/Equipe%20MATH-EDS) »

« [www.inria.fr/Equipe/Equipe%20MATH-EDS](http://www.inria.fr/Equipe/Equipe%20MATH-EDS) »

« [www.inria.fr/Equipe/Equipe%20MATH-EDS](http://www.inria.fr/Equipe/Equipe%20MATH-EDS) »

« [www.inria.fr/Equipe/Equipe%20MATH-EDS](http://www.inria.fr/Equipe/Equipe%20MATH-EDS) »



- A. Casamayou, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, N. M. Thiéry et P. Zimmermann,  
Calcul mathématique avec Sage,  
 Amazon, ISBN 9781481191043, adresse : <http://hal.inria.fr/inria-00540485>, 2013,  
 electronic version available under Creative Commons license.
- F. de Dinechin,  
 « Page personnelle », adresse : <http://perso.citi-lab.fr/fdedinec/>, 2014.
- L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Mpfr : A multiple-precision binary floating-point library with correct rounding »,  
ACM Trans. Math. Softw., vol. 33, **2**, juin 2007, ISSN 0098-3500, doi : 10.1145/1236463.1236468, adresse :  
<http://doi.acm.org/10.1145/1236463.1236468>.
- D. Goldberg,  
 « What every computer scientist should know about floating point arithmetic »,  
ACM Computing Surveys, vol. 23, **1**, p. 5–48, 1991.
- F. Goualard,  
 « Page personnelle », adresse : <http://goualard.frederic.free.fr/>, 2014.
- W. Kahan et J. D. Darcy,  
 « How Java's floating-point hurts everyone everywhere »,  
 Technical report, inst-BERKELEY-MATH-EECS, inst-BERKELEY-MATH-EECS adr, adresse :  
<http://www.cs.berkeley.edu/~wkahan/JAVAHurt.pdf>, juin 1998.



- A. Casamayou, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, N. M. Thiéry et P. Zimmermann,  
Calcul mathématique avec Sage,  
 Amazon, ISBN 9781481191043, adresse : <http://hal.inria.fr/inria-00540485>, 2013,  
 electronic version available under Creative Commons license.
- F. de Dinechin,  
 « Page personnelle », adresse : <http://perso.citi-lab.fr/fdedinec/>, 2014.
- L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Mpfr : A multiple-precision binary floating-point library with correct rounding »,  
ACM Trans. Math. Softw., vol. 33, ||2||, juin 2007, ISSN 0098-3500, doi : 10.1145/1236463.1236468, adresse :  
<http://doi.acm.org/10.1145/1236463.1236468>.
- D. Goldberg,  
 « What every computer scientist should know about floating point arithmetic »,  
ACM Computing Surveys, vol. 23, ||1||, p. 5–48, 1991.
- F. Goualard,  
 « Page personnelle », adresse : <http://goualard.frederic.free.fr/>, 2014.
- W. Kahan et J. D. Darcy,  
 « How Java's floating-point hurts everyone everywhere »,  
 Technical report, inst-BERKELEY-MATH-EECS, inst-BERKELEY-MATH-EECS :adr, adresse :  
<http://www.cs.berkeley.edu/~wkahan/JAVAhurt.pdf>, juin 1998.
- W. Kahan,  
 « Page personnelle », adresse : <http://www.cs.berkeley.edu/~wkahan/>, 2012.



- A. Casamayou, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, N. M. Thiéry et P. Zimmermann,  
Calcul mathématique avec Sage,  
 Amazon, ISBN 9781481191043, adresse : <http://hal.inria.fr/inria-00540485>, 2013,  
 electronic version available under Creative Commons license.
- F. de Dinechin,  
 « Page personnelle », adresse : <http://perso.citi-lab.fr/fdedinec/>, 2014.
- L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Mpfr : A multiple-precision binary floating-point library with correct rounding »,  
ACM Trans. Math. Softw., vol. 33, ||2||, juin 2007, ISSN 0098-3500, doi : 10.1145/1236463.1236468, adresse :  
<http://doi.acm.org/10.1145/1236463.1236468>.
- D. Goldberg,  
 « What every computer scientist should know about floating point arithmetic »,  
ACM Computing Surveys, vol. 23, ||1||, p. 5–48, 1991.
- F. Goualard,  
 « Page personnelle », adresse : <http://goualard.frederic.free.fr/>, 2014.
- W. Kahan et J. D. Darcy,  
 « How Java's floating-point hurts everyone everywhere »,  
 Technical report, inst-BERKELEY-MATH-EECS, inst-BERKELEY-MATH-EECS :adr, adresse :  
<http://www.cs.berkeley.edu/~wkahan/JAVAhurt.pdf>, juin 1998.
- W. Kahan,  
 « Page personnelle », adresse : <http://www.cs.berkeley.edu/~wkahan/>, 2012.
- J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé et S. Torres,  
Handbook of Floating-Point Arithmetic,  
 Birkhäuser Boston, 2010.  
 ACM G.1.0, G.1.2, G.4, B.2.0, B.2.4, F.2.1, ISBN 978-0-8176-4704-9



- A. Casamayou, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, N. M. Thiéry et P. Zimmermann,  
Calcul mathématique avec Sage,  
 Amazon, ISBN 9781481191043, adresse : <http://hal.inria.fr/inria-00540485>, 2013,  
 electronic version available under Creative Commons license.
- F. de Dinechin,  
 « Page personnelle », adresse : <http://perso.citi-lab.fr/fdedinec/>, 2014.
- L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Mpfr : A multiple-precision binary floating-point library with correct rounding »,  
ACM Trans. Math. Softw., vol. 33,  $\parallel 2 \parallel$ , juin 2007, ISSN 0098-3500, doi : 10.1145/1236463.1236468, adresse :  
<http://doi.acm.org/10.1145/1236463.1236468>.
- D. Goldberg,  
 « What every computer scientist should know about floating point arithmetic »,  
ACM Computing Surveys, vol. 23,  $\parallel 1 \parallel$ , p. 5–48, 1991.
- F. Goualard,  
 « Page personnelle », adresse : <http://goualard.frederic.free.fr/>, 2014.
- W. Kahan et J. D. Darcy,  
 « How Java's floating-point hurts everyone everywhere »,  
 Technical report, inst-BERKELEY-MATH-EECS, inst-BERKELEY-MATH-EECS :adr, adresse :  
<http://www.cs.berkeley.edu/~wkahan/JAVAhurt.pdf>, juin 1998.
- W. Kahan,  
 « Page personnelle », adresse : <http://www.cs.berkeley.edu/~wkahan/>, 2012.
- J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé et S. Torres,  
Handbook of Floating-Point Arithmetic,  
 Birkhäuser Boston, 2010,  
 ACM G.1.0 ; G.1.2 ; G.4 ; B.2.0 ; B.2.4 ; F.2.1., ISBN 978-0-8176-4704-9.



- A. Casamayou, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, N. M. Thiéry et P. Zimmermann,  
Calcul mathématique avec Sage,  
 Amazon, ISBN 9781481191043, adresse : <http://hal.inria.fr/inria-00540485>, 2013,  
 electronic version available under Creative Commons license.
- F. de Dinechin,  
 « Page personnelle », adresse : <http://perso.citi-lab.fr/fdedinec/>, 2014.
- L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Mpfr : A multiple-precision binary floating-point library with correct rounding »,  
ACM Trans. Math. Softw., vol. 33,  $\parallel 2 \parallel$ , juin 2007, ISSN 0098-3500, doi : 10.1145/1236463.1236468, adresse :  
<http://doi.acm.org/10.1145/1236463.1236468>.
- D. Goldberg,  
 « What every computer scientist should know about floating point arithmetic »,  
ACM Computing Surveys, vol. 23,  $\parallel 1 \parallel$ , p. 5–48, 1991.
- F. Goualard,  
 « Page personnelle », adresse : <http://goualard.frederic.free.fr/>, 2014.
- W. Kahan et J. D. Darcy,  
 « How Java's floating-point hurts everyone everywhere »,  
 Technical report, inst-BERKELEY-MATH-EECS, inst-BERKELEY-MATH-EECS :adr, adresse :  
<http://www.cs.berkeley.edu/~wkahan/JAVAhurt.pdf>, juin 1998.
- W. Kahan,  
 « Page personnelle », adresse : <http://www.cs.berkeley.edu/~wkahan/>, 2012.
- J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé et S. Torres,  
Handbook of Floating-Point Arithmetic,  
 Birkhäuser Boston, 2010,  
 ACM G.1.0 ; G.1.2 ; G.4 ; B.2.0 ; B.2.4 ; F.2.1., ISBN 978-0-8176-4704-9.



- A. Casamayou, G. Connan, T. Dumont, L. Fousse, F. Maltey, M. Meulien, M. Mezzarobba, C. Pernet, N. M. Thiéry et P. Zimmermann,  
Calcul mathématique avec Sage,  
 Amazon, ISBN 9781481191043, adresse : <http://hal.inria.fr/inria-00540485>, 2013,  
 electronic version available under Creative Commons license.
- F. de Dinechin,  
 « Page personnelle », adresse : <http://perso.citi-lab.fr/fdedinec/>, 2014.
- L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier et P. Zimmermann,  
 « Mpfr : A multiple-precision binary floating-point library with correct rounding »,  
ACM Trans. Math. Softw., vol. 33,  $\parallel 2 \parallel$ , juin 2007, ISSN 0098-3500, doi : 10.1145/1236463.1236468, adresse :  
<http://doi.acm.org/10.1145/1236463.1236468>.
- D. Goldberg,  
 « What every computer scientist should know about floating point arithmetic »,  
ACM Computing Surveys, vol. 23,  $\parallel 1 \parallel$ , p. 5–48, 1991.
- F. Goualard,  
 « Page personnelle », adresse : <http://goualard.frederic.free.fr/>, 2014.
- W. Kahan et J. D. Darcy,  
 « How Java's floating-point hurts everyone everywhere »,  
 Technical report, inst-BERKELEY-MATH-EECS, inst-BERKELEY-MATH-EECS :adr, adresse :  
<http://www.cs.berkeley.edu/~wkahan/JAVAhurt.pdf>, juin 1998.
- W. Kahan,  
 « Page personnelle », adresse : <http://www.cs.berkeley.edu/~wkahan/>, 2012.
- J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé et S. Torres,  
Handbook of Floating-Point Arithmetic,  
 Birkhäuser Boston, 2010,  
 ACM G.1.0 ; G.1.2 ; G.4 ; B.2.0 ; B.2.4 ; F.2.1., ISBN 978-0-8176-4704-9.

