

Nombres flottants

Ivan Noyer

Lycée Leconte de Lisle

- 1 Présentation
- 2 La norme IEEE754
- 3 Exemples et règles d'arrondis
- 4 Problèmes induits par la norme
- 5 Arithmétique psychédélique

Crédits

- Informatique pour tous en classes préparatoires aux grande écoles (Eyrolles)
- Nombres flottants : Wikipedia
- cppreference (en C)

- 1 Présentation
- 2 La norme IEEE754
- 3 Exemples et règles d'arrondis
- 4 Problèmes induits par la norme
- 5 Arithmétique psychédélique

Introduction

Objectif : représenter une partie des nombre rationnels : des nombres avec une quantité bornée de chiffres après la virgule.

- Idée : on se donne tous les chiffres composants le nombre plus la position de la virgule. Dans cet esprit $(1234, 2)$ représente $12,34$, c'est à dire $1.234 \cdot 10^1$.

Introduction

Objectif : représenter une partie des nombre rationnels : des nombres avec une quantité bornée de chiffres après la virgule.

- Idée : on se donne tous les chiffres composants le nombre plus la position de la virgule. Dans cet esprit $(1234, 2)$ représente $12,34$, c'est à dire $1.234 \cdot 10^1$.
- Moyen : tous les nombres réels non nuls peuvent de manière unique s'écrire sous la forme $(-1)^s(1 + m)2^e$ dans laquelle

Introduction

Objectif : représenter une partie des nombre rationnels : des nombres avec une quantité bornée de chiffres après la virgule.

- Idée : on se donne tous les chiffres composants le nombre plus la position de la virgule. Dans cet esprit $(1234, 2)$ représente $12,34$, c'est à dire $1.234 \cdot 10^1$.
- Moyen : tous les nombres réels non nuls peuvent de manière unique s'écrire sous la forme $(-1)^s(1 + m)2^e$ dans laquelle
 - $s \in \{0, 1\}$ est le *signe* du nombre

Introduction

Objectif : représenter une partie des nombre rationnels : des nombres avec une quantité bornée de chiffres après la virgule.

- Idée : on se donne tous les chiffres composants le nombre plus la position de la virgule. Dans cet esprit $(1234, 2)$ représente $12,34$, c'est à dire $1.234 \cdot 10^1$.
- Moyen : tous les nombres réels non nuls peuvent de manière unique s'écrire sous la forme $(-1)^s(1 + m)2^e$ dans laquelle
 - $s \in \{0, 1\}$ est le *signe* du nombre
 - $m \in [0, 1[$ est sa *mantisse*

Introduction

Objectif : représenter une partie des nombre rationnels : des nombres avec une quantité bornée de chiffres après la virgule.

- Idée : on se donne tous les chiffres composants le nombre plus la position de la virgule. Dans cet esprit $(1234, 2)$ représente $12,34$, c'est à dire $1.234 \cdot 10^1$.
- Moyen : tous les nombres réels non nuls peuvent de manière unique s'écrire sous la forme $(-1)^s(1 + m)2^e$ dans laquelle
 - $s \in \{0, 1\}$ est le *signe* du nombre
 - $m \in [0, 1[$ est sa *mantisse*
 - $e \in \mathbb{Z}$ est son *exposant*

Introduction

- Moyen : tous les nombres réels non nuls peuvent de manière unique s'écrire sous la forme $(-1)^s(1 + m)2^e$ dans laquelle $s \in \{0, 1\}$ est le *signe* du nombre, $m \in [0, 1[$ est sa *mantisse*, $e \in \mathbb{Z}$ est son *exposant*.

Introduction

- Moyen : tous les nombres réels non nuls peuvent de manière unique s'écrire sous la forme $(-1)^s(1 + m)2^e$ dans laquelle $s \in \{0, 1\}$ est le *signe* du nombre, $m \in [0, 1[$ est sa *mantisse*, $e \in \mathbb{Z}$ est son *exposant*.
- L'exposant du nombre représente la position de la virgule dans son expression en base 2.

Introduction

- Moyen : tous les nombres réels non nuls peuvent de manière unique s'écrire sous la forme $(-1)^s(1+m)2^e$ dans laquelle $s \in \{0, 1\}$ est le *signe* du nombre, $m \in [0, 1[$ est sa *mantisse*, $e \in \mathbb{Z}$ est son *exposant*.
- L'exposant du nombre représente la position de la virgule dans son expression en base 2.
- $1+m$ est un nombre entre 1 et 2 (exclu). Il s'écrit sous la forme $1.c_1 \dots c_i \dots$ avec les $c_i \in \{0, 1\}$. Comme les ordinateurs ne gèrent que des quantités finies, le nombre de chiffres binaires c_i est borné (en fait il est fixe, le plus souvent égal à 23 ou 52).

Pourquoi les nombres à virgule flottante ?

- L'avantage de la représentation en virgule flottante par rapport à la virgule fixe est que la virgule flottante est capable, à nombre de chiffres égal, de gérer un intervalle de nombres réels plus important.

Pourquoi les nombres à virgule flottante ?

- L'avantage de la représentation en virgule flottante par rapport à la virgule fixe est que la virgule flottante est capable, à nombre de chiffres égal, de gérer un intervalle de nombres réels plus important.
- Considérons une représentation en virgule fixe qui a 5 chiffres dont un après la virgule. Elle peut exprimer 10^5 nombres décimaux dans $[0; 9999.9]$. Avec 5 fois le chiffre 1, on ne représente que 1111.1.

Pourquoi les nombres à virgule flottante ?

- L'avantage de la représentation en virgule flottante par rapport à la virgule fixe est que la virgule flottante est capable, à nombre de chiffres égal, de gérer un intervalle de nombres réels plus important.
- Considérons une représentation en virgule fixe qui a 5 chiffres dont un après la virgule. Elle peut exprimer 10^5 nombres décimaux dans $[0; 9999.9]$. Avec 5 fois le chiffre 1, on ne représente que 1111.1.
- Avec une représentation en virgule flottante et 5 chiffres 1 : on peut représenter 11111, 1111.1, 111.11, 11.111, 1.1111, .11111, donc 6 fois plus d'expressions en virgule flottante qu'en virgule fixe. De plus l'intervalle de représentation est plus grand : $[0; 99999]$.

Pourquoi les nombres à virgule flottante ?

- L'avantage de la représentation en virgule flottante par rapport à la virgule fixe est que la virgule flottante est capable, à nombre de chiffres égal, de gérer un intervalle de nombres réels plus important.
- Considérons une représentation en virgule fixe qui a 5 chiffres dont un après la virgule. Elle peut exprimer 10^5 nombres décimaux dans $[0; 9999.9]$. Avec 5 fois le chiffre 1, on ne représente que 1111.1.
- Avec une représentation en virgule flottante et 5 chiffres 1 : 6 fois plus d'expressions en virgule flottante qu'en virgule fixe. De plus l'intervalle de représentation est plus grand : $[0; 99999]$.
- Cependant, il faut alors coder la position de la virgule. Cela demande donc un peu plus de place.

- 1 Présentation
- 2 La norme IEEE754**
- 3 Exemples et règles d'arrondis
- 4 Problèmes induits par la norme
- 5 Arithmétique psychédélique

IEEE754, présentation

IEEE Standard for Binary Floating-Point Arithmetic (ANSI/IEEE Std 754-1985) (standard IEEE pour l'arithmétique binaire en virgule flottante) : IEEE 754.

- Standard le plus employé actuellement pour le calcul des nombres à virgule flottante dans le domaine informatique, avec les CPU et les FPU.

IEEE754, présentation

IEEE Standard for Binary Floating-Point Arithmetic (ANSI/IEEE Std 754-1985) (standard IEEE pour l'arithmétique binaire en virgule flottante) : IEEE 754.

- Standard le plus employé actuellement pour le calcul des nombres à virgule flottante dans le domaine informatique, avec les CPU et les FPU.
- Le standard définit les formats de représentation des nombres à virgule flottante (signe, mantisse, exposant, nombres dénormalisés) et valeurs spéciales (infinis et NaN), en même temps qu'un ensemble d'opérations sur les nombres flottants.

IEEE754, présentation

IEEE Standard for Binary Floating-Point Arithmetic (ANSI/IEEE Std 754-1985) (standard IEEE pour l'arithmétique binaire en virgule flottante) : IEEE 754.

- Standard le plus employé actuellement pour le calcul des nombres à virgule flottante dans le domaine informatique, avec les CPU et les FPU.
- Le standard définit les formats de représentation des nombres à virgule flottante (signe, mantisse, exposant, nombres dénormalisés) et valeurs spéciales (infinis et NaN), en même temps qu'un ensemble d'opérations sur les nombres flottants.
- Il décrit aussi quatre modes d'arrondi et cinq exceptions (comprenant les conditions dans lesquelles une exception se produit, et ce qui se passe dans ce cas).

IEEE754, présentation (1)

- Les quatre modes d'arrondi :

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - ① Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - ① Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).
 - ② Vers plus l'infini (ex : l'arrondi à l'unité de -3.4 est -3)

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - 1 Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).
 - 2 Vers plus l'infini (ex : l'arrondi à l'unité de -3.4 est -3)
 - 3 Vers zéro (-3.4 arrondi en -3 , 3.4 arrondi en 3).

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - 1 Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).
 - 2 Vers plus l'infini (ex : l'arrondi à l'unité de -3.4 est -3)
 - 3 Vers zéro (-3.4 arrondi en -3 , 3.4 arrondi en 3).
 - 4 Au plus proche (avec le cas particulier de l'équidistance : le nb 1.5 doit-il être arrondi à l'unité à 1 ou 2 ?).

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - 1 Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).
 - 2 Vers plus l'infini (ex : l'arrondi à l'unité de -3.4 est -3)
 - 3 Vers zéro (-3.4 arrondi en -3 , 3.4 arrondi en 3).
 - 4 Au plus proche (avec le cas particulier de l'équidistance : le nb 1.5 doit-il être arrondi à l'unité à 1 ou 2 ?).
- La version 1985 de la norme IEEE 754 définit 4 formats pour représenter des nombres à virgule flottante :

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - 1 Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).
 - 2 Vers plus l'infini (ex : l'arrondi à l'unité de -3.4 est -3)
 - 3 Vers zéro (-3.4 arrondi en -3 , 3.4 arrondi en 3).
 - 4 Au plus proche (avec le cas particulier de l'équidistance : le nb 1.5 doit-il être arrondi à l'unité à 1 ou 2 ?).
- La version 1985 de la norme IEEE 754 définit 4 formats pour représenter des nombres à virgule flottante :
 - 1 simple précision (32 bits : 1 bit de signe, 8 bits d'exposant (-126 à 127), 23 bits de mantisse, avec bit 1 implicite),

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - ① Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).
 - ② Vers plus l'infini (ex : l'arrondi à l'unité de -3.4 est -3)
 - ③ Vers zéro (-3.4 arrondi en -3 , 3.4 arrondi en 3).
 - ④ Au plus proche (avec le cas particulier de l'équidistance : le nb 1.5 doit-il être arrondi à l'unité à 1 ou 2 ?).
- La version 1985 de la norme IEEE 754 définit 4 formats pour représenter des nombres à virgule flottante :
 - ① simple précision (32 bits : 1 bit de signe, 8 bits d'exposant (-126 à 127), 23 bits de mantisse, avec bit 1 implicite),
 - ② simple précision étendue (≥ 43 bits, obsolète, implémenté en pratique par la double précision),

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - ① Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).
 - ② Vers plus l'infini (ex : l'arrondi à l'unité de -3.4 est -3)
 - ③ Vers zéro (-3.4 arrondi en -3 , 3.4 arrondi en 3).
 - ④ Au plus proche (avec le cas particulier de l'équidistance : le nb 1.5 doit-il être arrondi à l'unité à 1 ou 2 ?).
- La version 1985 de la norme IEEE 754 définit 4 formats pour représenter des nombres à virgule flottante :
 - ① simple précision (32 bits : 1 bit de signe, 8 bits d'exposant (-126 à 127), 23 bits de mantisse, avec bit 1 implicite),
 - ② simple précision étendue (≥ 43 bits, obsolète, implémenté en pratique par la double précision),
 - ③ double précision (64 bits : 1 bit de signe, 11 bits d'exposant (-1022 à 1023), 52 bits de mantisse, avec bit 1 implicite),

IEEE754, présentation (1)

- Les quatre modes d'arrondi :
 - ① Vers moins l'infini (ex : l'arrondi de la partie entière, $\lfloor -3.4 \rfloor = -4$).
 - ② Vers plus l'infini (ex : l'arrondi à l'unité de -3.4 est -3)
 - ③ Vers zéro (-3.4 arrondi en -3 , 3.4 arrondi en 3).
 - ④ Au plus proche (avec le cas particulier de l'équidistance : le nb 1.5 doit-il être arrondi à l'unité à 1 ou 2 ?).
- La version 1985 de la norme IEEE 754 définit 4 formats pour représenter des nombres à virgule flottante :
 - ① simple précision (32 bits : 1 bit de signe, 8 bits d'exposant (-126 à 127), 23 bits de mantisse, avec bit 1 implicite),
 - ② simple précision étendue (≥ 43 bits, obsolète, implémenté en pratique par la double précision),
 - ③ double précision (64 bits : 1 bit de signe, 11 bits d'exposant (-1022 à 1023), 52 bits de mantisse, avec bit 1 implicite),
 - ④ double précision étendue (≥ 79 bits)

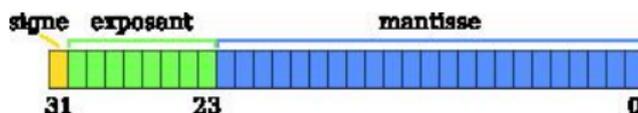
Bit 1, bit *implicite*

Explication

la mantisse représente un nombre décimal entre 1 et 2 (exclu), par exemple 1.1101100011. Rendre le bit 1 implicite consiste à écrire 1101100011, la partie « 1. » étant sous-entendue puisque toujours la même.

Représentation binaire des nombres flottants à précision simple

FIGURE – IEEE-754 simple précision



Plus généralement, les nombres sont écrits au format $(1, E, M)$

donc sur $1 + E + M$ bits :

Signe	Exposant décalé	Mantisse
(1 bit)	(E bits)	(M bits)

Représentation binaire des nombres flottants à précision simple

Ce format (1,8,23) obsolète est privilégié ici parce qu'on peut faire tenir 32 bits sur un transparent.

- bit de poids fort à 1 : négatif, à 0 : positif.

Représentation binaire des nombres flottants à précision simple

Ce format (1,8,23) obsolète est privilégié ici parce qu'on peut faire tenir 32 bits sur un transparent.

- bit de poids fort à 1 : négatif, à 0 : positif.
- Exposant : pas de représentation en complément à 2 (car comparer des nombres serait difficile). L'exposant est *décalé*, afin de le stocker sous forme d'un nombre non signé. En notant E le nombre de chiffres (toujours le même nombre) de l'exposant, on ajoute un *décalage* (ou *biais*) de $d = 2^{E-1} - 1$.

Représentation binaire des nombres flottants à précision simple

Ce format (1,8,23) obsolète est privilégié ici parce qu'on peut faire tenir 32 bits sur un transparent.

- bit de poids fort à 1 : négatif, à 0 : positif.
- Exposant : pas de représentation en complément à 2 (car comparer des nombres serait difficile). L'exposant est *décalé*, afin de le stocker sous forme d'un nombre non signé. En notant E le nombre de chiffres (toujours le même nombre) de l'exposant, on ajoute un *décalage* (ou *biais*) de $d = 2^{E-1} - 1$.
- Avec $E = 8$, $d = 127$. L'exposant est dans l'intervalle $[-127; 128]$, donc l'exposant décalé est dans $[0, 255]$; 0 et 255 ayant une signification spéciale.

Représentation binaire des nombres flottants à précision simple

- Ils sont longs de 4 octets (32 bits) : (1,8,23)

Représentation binaire des nombres flottants à précision simple

- Ils sont longs de 4 octets (32 bits) : (1,8,23)
- La mantisse complète, le significande, doit être considérée comme une valeur sur 24 bits. Si la mantisse avec bit 1 implicite est 101000... alors le significande en base 2 est 1.101000....

Représentation binaire des nombres flottants à précision simple

- Ils sont longs de 4 octets (32 bits) : (1,8,23)
- La mantisse complète, le significande, doit être considérée comme une valeur sur 24 bits. Si la mantisse avec bit 1 implicite est 101000... alors le significande en base 2 est 1.101000....
- La quantité de nombres représentables au format $(1, E, M)$ est grande mais finie. L'ordinateur travaille donc avec des valeurs en général approchées :

Par exemple avec **float x = 0.1** ;, l'ordinateur travaille en interne avec $0\ 01111011\ 10011001100110011001101_2$ c'est à dire 0.100000001490116_{10} qui est égal à $\frac{13421773}{134217728}$.

Affichage

- Quand on entre un nombre au clavier, l'ordinateur en calcule une représentation en virgule flottante. Du fait des arrondis, une infinité de nombres peuvent avoir la même représentation comme flottant.

Affichage

- Quand on entre un nombre au clavier, l'ordinateur en calcule une représentation en virgule flottante. Du fait des arrondis, une infinité de nombres peuvent avoir la même représentation comme flottant.
- La relation \ll a la même représentation que \gg est une relation d'équivalence.

Affichage

- Quand on entre un nombre au clavier, l'ordinateur en calcule une représentation en virgule flottante. Du fait des arrondis, une infinité de nombres peuvent avoir la même représentation comme flottant.
- La relation \ll a la même représentation que \gg est une relation d'équivalence.
- Pour le confort de lecture, l'ordinateur affiche l'unique représentant de cette classe d'équivalence qui nécessite le moins de chiffres décimaux.

Affichage

- Quand on entre un nombre au clavier, l'ordinateur en calcule une représentation en virgule flottante. Du fait des arrondis, une infinité de nombres peuvent avoir la même représentation comme flottant.
- La relation \ll a la même représentation que \gg est une relation d'équivalence.
- Pour le confort de lecture, l'ordinateur affiche l'unique représentant de cette classe d'équivalence qui nécessite le moins de chiffres décimaux.
- C'est la raison pour laquelle l'ordinateur affiche **0.1** et non le nombre rationnel avec lequel il travaille effectivement.

Affichage

- Quand on entre un nombre au clavier, l'ordinateur en calcule une représentation en virgule flottante. Du fait des arrondis, une infinité de nombres peuvent avoir la même représentation comme flottant.
- La relation \ll a la même représentation que \gg est une relation d'équivalence.
- Pour le confort de lecture, l'ordinateur affiche l'unique représentant de cette classe d'équivalence qui nécessite le moins de chiffres décimaux.
- C'est la raison pour laquelle l'ordinateur affiche **0.1** et non le nombre rationnel avec lequel il travaille effectivement.
- Un nombre qui est égal à sa représentation en flottant est dit *représentable exactement* en machine. **0.1** n'est pas représentable exactement en machine ; **1.0** et **3.75** oui.

À propos de l'exposant

Format $(1, E, M)$

- En fonction de la valeur e_d du champ exposant décalé (0, 255, ou autre), certains nombres peuvent avoir une signification spéciale. Ils peuvent être : Des nombres dénormalisés ($e_d = 0$); Zéro ($e_d = 0$); Infini ($e_d = 2^{N-1} - 1$); NaN (*Not a Number* ($e_d = 2^{N-1} - 1$) « pas un nombre », comme $0/0$ ou $\sqrt{-1}$).

À propos de l'exposant

Format $(1, E, M)$

- En fonction de la valeur e_d du champ exposant décalé (0, 255, ou autre), certains nombres peuvent avoir une signification spéciale. Ils peuvent être : Des nombres dénormalisés ($e_d = 0$); Zéro ($e_d = 0$); Infini ($e_d = 2^{N-1} - 1$); NaN (*Not a Number*) ($e_d = 2^{N-1} - 1$) « pas un nombre », comme $0/0$ ou $\sqrt{-1}$.
- L'exposant décalé est dans $\llbracket 0, 2^E - 1 \rrbracket$, le décalage est $2^{E-1} - 1$. Les nombres « normalisés » ont un exposant décalé dans $\llbracket 1, 2^E - 2 \rrbracket$.

À propos de l'exposant

Format $(1, E, M)$

- En fonction de la valeur e_d du champ exposant décalé (0, 255, ou autre), certains nombres peuvent avoir une signification spéciale. Ils peuvent être : Des nombres dénormalisés ($e_d = 0$); Zéro ($e_d = 0$); Infini ($e_d = 2^{N-1} - 1$); NaN (*Not a Number*) ($e_d = 2^{N-1} - 1$) « pas un nombre », comme $0/0$ ou $\sqrt{-1}$.
- L'exposant décalé est dans $\llbracket 0, 2^E - 1 \rrbracket$, le décalage est $2^{E-1} - 1$. Les nombres « normalisés » ont un exposant décalé dans $\llbracket 1, 2^E - 2 \rrbracket$.
- Le bit implicite de la mantisse est déterminé par la valeur de l'exposant décalé. Il vaut 0 si l'exposant décalé est égal à 0 et 1 sinon.

À propos de l'exposant

Nombres normalisés et dé-normalisés

Soit le format $(1, E, M)$

- Si l'exposant décalé est différent de 0 et de $2^E - 1$, le bit implicite est 1, et le nombre est dit *normalisé* :
 $(-1)^s(1 + m)2^e$ avec $m \in [0; 1[$.

À propos de l'exposant

Nombres normalisés et dé-normalisés

Soit le format $(1, E, M)$

- Si l'exposant décalé est différent de 0 et de $2^E - 1$, le bit implicite est 1, et le nombre est dit *normalisé* :
 $(-1)^s(1 + m)2^e$ avec $m \in [0; 1[$.
- Si l'exposant décalé est nul et la mantisse non nulle, par convention, le bit implicite vaut 0. Le nombre est dit *dé-normalisé*. La représentation au format dé-normalisé est destinée aux très petites quantités en valeur absolue :
 $(-1)^s(0 + m)2^{-d+1}$.

À propos de l'exposant

Nombres normalisés et dé-normalisés

Soit le format $(1, E, M)$

- Si l'exposant décalé est différent de 0 et de $2^E - 1$, le bit implicite est 1, et le nombre est dit *normalisé*.
- Si l'exposant décalé est nul et la mantisse non nulle, par convention, le bit implicite vaut 0. Le nombre est dit *dé-normalisé*. La représentation au format dé-normalisé est destinée aux très petites quantités en valeur absolue.
- La quantité de nombres à virgules flottante sur une machine est grande mais finie. Chaque nombre positif (sauf le plus grand) a un *successeur* et un *prédécesseur* (s'il n'est pas nul) positif.

À propos de l'exposant

Nombres normalisés et dé-normalisés

Soit le format $(1, E, M)$

- Si l'exposant décalé est différent de 0 et de $2^E - 1$, le bit implicite est 1, et le nombre est dit *normalisé*.
- Si l'exposant décalé est nul et la mantisse non nulle, par convention, le bit implicite vaut 0. Le nombre est dit *dé-normalisé*. La représentation au format dé-normalisé est destinée aux très petites quantités en valeur absolue.
- La quantité de nombres à virgules flottante sur une machine est grande mais finie. Chaque nombre positif (sauf le plus grand) a un *successeur* et un *prédécesseur* (s'il n'est pas nul) positif.
- Le successeur du nombre dé-normalisé positif le plus grand est le plus petit nombre normalisé positif.

À propos de l'exposant

Nombres normalisés et dé-normalisés

Soit le format $(1, E, M)$

- Si l'exposant décalé est différent de 0 et de $2^E - 1$, le bit implicite est 1, et le nombre est dit *normalisé*.
- Si l'exposant décalé est nul et la mantisse non nulle, par convention, le bit implicite vaut 0. Le nombre est dit *dé-normalisé*. La représentation au format dé-normalisé est destinée aux très petites quantités en valeur absolue.
- La quantité de nombres à virgules flottante sur une machine est grande mais finie. Chaque nombre positif (sauf le plus grand) a un *successeur* et un *prédécesseur* (s'il n'est pas nul) positif.
- Le successeur du nombre dé-normalisé positif le plus grand est le plus petit nombre normalisé positif.
- Zéro : ni normalisé ni dé-normalisé. Deux écritures : $+0$ et -0 .

Nombres dé-normalisés

Pour un format $1, E, M$:

- Mantisse non nulle, champ exposant décalé : E bits à zéro.

Nombres dé-normalisés

Pour un format 1, E , M :

- Mantisse non nulle, champ exposant décalé : E bits à zéro.
- Tous les nombres dé-normalisés ont le même exposant

Nombres dé-normalisés

Pour un format $1, E, M$:

- Mantisse non nulle, champ exposant décalé : E bits à zéro.
- Tous les nombres dé-normalisés ont le même exposant
- Si la règle était la même que pour les nombres normalisés, l'exposant serait donc de 0 – décalage soit $-2^{E-1} + 1$, donc -1023 pour les flottants double précision.

Nombres dé-normalisés

Pour un format 1, E , M :

- Mantisse non nulle, champ exposant décalé : E bits à zéro.
- Tous les nombres dé-normalisés ont le même exposant
- Si la règle était la même que pour les nombres normalisés, l'exposant serait donc de 0 – décalage soit $-2^{E-1} + 1$, donc -1023 pour les flottants double précision.
- Mais, par convention, l'exposant pour les nombres dé-normalisés est en fait égal au plus petit exposant de nombre normalisé soit $-2^{E-1} + 1 + 1$. Ce qui change, c'est le bit implicite 0 (dé-normalisé) ou 1 (normalisé) :

Nombres dé-normalisés

Pour un format 1, E , M :

- Mantisse non nulle, champ exposant décalé : E bits à zéro.
- Tous les nombres dé-normalisés ont le même exposant
- Si la règle était la même que pour les nombres normalisés, l'exposant serait donc de 0 – décalage soit $-2^{E-1} + 1$, donc -1023 pour les flottants double précision.
- Mais, par convention, l'exposant pour les nombres dé-normalisés est en fait égal au plus petit exposant de nombre normalisé soit $-2^{E-1} + 1 + 1$. Ce qui change, c'est le bit implicite 0 (dé-normalisé) ou 1 (normalisé) :
 - Plus petit normalisé positif $(1 + m)2^{-2^{E-1}+1+1}$ avec $m = 0$

Nombres dé-normalisés

Pour un format 1, E , M :

- Mantisse non nulle, champ exposant décalé : E bits à zéro.
- Tous les nombres dé-normalisés ont le même exposant
- Si la règle était la même que pour les nombres normalisés, l'exposant serait donc de 0 – décalage soit $-2^{E-1} + 1$, donc -1023 pour les flottants double précision.
- Mais, par convention, l'exposant pour les nombres dé-normalisés est en fait égal au plus petit exposant de nombre normalisé soit $-2^{E-1} + 1 + 1$. Ce qui change, c'est le bit implicite 0 (dé-normalisé) ou 1 (normalisé) :
 - Plus petit normalisé positif $(1 + m)2^{-2^{E-1}+1+1}$ avec $m = 0$
 - Plus petit dé-normalisé positif $(0 + m)2^{-2^{E-1}+1+1}$ avec m non nul minimum. Ainsi, m s'écrit avec $M - 1$ bits à 0 suivis de 1 donc $m = 2^{-M}$. Conclusion $2^{-M}2^{-2^{E-1}+1+1}$.

Exposant pour un format $(1, E, M)$

Tableau récapitulatif :

Type	Exposant décalé	Mantisse
Zéros : ± 0	0	0
Nombres dénormalisés	0	$\neq 0$, 0. implicite
Nombres normalisés	1 à $2^E - 2$	quelconque, 1. implicite
Infinis $\pm\infty$	$2^E - 1$	0
NaN (<i>Not a Number</i>)	$2^E - 1$	différente de 0

Exemple (Exposant e d'un nombre normalisé)

Si $E = 8$, alors $e \in [-126; 127]$. L'exposant -127 (qui est décalé vers la valeur 0) est réservé pour zéro et les nombres dénormalisés, tandis que l'exposant 128 (décalé vers 255) est réservé pour coder les infinis et les NaN.

Real 2 float

Calcul du triplet (s, e, m)

Cas des nombres normalisés. On calcule (s, e, m) en base 10

- **Ecriture de 0** : 32 bits à 0 ou 1 suivi de 31 bits nuls.

Real 2 float

Calcul du triplet (s, e, m)

Cas des nombres normalisés. On calcule (s, e, m) en base 10

- **Ecriture de 0** : 32 bits à 0 ou 1 suivi de 31 bits nuls.
- Ecriture sous forme scientifique au standard décimal : Si $X \neq 0$, $X = (-1)^s \cdot 2^e \cdot (1 + m)$ avec

Il faudra exprimer ces nombres au format binaire au moyen d'un triplet (s_2, e_2, m_2) .

Real 2 float

Calcul du triplet (s, e, m)

Cas des nombres normalisés. On calcule (s, e, m) en base 10

- **écriture de 0** : 32 bits à 0 ou 1 suivi de 31 bits nuls.
- Écriture sous forme scientifique au standard décimal : Si $X \neq 0$, $X = (-1)^s \cdot 2^e \cdot (1 + m)$ avec
 - $s \in \{0, 1\}$

Il faudra exprimer ces nombres au format binaire au moyen d'un triplet (s_2, e_2, m_2) .

Real 2 float

Calcul du triplet (s, e, m)

Cas des nombres normalisés. On calcule (s, e, m) en base 10

- **écriture de 0** : 32 bits à 0 ou 1 suivi de 31 bits nuls.
- Écriture sous forme scientifique au standard décimal : Si $X \neq 0$, $X = (-1)^s \cdot 2^e \cdot (1 + m)$ avec
 - $s \in \{0, 1\}$
 - $e \in \mathbb{Z}$

Il faudra exprimer ces nombres au format binaire au moyen d'un triplet (s_2, e_2, m_2) .

Real 2 float

Calcul du triplet (s, e, m)

Cas des nombres normalisés. On calcule (s, e, m) en base 10

- **écriture de 0** : 32 bits à 0 ou 1 suivi de 31 bits nuls.
- Écriture sous forme scientifique au standard décimal : Si $X \neq 0$, $X = (-1)^s \cdot 2^e \cdot (1 + m)$ avec
 - $s \in \{0, 1\}$
 - $e \in \mathbb{Z}$
 - $m \in [0, 1[$

Il faudra exprimer ces nombres au format binaire au moyen d'un triplet (s_2, e_2, m_2) .

Real 2 float

Calcul du triplet (s, e, m)

- **Recherche de s** : 0 si X est positif ou nul, 1 sinon.

Real 2 float

Calcul du triplet (s, e, m)

- Recherche de s : 0 si X est positif ou nul, 1 sinon.
- Recherche de e :

Real 2 float

Calcul du triplet (s, e, m)

- **Recherche de s** : 0 si X est positif ou nul, 1 sinon.
- **Recherche de e** :
 - si $|X| \geq 2$, diviser par 2 la valeur absolue de X autant de fois que nécessaire jusqu'à obtenir un entier de l'intervalle $[1; 2[$. e est le nombre de divisions effectuées.

Real 2 float

Calcul du triplet (s, e, m)

- **Recherche de s** : 0 si X est positif ou nul, 1 sinon.
- **Recherche de e** :
 - si $|X| \geq 2$, diviser par 2 la valeur absolue de X autant de fois que nécessaire jusqu'à obtenir un entier de l'intervalle $[1;2[$. e est le nombre de divisions effectuées.
 - si $|X| < 1$, multiplier par 2 la valeur absolue de X par 2 autant de fois que nécessaire jusqu'à obtenir un entier de l'intervalle $[1;2[$. e est l'opposé du nombre de multiplications.

Real 2 float

Calcul du triplet (s, e, m)

- **Recherche de s** : 0 si X est positif ou nul, 1 sinon.
- **Recherche de e** :
 - si $|X| \geq 2$, diviser par 2 la valeur absolue de X autant de fois que nécessaire jusqu'à obtenir un entier de l'intervalle $[1;2[$. e est le nombre de divisions effectuées.
 - si $|X| < 1$, multiplier par 2 la valeur absolue de X par 2 autant de fois que nécessaire jusqu'à obtenir un entier de l'intervalle $[1;2[$. e est l'opposé du nombre de multiplications.
 - si $2 > |X| \geq 1$ alors $e = 0$.

Real 2 float

Calcul du triplet (s, e, m)

- **Recherche de s** : 0 si X est positif ou nul, 1 sinon.
- **Recherche de e** :
 - si $|X| \geq 2$, diviser par 2 la valeur absolue de X autant de fois que nécessaire jusqu'à obtenir un entier de l'intervalle $[1; 2[$. e est le nombre de divisions effectuées.
 - si $|X| < 1$, multiplier par 2 la valeur absolue de X par 2 autant de fois que nécessaire jusqu'à obtenir un entier de l'intervalle $[1; 2[$. e est l'opposé du nombre de multiplications.
 - si $2 > |X| \geq 1$ alors $e = 0$.
- **Recherche de m** : Si on connaît X , s et e alors il est facile de trouver m :

$$m = \frac{(-1)^s}{2^e} X - 1.$$

Real 2 float

Triplet (s_2, e_2, m_2) en base 2

- Exemple : $X = -9.6$; $s_{10} = 1$; $9.6 \times 2^{-3} = 1.2 \in [1; 2[$ donc $e_{10} = 3$; $m_{10} = 1.2 - 1$.

Real 2 float

Triplet (s_2, e_2, m_2) en base 2

- Exemple : $X = -9.6$; $s_{10} = 1$; $9.6 \times 2^{-3} = 1.2 \in [1; 2[$ donc $e_{10} = 3$; $m_{10} = 1.2 - 1$.
- Connaissant s, e, m de la notation scientifique décimale, on veut s_2, e_2, m_2 tels que

Real 2 float

Triplet (s_2, e_2, m_2) en base 2

- Exemple : $X = -9.6$; $s_{10} = 1$; $9.6 \times 2^{-3} = 1.2 \in [1; 2[$ donc $e_{10} = 3$; $m_{10} = 1.2 - 1$.
- Connaissant s, e, m de la notation scientifique décimale, on veut s_2, e_2, m_2 tels que
 - bit de signe $s_2 = s_{10}$,

Real 2 float

Triplet (s_2, e_2, m_2) en base 2

- Exemple : $X = -9.6$; $s_{10} = 1$; $9.6 \times 2^{-3} = 1.2 \in [1; 2[$ donc $e_{10} = 3$; $m_{10} = 1.2 - 1$.
- Connaissant s, e, m de la notation scientifique décimale, on veut s_2, e_2, m_2 tels que
 - bit de signe $s_2 = s_{10}$,
 - bits d'exposant $e_2 = \text{bin}(e + 127)$ (conversion de l'exposant *décalé* en base 2),

Real 2 float

Triplet (s_2, e_2, m_2) en base 2

- Exemple : $X = -9.6$; $s_{10} = 1$; $9.6 \times 2^{-3} = 1.2 \in [1; 2[$ donc $e_{10} = 3$; $m_{10} = 1.2 - 1$.
- Connaissant s, e, m de la notation scientifique décimale, on veut s_2, e_2, m_2 tels que
 - bit de signe $s_2 = s_{10}$,
 - bits d'exposant $e_2 = \text{bin}(e + 127)$ (conversion de l'exposant *décalé* en base 2),
 - Pour la mantisse sur 23 bits

Recommencer à **I** avec la partie décimale du résultat jusqu'à avoir 23 bits de mantisse (et même un peu plus pour arrondir)

Real 2 float

Triplet (s_2, e_2, m_2) en base 2

- Exemple : $X = -9.6$; $s_{10} = 1$; $9.6 \times 2^{-3} = 1.2 \in [1; 2[$ donc $e_{10} = 3$; $m_{10} = 1.2 - 1$.
- Connaissant s, e, m de la notation scientifique décimale, on veut s_2, e_2, m_2 tels que
 - bit de signe $s_2 = s_{10}$,
 - bits d'exposant $e_2 = \text{bin}(e + 127)$ (conversion de l'exposant *décalé* en base 2),
 - Pour la mantisse sur 23 bits
 - | Multiplier m_{10} par 2,

Recommencer à | avec la partie décimale du résultat jusqu'à avoir 23 bits de mantisse (et même un peu plus pour arrondir)

Real 2 float

Triplet (s_2, e_2, m_2) en base 2

- Exemple : $X = -9.6$; $s_{10} = 1$; $9.6 \times 2^{-3} = 1.2 \in [1; 2[$ donc $e_{10} = 3$; $m_{10} = 1.2 - 1$.
- Connaissant s, e, m de la notation scientifique décimale, on veut s_2, e_2, m_2 tels que
 - bit de signe $s_2 = s_{10}$,
 - bits d'exposant $e_2 = \text{bin}(e + 127)$ (conversion de l'exposant *décalé* en base 2),
 - Pour la mantisse sur 23 bits
 - I Multiplier m_{10} par 2,
 - II Calculer partie entière (0 ou 1) ; partie décimale.

Recommencer à I avec la partie décimale du résultat jusqu'à avoir 23 bits de mantisse (et même un peu plus pour arrondir)

Real 2 float

Triplet (s_2, e_2, m_2) en base 2

- Exemple : $X = -9.6$; $s_{10} = 1$; $9.6 \times 2^{-3} = 1.2 \in [1; 2[$ donc $e_{10} = 3$; $m_{10} = 1.2 - 1$.
- Connaissant s, e, m de la notation scientifique décimale, on veut s_2, e_2, m_2 tels que
 - bit de signe $s_2 = s_{10}$,
 - bits d'exposant $e_2 = \text{bin}(e + 127)$ (conversion de l'exposant *décalé* en base 2),
 - Pour la mantisse sur 23 bits
 - I Multiplier m_{10} par 2,
 - II Calculer partie entière (0 ou 1) ; partie décimale.
 - III partie entière : un nouveau bit de la représentation.

Recommencer à I avec la partie décimale du résultat jusqu'à avoir 23 bits de mantisse (et même un peu plus pour arrondir)

Real 2 float

Exemple

- $X = -9.6$, $s_{10} = 1$, $e_{10} = 3$, $m_{10} = 0.2$.

Real 2 float

Exemple

- $X = -9.6$, $s_{10} = 1$, $e_{10} = 3$, $m_{10} = 0.2$.
- $e_{10} + 127_{10} = 130_{10}$, en base 2 : $e_2 = 10000010_2$

Real 2 float

Exemple

- $X = -9.6$, $s_{10} = 1$, $e_{10} = 3$, $m_{10} = 0.2$.
- $e_{10} + 127_{10} = 130_{10}$, en base 2 : $e_2 = 10000010_2$
- Pour la mantisse :

Real 2 float

Exemple

- $X = -9.6$, $s_{10} = 1$, $e_{10} = 3$, $m_{10} = 0.2$.
- $e_{10} + 127_{10} = 130_{10}$, en base 2 : $e_2 = 10000010_2$
- Pour la mantisse :
 - ① $0.2 \times 2 = 0.4 = 0 + 0.4$

Real 2 float

Exemple

- $X = -9.6$, $s_{10} = 1$, $e_{10} = 3$, $m_{10} = 0.2$.
- $e_{10} + 127_{10} = 130_{10}$, en base 2 : $e_2 = 10000010_2$
- Pour la mantisse :
 - 1 $0.2 \times 2 = 0.4 = 0 + 0.4$
 - 2 $0.4 \times 2 = 0.8 = 0 + 0.8$

Real 2 float

Exemple

- $X = -9.6$, $s_{10} = 1$, $e_{10} = 3$, $m_{10} = 0.2$.
- $e_{10} + 127_{10} = 130_{10}$, en base 2 : $e_2 = 10000010_2$
- Pour la mantisse :
 - 1 $0.2 \times 2 = 0.4 = 0 + 0.4$
 - 2 $0.4 \times 2 = 0.8 = 0 + 0.8$
 - 3 $0.8 \times 2 = 1.6 = 1 + 0.6$

Real 2 float

Exemple

- $X = -9.6$, $s_{10} = 1$, $e_{10} = 3$, $m_{10} = 0.2$.
- $e_{10} + 127_{10} = 130_{10}$, en base 2 : $e_2 = 10000010_2$
- Pour la mantisse :
 - 1 $0.2 \times 2 = 0.4 = 0 + 0.4$
 - 2 $0.4 \times 2 = 0.8 = 0 + 0.8$
 - 3 $0.8 \times 2 = 1.6 = 1 + 0.6$
 - 4 $0.6 \times 2 = 1.2 = 1 + 0.2$. On est revenu à 0.2 : séquence infinie.
Mantisse : **0011 0011 0011 0011 0011 0011 0011...**

Real 2 float

Exemple

2^{-1}	2^{-2}	2^{-3}	2^{-4}	2^{-5}	2^{-6}	2^{-7}	2^{-8}	2^{-9}	2^{-10}
0	0	1	1	0	0	1	1	0	0
2^{-11}	2^{-12}	2^{-13}	2^{-14}	2^{-15}	2^{-16}	2^{-17}	2^{-18}	2^{-19}	2^{-20}
1	1	0	0	1	1	0	0	1	1
2^{-21}	2^{-22}	2^{-23}	2^{-24}	2^{-25}	2^{-26}	2^{-27}	...		
0	0	1	1	0	0	1	...		

Real 2 float

Exemple

2^{-1}	2^{-2}	2^{-3}	2^{-4}	2^{-5}	2^{-6}	2^{-7}	2^{-8}	2^{-9}	2^{-10}
0	0	1	1	0	0	1	1	0	0
2^{-11}	2^{-12}	2^{-13}	2^{-14}	2^{-15}	2^{-16}	2^{-17}	2^{-18}	2^{-19}	2^{-20}
1	1	0	0	1	1	0	0	1	1
2^{-21}	2^{-22}	2^{-23}	2^{-24}	2^{-25}	2^{-26}	2^{-27}	...		
0	0	1	1	0	0	1	...		

- On a donc la mantisse comme une somme infinie de

coefficients $\sum_{i=1}^{+\infty} \frac{a_i}{2^i}$ mais on veut une somme finie.

Real 2 float

Exemple

2^{-1}	2^{-2}	2^{-3}	2^{-4}	2^{-5}	2^{-6}	2^{-7}	2^{-8}	2^{-9}	2^{-10}
0	0	1	1	0	0	1	1	0	0
2^{-11}	2^{-12}	2^{-13}	2^{-14}	2^{-15}	2^{-16}	2^{-17}	2^{-18}	2^{-19}	2^{-20}
1	1	0	0	1	1	0	0	1	1
2^{-21}	2^{-22}	2^{-23}	2^{-24}	2^{-25}	2^{-26}	2^{-27}	...		
0	0	1	1	0	0	1	...		

- On a donc la mantisse comme une somme infinie de coefficients $\sum_{i=1}^{+\infty} \frac{a_i}{2^i}$ mais on veut une somme finie.
- Question : $\frac{1}{2^{24}} + 0 \times \frac{1}{2^{25}} + 0 \times \frac{1}{2^{26}} + 1 \times \frac{1}{2^{27}}$ est-il plus proche de $\frac{1}{2^{23}}$ ou de 0 ? Donc m est-il plus proche de $\sum_{i=1}^{23} \frac{a_i}{2^i}$ ou de

$$\left(\sum_{i=1}^{23} \frac{a_i}{2^i} \right) + \frac{1}{2^{24}} + \frac{1}{2^{27}} ?$$

Real 2 float

Exemple

- $\frac{1}{2^{24}} + 0 \cdot \frac{1}{2^{25}} + 0 \cdot \frac{1}{2^{26}} + 1 \cdot \frac{1}{2^{27}}$ est plus proche de $\frac{1}{2^{23}}$: on ajoute donc 1 au bit de poids faible (le 23ème) et on tient compte des retenues.

Attention : Dans le pire cas (mais pas ici), $m = \sum_{i=1}^{23} \frac{1}{2^i}$ et on ajoute $\frac{1}{2^{23}}$. Alors $m + \frac{1}{2^{23}} = \frac{\frac{1}{2} - \frac{1}{2^{24}}}{1 - \frac{1}{2}} + \frac{1}{2^{23}} = 1$ donc $1 + m = 2$ et il faut changer l'exposant !!

Real 2 float

Exemple

- $\frac{1}{2^{24}} + 0 \cdot \frac{1}{2^{25}} + 0 \cdot \frac{1}{2^{26}} + 1 \cdot \frac{1}{2^{27}}$ est plus proche de $\frac{1}{2^{23}}$: on ajoute donc 1 au bit de poids faible (le 23ème) et on tient compte des retenues.

Attention : Dans le pire cas (mais pas ici), $m = \sum_{i=1}^{23} \frac{1}{2^i}$ et on ajoute $\frac{1}{2^{23}}$. Alors $m + \frac{1}{2^{23}} = \frac{\frac{1}{2} - \frac{1}{2^{24}}}{1 - \frac{1}{2}} + \frac{1}{2^{23}} = 1$ donc $1 + m = 2$ et il faut changer l'exposant !!

- Arrondi au plus proche sur 23 bits :

$$M = 0011\ 0011\ 0011\ 0011\ 0011\ 010$$

Real 2 float

Exemple

- $\frac{1}{2^{24}} + 0 \cdot \frac{1}{2^{25}} + 0 \cdot \frac{1}{2^{26}} + 1 \cdot \frac{1}{2^{27}}$ est plus proche de $\frac{1}{2^{23}}$: on ajoute donc 1 au bit de poids faible (le 23ème) et on tient compte des retenues.

Attention : Dans le pire cas (mais pas ici), $m = \sum_{i=1}^{23} \frac{1}{2^i}$ et on ajoute $\frac{1}{2^{23}}$. Alors $m + \frac{1}{2^{23}} = \frac{\frac{1}{2} - \frac{1}{2^{24}}}{1 - \frac{1}{2}} + \frac{1}{2^{23}} = 1$ donc $1 + m = 2$ et il faut changer l'exposant !!

- Arrondi au plus proche sur 23 bits :

$$M = 0011\ 0011\ 0011\ 0011\ 0011\ 010$$

- Finalement : -9.6_{10} est représenté par

$$1_2\ 10000010_2\ 0011\ 0011\ 0011\ 0011\ 0011\ 010_2$$

Real 2 float : arrondi si $M = 23$

- Dans l'exemple étudié $X = -9.6$, on pousse le calcul des décimales jusqu'au 27ème bit et on peut déterminer sans erreur quel est l'arrondi au plus proche de la matrice.

Real 2 float : arrondi si $M = 23$

- Dans l'exemple étudié $X = -9.6$, on pousse le calcul des décimales jusqu'au 27ème bit et on peut déterminer sans erreur quel est l'arrondi au plus proche de la matrice.
- Les bits 24,25,26 s'écrivent en effet **100** et leur connaissance seule ne permet pas de trancher la question : faut-il arrondir par défaut ou par excès (exactement comme arrondir en base 10 le nombre **9.5** à l'unité au plus proche a 2 réponses).

Real 2 float : arrondi si $M = 23$

- Dans l'exemple étudié $X = -9.6$, on pousse le calcul des décimales jusqu'au 27ème bit et on peut déterminer sans erreur quel est l'arrondi au plus proche de la matrice.
- Les bits 24,25,26 s'écrivent en effet **100** et leur connaissance seule ne permet pas de trancher la question : faut-il arrondir par défaut ou par excès (exactement comme arrondir en base 10 le nombre **9.5** à l'unité au plus proche a 2 réponses).
- Dans nos exemples, on ne pousse pas les calculs trop loin après le dernier bit maintenu (le 23ème ici) et on préfère calculer seulement 3 bits supplémentaires. Le cas où ces 3 bits s'écrivent **100** est géré par la règle dite de l'*arrondi au plus proche pair* (cf. plus loin).

Real 2 float : arrondi si $M = 23$

- Dans l'exemple étudié $X = -9.6$, on pousse le calcul des décimales jusqu'au 27ème bit et on peut déterminer sans erreur quel est l'arrondi au plus proche de la matrice.
- Les bits 24,25,26 s'écrivent en effet **100** et leur connaissance seule ne permet pas de trancher la question : faut-il arrondir par défaut ou par excès (exactement comme arrondir en base 10 le nombre **9.5** à l'unité au plus proche a 2 réponses).
- Dans nos exemples, on ne pousse pas les calculs trop loin après le dernier bit maintenu (le 23ème ici) et on préfère calculer seulement 3 bits supplémentaires. Le cas où ces 3 bits s'écrivent **100** est géré par la règle dite de l'*arrondi au plus proche pair* (cf. plus loin).
- Cas dégénéré : On peut aussi arrondir sans utiliser les bits au delà du 23ème : (cf. arrondi au plus proche pair).

Float 2 real

Soit le flottant 0100 0000 1011 1000 0000 0000 0000 0000. À quel réel correspond-il ?



Signe	Exposant décalé	Bit caché + mantisse
0	1000 0001	(1) 011 1000 0000 0000 0000 0000

Float 2 real

Soit le flottant 0100 0000 1011 1000 0000 0000 0000 0000. À quel réel correspond-il ?



Signe	Exposant décalé	Bit caché + mantisse
0	1000 0001	(1) 011 1000 0000 0000 0000 0000

- Le signe est 0, le nombre est donc positif. Le champ exposant décalé est $e_2 = 10000001_2$, autrement dit $e_{10} = 129$. La valeur réelle de l'exposant est donc $e_{10} - d = 129 - 127 = 2$. Le significande (donc avec le bit implicite) est $1.0111000000000000000000_2$.

Float 2 real

Soit le flottant 0100 0000 1011 1000 0000 0000 0000 0000. À quel réel correspond-il ?



Signe	Exposant décalé	Bit caché + mantisse
0	1000 0001	(1) 011 1000 0000 0000 0000 0000

- Le signe est 0, le nombre est donc positif. Le champ exposant décalé est $e_2 = 10000001_2$, autrement dit $e_{10} = 129$. La valeur réelle de l'exposant est donc $e_{10} - d = 129 - 127 = 2$. Le significande (donc avec le bit implicite) est $1.01110000000000000000000_2$.
- Conversion :

$$(-1)^0 \times 2^2 \times \left(\underset{\text{implicite}}{1} \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} + 1 \times 2^{-4} \right) = 5.75$$

float 2 real

Plus généralement, partant d'un flottant simple précision normalisé :

- L'écrire en binaire et retrouver chaque champs.

float 2 real

Plus généralement, partant d'un flottant simple précision normalisé :

- L'écrire en binaire et retrouver chaque champs.
- Signe s : bit de poids fort.

float 2 real

Plus généralement, partant d'un flottant simple précision normalisé :

- L'écrire en binaire et retrouver chaque champs.
- Signe s : bit de poids fort.
- Convertir le binaire du champs exposant en un entier e , lui retrancher le décalage $d = 127 (2^{8-1} - 1)$.

float 2 real

Plus généralement, partant d'un flottant simple précision normalisé :

- L'écrire en binaire et retrouver chaque champs.
- Signe s : bit de poids fort.
- Convertir le binaire du champs exposant en un entier e , lui retrancher le décalage $d = 127 (2^8 - 1)$.
- Partie décimale m (indiquée par la mantisse $b_1 b_2 \dots b_{22} b_{23}$).

$$m = \sum_{i=0}^{22} b_i 2^{-i} = b_1 2^{-1} + b_2 2^{-2} + \dots + b_{22} 2^{-22} + b_{23} 2^{-23}.$$

float 2 real

Plus généralement, partant d'un flottant simple précision normalisé :

- L'écrire en binaire et retrouver chaque champs.
- Signe s : bit de poids fort.
- Convertir le binaire du champs exposant en un entier e , lui retrancher le décalage $d = 127$ ($2^8 - 1$).
- Partie décimale m (indiquée par la mantisse $b_1 b_2 \dots b_{22} b_{23}$).

$$m = \sum_{i=0}^{22} b_i 2^{-i} = b_1 2^{-1} + b_2 2^{-2} + \dots + b_{22} 2^{-22} + b_{23} 2^{-23}.$$

- Le nombre réel correspondant est $(-1)^s \left(\underset{\text{bit implicite}}{\mathbf{1}} + m \right) 2^{e-127}$.

FIGURE – Quelques nombres positifs (Wikipedia)

Type	Exposant	Mantisse	Valeur approchée	Écart / préc
Zéro	0000 0000	000 0000 0000 0000 0000 0000	0,0	
Plus petit nombre dénormalisé	0000 0000	000 0000 0000 0000 0000 0001	$1,4 \times 10^{-45}$	$1,4 \times 10^{-45}$
Nombre dénormalisé suivant	0000 0000	000 0000 0000 0000 0000 0010	$2,8 \times 10^{-45}$	$1,4 \times 10^{-45}$
Nombre dénormalisé suivant	0000 0000	000 0000 0000 0000 0000 0011	$4,2 \times 10^{-45}$	$1,4 \times 10^{-45}$
Autre nombre dénormalisé	0000 0000	100 0000 0000 0000 0000 0000	$5,9 \times 10^{-39}$	
Plus grand nombre dénormalisé	0000 0000	111 1111 1111 1111 1111 1111	$1.17549421 \times 10^{-38}$	
Plus petit nombre normalisé	0000 0001	000 0000 0000 0000 0000 0000	$1.17549435 \times 10^{-38}$	$1,4 \times 10^{-45}$
Nombre normalisé suivant	0000 0001	000 0000 0000 0000 0000 0001	$1.17549449 \times 10^{-38}$	$1,4 \times 10^{-45}$
Presque le double	0000 0001	111 1111 1111 1111 1111 1111	$2,35098856 \times 10^{-38}$	$1,4 \times 10^{-45}$
Nombre normalisé suivant	0000 0010	000 0000 0000 0000 0000 0000	$2,35098870 \times 10^{-38}$	$1,4 \times 10^{-45}$
Nombre normalisé suivant	0000 0010	000 0000 0000 0000 0000 0001	$2.35098898 \times 10^{-38}$	$2,8 \times 10^{-45}$
Presque 1	0111 1110	111 1111 1111 1111 1111 1111	0,99999994	$0,6 \times 10^{-7}$
1	0111 1111	000 0000 0000 0000 0000 0000	1,00000000	
Nombre suivant 1	0111 1111	000 0000 0000 0000 0000 0001	1,00000012	$1,2 \times 10^{-7}$
Presque le plus grand nombre	1111 1110	111 1111 1111 1111 1111 1110	$3,40282326 \times 10^{38}$	
Plus grand nombre normalisé	1111 1110	111 1111 1111 1111 1111 1111	$3,40282346 \times 10^{38}$	2×10^{31}

- 1 Présentation
- 2 La norme IEEE754
- 3 Exemples et règles d'arrondis
- 4 Problèmes induits par la norme
- 5 Arithmétique psychédélique

Un exemple d'arrondi au plus proche

- Soit un nombre de mantisse $110110000000\dots$

Un exemple d'arrondi au plus proche

- Soit un nombre de mantisse $110110000000\dots$
- Le significande complet avec bit caché est $1.1101100\dots$

Un exemple d'arrondi au plus proche

- Soit un nombre de mantisse $110110000000\dots$
- Le significande complet avec bit caché est $1.1101100\dots$
- On veut l'arrondir à 3 chiffres après la virgule. On a le choix entre 1.110 ou $1.110 + 0.001 = 1.111$.

Un exemple d'arrondi au plus proche

- Soit un nombre de mantisse $110110000000\dots$

Un exemple d'arrondi au plus proche

- Soit un nombre de mantisse 1101100000000....
-

$$2^0 + \frac{1}{2} + \frac{1}{2^2} + 0 + \underbrace{\frac{1}{2^4} + \frac{1}{2^5}}_{\text{partie à arrondir}} \sim 111011$$

$$\text{Arrondi par défaut : } 2^0 + \frac{1}{2^1} + \frac{1}{2^2} + 0 \sim 1.110$$

$$\text{Arrondi par excès : } 2^0 + \frac{1}{2^1} + \frac{1}{2^2} + \frac{1}{2^3} \sim 1.111.$$

Un exemple d'arrondi au plus proche

- Soit un nombre de mantisse $110110000000\dots$
-

$$2^0 + \frac{1}{2} + \frac{1}{2^2} + 0 + \underbrace{\frac{1}{2^4} + \frac{1}{2^5}}_{\text{partie à arrondir}} \sim 111011$$

$$\text{Arrondi par défaut : } 2^0 + \frac{1}{2^1} + \frac{1}{2^2} + 0 \sim 1.110$$

$$\text{Arrondi par excès : } 2^0 + \frac{1}{2^1} + \frac{1}{2^2} + \frac{1}{2^3} \sim 1.111.$$

- $r = \frac{1}{2^4} + \frac{1}{2^5}$ est-il plus proche de 0 ou de $\frac{1}{2^3}$?

Un exemple d'arrondi au plus proche

- Soit un nombre de mantisse $110110000000\dots$
-

$$2^0 + \frac{1}{2} + \frac{1}{2^2} + 0 + \underbrace{\frac{1}{2^4} + \frac{1}{2^5}}_{\text{partie à arrondir}} \sim 111011$$

$$\text{Arrondi par défaut : } 2^0 + \frac{1}{2^1} + \frac{1}{2^2} + 0 \sim 1.110$$

$$\text{Arrondi par excès : } 2^0 + \frac{1}{2^1} + \frac{1}{2^2} + \frac{1}{2^3} \sim 1.111.$$

- $r = \frac{1}{2^4} + \frac{1}{2^5}$ est-il plus proche de 0 ou de $\frac{1}{2^3}$?
- Plus proche de $\frac{1}{2^3}$.

Un exemple d'arrondi au plus proche

- Soit un nombre de mantisse $1101100000000\dots$
-

$$2^0 + \frac{1}{2} + \frac{1}{2^2} + 0 + \underbrace{\frac{1}{2^4} + \frac{1}{2^5}}_{\text{partie à arrondir}} \sim 111011$$

$$\text{Arrondi par défaut : } 2^0 + \frac{1}{2^1} + \frac{1}{2^2} + 0 \sim 1.110$$

$$\text{Arrondi par excès : } 2^0 + \frac{1}{2^1} + \frac{1}{2^2} + \frac{1}{2^3} \sim 1.111.$$

- $r = \frac{1}{2^4} + \frac{1}{2^5}$ est-il plus proche de 0 ou de $\frac{1}{2^3}$?
- Plus proche de $\frac{1}{2^3}$.
- Réponse : 1.111.

Choix du nombre le plus proche

Cas de l'examen de 3 bits après le dernier maintenu

- L'exemple précédent était facile car il y avait une seule réponse possible.

Choix du nombre le plus proche

Cas de l'examen de 3 bits après le dernier maintenu

- L'exemple précédent était facile car il y avait une seule réponse possible.
- Mais que se passe-t-il quand on a le choix ? Par exemple, en base 10, comment arrondir à l'unité 1.500 qui est aussi proche de 1 que de 2 ?
Arrondir au plus proche pair signifie choisir 2 plutôt que 1.

Choix du nombre le plus proche

Cas de l'examen de 3 bits après le dernier maintenu

- L'exemple précédent était facile car il y avait une seule réponse possible.
- Mais que se passe-t-il quand on a le choix ? Par exemple, en base 10, comment arrondir à l'unité 1.500 qui est aussi proche de 1 que de 2 ?
Arrondir au plus proche pair signifie choisir 2 plutôt que 1.
- En base deux, le problème se pose lorsque le nombre après le dernier bit maintenu est 100 (si 3 bits au delà du dernier maintenu).

Choix du nombre le plus proche

Cas de l'examen de 3 bits après le dernier maintenu

- L'exemple précédent était facile car il y avait une seule réponse possible.
- Mais que se passe-t-il quand on a le choix ? Par exemple, en base 10, comment arrondir à l'unité 1.500 qui est aussi proche de 1 que de 2 ?
Arrondir au plus proche pair signifie choisir 2 plutôt que 1.
- En base deux, le problème se pose lorsque le nombre après le dernier bit maintenu est 100 (si 3 bits au delà du dernier maintenu).
- En base deux, on prend en général *l'arrondi au plus proche pair*. Il faut que le dernier chiffre de l'écriture binaire soit pair.

Choix du nombre le plus proche

Cas de l'examen de 3 bits après le dernier maintenu

- L'exemple précédent était facile car il y avait une seule réponse possible.
- Mais que se passe-t-il quand on a le choix ? Par exemple, en base 10, comment arrondir à l'unité 1.500 qui est aussi proche de 1 que de 2 ?
Arrondir au plus proche pair signifie choisir 2 plutôt que 1.
- En base deux, le problème se pose lorsque le nombre après le dernier bit maintenu est 100 (si 3 bits au delà du dernier maintenu).
- En base deux, on prend en général *l'arrondi au plus proche pair*. Il faut que le dernier chiffre de l'écriture binaire soit pair.
- Arrondir au plus proche pair revient, lorsqu'on a le choix, à privilégier les écritures qui se terminent par 0.

Exemples d'arrondis au plus proche pair

- Arrondir 1.100100 à 3 chiffres après la virgule : plus proche pair 1.100 (0 est pair).

Exemples d'arrondis au plus proche pair

- Arrondir 1.100100 à 3 chiffres après la virgule : plus proche pair 1.100 (0 est pair).
- Arrondir 1.101100 à 3 chiffres après la virgule. 1 est impair.
Plus proche pair : $1.101 + 0.001 = 1.110$.

Exemples d'arrondis au plus proche pair

- Arrondir 1.100100 à 3 chiffres après la virgule : plus proche pair 1.100 (0 est pair).
- Arrondir 1.101100 à 3 chiffres après la virgule. 1 est impair. Plus proche pair : $1.101 + 0.001 = 1.110$.
- Arrondir 1.111100 à 3 chiffres après la virgule. 1 est impair. Plus proche pair : $1.111 + 0.001 = 10.000$. Il faut changer l'exposant (ajouter 1 à l'exposant)!!

Exemples d'arrondis au plus proche pair

- Arrondir 1.100100 à 3 chiffres après la virgule : plus proche pair 1.100 (0 est pair).
- Arrondir 1.101100 à 3 chiffres après la virgule. 1 est impair. Plus proche pair : $1.101 + 0.001 = 1.110$.
- Arrondir 1.111100 à 3 chiffres après la virgule. 1 est impair. Plus proche pair : $1.111 + 0.001 = 10.000$. Il faut changer l'exposant (ajouter 1 à l'exposant)!!
- Lorsqu'on est dans le cas de figure où il faut changer l'exposant, et que l'exposant est lui même maximum ($254=127+127$ pour les nombres sur 32 bits), on se retrouve avec un nombre considéré comme infini !

Règle d'arrondi au plus proche pair

- Exemple d'arrondi au troisième chiffre après la virgule :

$$1.01110011 = \underbrace{1.011}_{\text{bits maintenus}} \underbrace{10011}_{\text{bits tronqués}} .$$

Règle d'arrondi au plus proche pair

- Exemple d'arrondi au troisième chiffre après la virgule :

$$1.01110011 = \underbrace{1.011}_{\text{bits maintenus}} \underbrace{10011}_{\text{bits tronqués}} .$$

- On considère les trois chiffres après le dernier bit maintenu.

Règle d'arrondi au plus proche pair

- Exemple d'arrondi au troisième chiffre après la virgule :

$$1.01110011 = \underbrace{1.011}_{\text{bits maintenus}} \underbrace{10011}_{\text{bits tronqués}} .$$

- On considère les trois chiffres après le dernier bit maintenu.
 - 1 **0xy** : juste tronquer l'expression (**x,y** sont quelconques).

Règle d'arrondi au plus proche pair

- Exemple d'arrondi au troisième chiffre après la virgule :

$$1.01110011 = \underbrace{1.011}_{\text{bits maintenus}} \underbrace{10011}_{\text{bits tronqués}} .$$

- On considère les trois chiffres après le dernier bit maintenu.
 - 0xy** : juste tronquer l'expression (**x,y** sont quelconques).
 - 100**. Arrondir au plus proche pair :

Règle d'arrondi au plus proche pair

- Exemple d'arrondi au troisième chiffre après la virgule :

$$1.01110011 = \underbrace{1.011}_{\text{bits maintenus}} \underbrace{10011}_{\text{bits tronqués}} .$$

- On considère les trois chiffres après le dernier bit maintenu.
 - 0xy** : juste tronquer l'expression (**x,y** sont quelconques).
 - 100**. Arrondir au plus proche pair :
 - Si le dernier bit mantenu vaut 0 : ne rien faire.

Règle d'arrondi au plus proche pair

- Exemple d'arrondi au troisième chiffre après la virgule :

$$1.01110011 = \underbrace{1.011}_{\text{bits maintenus}} \underbrace{10011}_{\text{bits tronqués}} .$$

- On considère les trois chiffres après le dernier bit maintenu.
 - 0xy** : juste tronquer l'expression (**x,y** sont quelconques).
 - 100**. Arrondir au plus proche pair :
 - Si le dernier bit maintenu vaut 0 : ne rien faire.
 - Sinon, ajouter 1 au dernier bit maintenu en tenant compte des retenus.

Règle d'arrondi au plus proche pair

- Exemple d'arrondi au troisième chiffre après la virgule :

$$1.01110011 = \underbrace{1.011}_{\text{bits maintenus}} \underbrace{10011}_{\text{bits tronqués}} .$$

- On considère les trois chiffres après le dernier bit maintenu.
 - 0xy** : juste tronquer l'expression (**x,y** sont quelconques).
 - 100**. Arrondir au plus proche pair :
 - Si le dernier bit mantenu vaut 0 : ne rien faire.
 - Sinon, ajouter 1 au dernier bit maintenu en tenant compte des retenus.
 - 1xy**, avec $x + y > 0$: ajouter 1 au dernier bit maintenu.

Règle d'arrondi au plus proche pair

- Exemple d'arrondi au troisième chiffre après la virgule :

$$1.01110011 = \underbrace{1.011}_{\text{bits maintenus}} \underbrace{10011}_{\text{bits tronqués}} .$$

- On considère les trois chiffres après le dernier bit maintenu.
 - 0xy** : juste tronquer l'expression (**x,y** sont quelconques).
 - 100**. Arrondir au plus proche pair :
 - Si le dernier bit mantenu vaut 0 : ne rien faire.
 - Sinon, ajouter 1 au dernier bit maintenu en tenant compte des retenus.
 - 1xy**, avec $x + y > 0$: ajouter 1 au dernier bit maintenu.
- Dans l'exemple, les trois chiffres après le dernier bit maintenu forment 100 (équidistance). L'arrondi au plus proche pair est $1.011 + 0.001 = 1.100$.

- 1 Présentation
- 2 La norme IEEE754
- 3 Exemples et règles d'arrondis
- 4 Problèmes induits par la norme**
- 5 Arithmétique psychédélique

Expressions infinies

- Les nombres flottants représentent des rationnels ayant une *expression finie*. Quid des expressions illimitées ?

Expressions infinies

- Les nombres flottants représentent des rationnels ayant une *expression finie*. Quid des expressions illimitées ?
- $\frac{1}{5} = 0.2$ admet une *expression finie en base 10* (EF10), mais pas $\frac{1}{3}$.

Expressions infinies

- Les nombres flottants représentent des rationnels ayant une *expression finie*. Quid des expressions illimitées ?
- $\frac{1}{5} = 0.2$ admet une *expression finie en base 10* (EF10), mais pas $\frac{1}{3}$.
- En base 2, $\frac{1}{a}$ admet une EF2 si et seulement si $a = 2^n$ ($n > 0$).

Expressions infinies

- Les nombres flottants représentent des rationnels ayant une *expression finie*. Quid des expressions illimitées ?
- $\frac{1}{5} = 0.2$ admet une *expression finie en base 10* (EF10), mais pas $\frac{1}{3}$.
- En base 2, $\frac{1}{a}$ admet une EF2 si et seulement si $a = 2^n$ ($n > 0$).
- $\frac{1}{10} = 0.1$ en base 10. Donc $\frac{1}{10}$ a une EF10 mais pas une EF2 car $10 = 2 \times 5$ et que 5 est premier avec 2.

Expressions infinies

- Les nombres flottants représentent des rationnels ayant une *expression finie*. Quid des expressions illimitées ?
- $\frac{1}{5} = 0.2$ admet une *expression finie en base 10* (EF10), mais pas $\frac{1}{3}$.
- En base 2, $\frac{1}{a}$ admet une EF2 si et seulement si $a = 2^n$ ($n > 0$).
- $\frac{1}{10} = 0.1$ en base 10. Donc $\frac{1}{10}$ a une EF10 mais pas une EF2 car $10 = 2 \times 5$ et que 5 est premier avec 2.
- Il faut donc arrondir. Le standard IEEE-754 prévoit 5 méthodes.

Expressions infinies

- Les nombres flottants représentent des rationnels ayant une *expression finie*. Quid des expressions illimitées ?
- $\frac{1}{5} = 0.2$ admet une *expression finie en base 10* (EF10), mais pas $\frac{1}{3}$.
- En base 2, $\frac{1}{a}$ admet une EF2 si et seulement si $a = 2^n$ ($n > 0$).
- $\frac{1}{10} = 0.1$ en base 10. Donc $\frac{1}{10}$ a une EF10 mais pas une EF2 car $10 = 2 \times 5$ et que 5 est premier avec 2.
- Il faut donc arrondir. Le standard IEEE-754 prévoit 5 méthodes.
- Règle de l'*Arrondi correct* : Une fois un mode d'arrondi choisi, le résultat d'une opération est déterministe : un seul résultat possible.

Un exemple d'arrondi

Un dixième

- 0.1 en base 10, correspond à la séquence suivante :
 $s = 0 = S$, $e = -4 + 127$ donc $E = 01111011$, M est une
séquence infinie $1.1001\ 1001\ 1001\ 1001\ 1001\ 100\ 110\ 0\dots$
Alors $1 + m = 1.10011001100110011001101_2$.

Un exemple d'arrondi

Un dixième

- 0.1 en base 10, correspond à la séquence suivante :
 $s = 0 = S$, $e = -4 + 127$ donc $E = 01111011$, M est une séquence infinie $1.1001\ 1001\ 1001\ 1001\ 1001\ 100\ 110\ 0\dots$
Alors $1 + m = 1.10011001100110011001101_2$.
- Du fait des arrondis, le nombre qui représente 0.1 sur un ordinateur avec flottants 32 bits est en fait le nombre $\frac{13421773}{134217728}$ soit 0.100000001490116 .

Règle de l'arrondi correct

La norme IEEE 754 impose l'arrondi correct pour les 5 opérations de base et la racine carrée : Un programme les utilisant donne le même résultat sur toute configuration (machine, système, processeur). Sous réserve :

- 1 qu'il 'y ait pas de précision intermédiaire étendue (ou alors désactivée). Ca veut dire que les résultats intermédiaires du calcul d'une expression ne doivent pas être calculés avec une précision plus grande que celle attendue pour le résultat.

Règle de l'arrondi correct

La norme IEEE 754 impose l'arrondi correct pour les 5 opérations de base et la racine carrée : Un programme les utilisant donne le même résultat sur toute configuration (machine, système, processeur). Sous réserve :

- 1 qu'il 'y ait pas de précision intermédiaire étendue (ou alors désactivée). Ca veut dire que les résultats intermédiaires du calcul d'une expression ne doivent pas être calculés avec une précision plus grande que celle attendue pour le résultat.
- 2 le compilateur ne doit pas changer l'ordre des opérations si cela peut conduire à un résultat différent.

Problème d'arrondi célèbre

Patriot

- En 1991, un anti-missile Patriot rate l'interception d'un missile irakien Skud, lequel blesse 100 personnes, en tue 28.

Problème d'arrondi célèbre

Patriot

- En 1991, un anti-missile Patriot rate l'interception d'un missile irakien Skud, lequel blesse 100 personnes, en tue 28.
- Un micro-processeur interne calcule l'heure en multiples de dixièmes de secondes.

Problème d'arrondi célèbre

Patriot

- En 1991, un anti-missile Patriot rate l'interception d'un missile irakien Skud, lequel blesse 100 personnes, en tue 28.
- Un micro-processeur interne calcule l'heure en multiples de dixièmes de secondes.
- Le nombre de dixièmes de secondes depuis le démarrage est stocké dans un registre entier puis multiplié par une approx. de $\frac{1}{10}$ sur 24 bits pour obtenir le temps en seconde.

Problème d'arrondi célèbre

Patriot

- En 1991, un anti-missile Patriot rate l'interception d'un missile irakien Skud, lequel blesse 100 personnes, en tue 28.
- Un micro-processeur interne calcule l'heure en multiples de dixièmes de secondes.
- Le nombre de dixièmes de secondes depuis le démarrage est stocké dans un registre entier puis multiplié par une approx. de $\frac{1}{10}$ sur 24 bits pour obtenir le temps en seconde.
- Approximation : $209715 \cdot 2^{-21}$, erreur 10^{-7} .

Problème d'arrondi célèbre

Patriot

- En 1991, un anti-missile Patriot rate l'interception d'un missile irakien Skud, lequel blesse 100 personnes, en tue 28.
- Un micro-processeur interne calcule l'heure en multiples de dixièmes de secondes.
- Le nombre de dixièmes de secondes depuis le démarrage est stocké dans un registre entier puis multiplié par une approx. de $\frac{1}{10}$ sur 24 bits pour obtenir le temps en seconde.
- Approximation : $209715 \cdot 2^{-21}$, erreur 10^{-7} .
- Processeur démarré 100h auparavant. Erreur de 0.34 secondes pendant lequel le Skud parcourt 500 m.

Problème d'arrondi célèbre

Patriot

- En 1991, un anti-missile Patriot rate l'interception d'un missile irakien Skud, lequel blesse 100 personnes, en tue 28.
- Un micro-processeur interne calcule l'heure en multiples de dixièmes de secondes.
- Le nombre de dixièmes de secondes depuis le démarrage est stocké dans un registre entier puis multiplié par une approx. de $\frac{1}{10}$ sur 24 bits pour obtenir le temps en seconde.
- Approximation : $209715 \cdot 2^{-21}$, erreur 10^{-7} .
- Processeur démarré 100h auparavant. Erreur de 0.34 secondes pendant lequel le Skud parcourt 500 m.
- Le Patriot rate sa cible, pas le Skud.

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .
- Soit x' arrondi en précision $q < p$ de y .

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .
- Soit x' arrondi en précision $q < p$ de y .
- x' n'est pas toujours l'arrondi en précision q de x !

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .
- Soit x' arrondi en précision $q < p$ de y .
- x' n'est pas toujours l'arrondi en précision q de x !
 - ① $x = 1.0110100101$. Arrondir à 7 chiffres après la virgule.

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .
- Soit x' arrondi en précision $q < p$ de y .
- x' n'est pas toujours l'arrondi en précision q de x !
 - ① $x = 1.0110100101$. Arrondir à 7 chiffres après la virgule.
 - ② Après la décimale 7 de x on a : 101. L'arrondi y à 7 chiffres après la virgule de x est : 1.0110101

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .
- Soit x' arrondi en précision $q < p$ de y .
- x' n'est pas toujours l'arrondi en précision q de x !
 - ① $x = 1.0110100101$. Arrondir à 7 chiffres après la virgule.
 - ② Après la décimale 7 de x on a : 101. L'arrondi y à 7 chiffres après la virgule de x est : 1.0110101
 - ③ Après la décimale 4 de y : 101. L'arrondi x' à 4 décimales de y est : 1.0111

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .
- Soit x' arrondi en précision $q < p$ de y .
- x' n'est pas toujours l'arrondi en précision q de x !
 - ① $x = 1.0110100101$. Arrondir à 7 chiffres après la virgule.
 - ② Après la décimale 7 de x on a : 101. L'arrondi y à 7 chiffres après la virgule de x est : 1.0110101
 - ③ Après la décimale 4 de y : 101. L'arrondi x' à 4 décimales de y est : 1.0111
 - ④ MAIS, après la décimale 4 de x on a 100. suivant la règle de l'arrondi pair, L'arrondi z à 4 décimales de x est 1.0110.

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .
- Soit x' arrondi en précision $q < p$ de y .
- x' n'est pas toujours l'arrondi en précision q de x !
 - ① $x = 1.0110100101$. Arrondir à 7 chiffres après la virgule.
 - ② Après la décimale 7 de x on a : 101. L'arrondi y à 7 chiffres après la virgule de x est : 1.0110101
 - ③ Après la décimale 4 de y : 101. L'arrondi x' à 4 décimales de y est : 1.0111
 - ④ MAIS, après la décimale 4 de x on a 100. suivant la règle de l'arrondi pair, L'arrondi z à 4 décimales de x est 1.0110.
 - ⑤ Et on a $z \neq x'$!

Double arrondi

Cas d'examen de 3 bits après le dernier maintenu

- Soit x réel, y l'arrondi en précision p de x .
- Soit x' arrondi en précision $q < p$ de y .
- x' n'est pas toujours l'arrondi en précision q de x !
 - ① $x = 1.011010\mathbf{0}101$. Arrondir à 7 chiffres après la virgule.
 - ② Après la décimale 7 de x on a : 101. L'arrondi y à 7 chiffres après la virgule de x est : $1.011\mathbf{0}101$
 - ③ Après la décimale 4 de y : 101. L'arrondi x' à 4 décimales de y est : 1.0111
 - ④ MAIS, après la décimale 4 de x on a 100. suivant la règle de l'arrondi pair, L'arrondi z à 4 décimales de x est 1.0110.
 - ⑤ Et on a $z \neq x'$!
- On peut montrer que le problème du double arrondi n'intervient que si on choisit *l'arrondi au plus proche pair*. Pas de chance : c'est le mode d'arrondi le plus répandu !

Pourquoi privilégier l'arrondi au plus proche pair ?

- La méthode de « l'arrondi bancaire » (autre nom pour l'arrondi au plus proche pair) est employée pour éliminer le biais qui surviendrait en arrondissant à chaque fois par excès les nombres dont les trois derniers chiffres seraient 100.

Pourquoi privilégier l'arrondi au plus proche pair ?

- La méthode de « l'arrondi bancaire » (autre nom pour l'arrondi au plus proche pair) est employée pour éliminer le biais qui surviendrait en arrondissant à chaque fois par excès les nombres dont les trois derniers chiffres seraient 100.
- Transposons en base 10 : Imaginons une multinationale qui reçoit un milliards de virements exprimés en centimes d'euros (sur une certaine période) arrondis au dixième de centime sur un de ses comptes en banque.

Pourquoi privilégier l'arrondi au plus proche pair ?

- Supposons que pour un millième de ces virements, la partie fractionnaire soit de la forme $.x500$ où $x \in \llbracket 0, 9 \rrbracket$.

Pourquoi privilégier l'arrondi au plus proche pair ?

- Supposons que pour un millième de ces virements, la partie fractionnaire soit de la forme $.x500$ où $x \in \llbracket 0, 9 \rrbracket$.
- Si la banque arrondi le montant de ces virements au dixième de centime supérieur ($.x + 0.1$, puis répercussion de la retenue), la multinationale gagne 0.05 centime de plus par virement que ce qu'elle aurait dû toucher.
Au total cela fait $10^6 \times 5 \times 10^{-2} = 5 \times 10^4 = 50\,000$ euros que la multinationale a gagnés au détriment de la banque ! **Au centime inférieur, ce serait la banque qui gagnerait de l'argent.**

Pourquoi privilégier l'arrondi au plus proche pair ?

- Supposons que pour un millième de ces virements, la partie fractionnaire soit de la forme $.x500$ où $x \in \llbracket 0, 9 \rrbracket$.
- Si la banque arrondi le montant de ces virements au dixième de centime supérieur ($.x + 0.1$, puis répercussion de la retenue), la multinationale gagne 0.05 centime de plus par virement que ce qu'elle aurait dû toucher.
Au total cela fait $10^6 \times 5 \times 10^{-2} = 5 \times 10^4 = 50\,000$ euros que la multinationale a gagnés au détriment de la banque ! **Au centime inférieur, ce serait la banque qui gagnerait de l'argent.**
- D'où la nécessité d'arrondir certains montants au centime supérieur et d'autres au centime inférieur pour équilibrer, comme avec l'arrondi au plus proche pair.

Cas des exceptions

En cas de problème, la norme impose de signaler des *Exceptions*

- Diviser un nombre différent de 0 par 0 donne $\pm\infty$

Cas des exceptions

En cas de problème, la norme impose de signaler des *Exceptions*

- Diviser un nombre différent de 0 par 0 donne $\pm\infty$
- Diviser zéro par zéro, ou calculer le logarithme d'un nombre négatif conduisent à générer des **NaN** qu'on peut décider de considérer comme des exceptions.

Cas des exceptions

En cas de problème, la norme impose de signaler des *Exceptions*

- Diviser un nombre différent de 0 par 0 donne $\pm\infty$
- Diviser zéro par zéro, ou calculer le logarithme d'un nombre négatif conduisent à générer des **NaN** qu'on peut décider de considérer comme des exceptions.
- Nombre entier positif plus grand que le plus grand entier représentable (*overflow*). Ou plus petit que le plus petit entier représentable (*underflow*).

- 1 Présentation
- 2 La norme IEEE754
- 3 Exemples et règles d'arrondis
- 4 Problèmes induits par la norme
- 5 Arithmétique psychédélique

Attention aux test d'égalité

```
1 | # 1.2*.3.;;  
2 | - : float = 3.59999999999999964  
3 | # 0.1+.0.2;;  
4 | - : float = 0.300000000000000044
```

- Il est risqué d'écrire un programme avec des tests d'égalité entre flottants. Eviter `if x = y then ...`

Attention aux test d'égalité

```
1 | # 1.2*.3.;;  
2 | - : float = 3.59999999999999964  
3 | # 0.1+.0.2;;  
4 | - : float = 0.3000000000000000044
```

- Il est risqué d'écrire un programme avec des tests d'égalité entre flottants. Eviter `if x = y then ...`
- Il vaut mieux utiliser `if abs_float (x-y) < eps`. Au format double précision, prendre $\varepsilon = 10^{-12}$ pour comparer des flottants d'ordre de grandeur 1 est tout à fait raisonnable.

Arrondi d'un calcul, commutativité, associativité

- L'arrondi de la somme n'est pas la somme des arrondis.
Comparer **0.1+0.2** et **0.3**

Arrondi d'un calcul, commutativité, associativité

- L'arrondi de la somme n'est pas la somme des arrondis.
Comparer **0.1+0.2** et **0.3**
- La multiplication et l'addition restent commutatives

Arrondi d'un calcul, commutativité, associativité

- L'arrondi de la somme n'est pas la somme des arrondis.
Comparer **0.1+0.2** et **0.3**
- La multiplication et l'addition restent commutatives
- L'associativité se perd

```
1 | # let a = 9e307 and b = 9e307 and c = -2e306 in (a
   |   +.b)+.c;;
2 | - : float = infinity
3 | # let a = 9e307 and b = 9e307 and c = -2e306 in a
   |   +.(b+.c));;
4 | - : float = 1.780000000000000023e+308
```

Distributivité

La distributivité se perd

```
1 | # let a = 1e-10 and b = 9e307 and c = 9e307
2 |   in a *. (b +. c);;
3 | - : float = infinity
4 | # let a = 1e-10 and b = 9e307 and c = 9e307 in
5 |   a *. b +. a *. c;;
6 | - : float = 1.800000000000000021e+298
```

Relation d'ordre

Dans \mathbb{R} , l'addition est compatible avec la relation d'ordre : si $a \leq b$ alors $a + c \leq b + c$.

Nous avons vu que la représentation en machine de **0.1** est strictement supérieure à **0.1**. En toute logique additionner 10 fois **0.1** devrait donner un résultat plus grand que 1 :

```
1 | # let rec add n x = if n = 1 then x
2 |   else add (n-1) (x +. 0.1)
3 |   in add 10 0.1;;
4 | - : float = 0.9999999999999999889
```

Convergence et divergence

Dans le cours de maths, la série harmonique $\sum_{k \geq 1} \frac{1}{k}$ diverge quel que soit l'ordre des opérations. De plus, du fait de la commutativité, les sommes partielles $\sum_{k=1}^n \frac{1}{k}$ et $\sum_{k=1}^n \frac{1}{n-k+1}$ sont égales.

Avec l'ordre décroissant

```
1 # let rec harmonique_dec n x = if n = 0 then x
2   else harmonique_dec (n-1) (1./float_of_int n +.x)
3   ;;
4   val harmonique_dec : int -> float -> float = <fun>
5 # harmonique_dec (int_of_float 1e5) 0.;;
6 - : float = 12.0901461298634079
```

Dernières décimales 4079

Convergence et divergence

Et l'ordre croissant

```
1 | # let rec harmonique_cr n m x = if n = m+1 then x
2 |   else harmonique_cr (n+1) m (1. /. float_of_int n +.x
   |   );;
3 |   val harmonique_cr : int -> int -> float -> float =
   |   <fun>
4 | # harmonique_cr 1 (int_of_float 1e5) 0.;;
5 | - : float = 12.090146129863335
```

- Les dernières décimales des deux résultats ne sont pas égales alors même que l'addition est commutative...

Convergence et divergence

Et l'ordre croissant

```
1 | # let rec harmonique_cr n m x = if n = m+1 then x
2 |   else harmonique_cr (n+1) m (1. /. float_of_int n +.x
   |   );;
3 |   val harmonique_cr : int -> int -> float -> float =
   |   <fun>
4 | # harmonique_cr 1 (int_of_float 1e5) 0.;;
5 | - : float = 12.090146129863335
```

- Les dernières décimales des deux résultats ne sont pas égales alors même que l'addition est commutative...
- De plus cette série converge pour les nombres flottants. Après un certain rang, $\frac{1}{n}$ devient plus petit que le plus petit nombre dénormalisé et est compté comme 0.

Infinis et NaN

```
1 | # 5./0., -5./0.;;  
2 | - : float * float = (infinity, neg_infinity)  
3 | # sqrt (-1.);;  
4 | - : float = nan  
5 | # max_float;;  
6 | - : float = 1.79769313486231571e+308  
7 | # max_float +. 10.;;  
8 | - : float = 1.79769313486231571e+308  
9 | # max_float +. max_float;;  
10 | - : float = infinity
```

Écart avec le successeur

- L'écart entre x et son successeur est croissant avec x

```
1 # 1. +. 1e-16, 1. +. 1e-15;;  
2 - : float * float = (1., 1.000000000000000111)  
3 # 1e16 +. 1., 1e16 +. 10.;;  
4 - : float * float = (1e+16, 10000000000000010.)  
5 # min_float;;  
6 - : float = 2.22507385850720138e-308
```

Écart avec le successeur

- L'écart entre x et son successeur est croissant avec x

```
1 # 1. +. 1e-16, 1. +. 1e-15;;  
2 - : float * float = (1., 1.000000000000000111)  
3 # 1e16 +. 1., 1e16 +. 10.;;  
4 - : float * float = (1e+16, 10000000000000010.)  
5 # min_float;;  
6 - : float = 2.22507385850720138e-308
```

- Ainsi, si on veut de la précision dans des calculs faisant intervenir de grands nombres, on a intérêt à travailler avec leurs inverses (qui seront sans doute des nombres dénormalisés).

Écart avec le successeur

- L'écart entre x et son successeur est croissant avec x

```
1 # 1. +. 1e-16, 1. +. 1e-15;;  
2 - : float * float = (1., 1.000000000000000111)  
3 # 1e16 +. 1., 1e16 +. 10.;;  
4 - : float * float = (1e+16, 10000000000000010.)  
5 # min_float;;  
6 - : float = 2.22507385850720138e-308
```

- Ainsi, si on veut de la précision dans des calculs faisant intervenir de grands nombres, on a intérêt à travailler avec leurs inverses (qui seront sans doute des nombres dénormalisés).
- Ajoutons que, suivant les machines, les opérations sur les nombres dénormalisés sont effectués au niveau matériel ou au niveau logiciel. Cela a une importance en terme de rapidité de calcul