La précision numérique de calculs

<u>Ou</u>

Pourquoi il ne faut pas trop faire confiance à une machine







Sleipner A Patriot Ariane 5

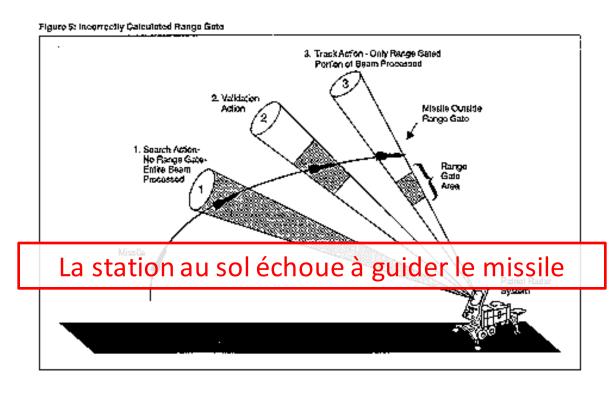
1991 : le Missile Patriot américain



- Première guerre du Golfe
- Missile pour contrer les SCUD irakiens
- 3 M\$ par unité

1991 : le Missile Patriot américain





- Première guerre du Golfe
- Missile pour contrer les SCUD irakiens
- 3 M\$ par unité

- Un Patriot rate un SCUD
- ♦ 28 soldats tués
- 98 soldats blessés

Un ordinateur ne sait pas vraiment calculer

Exemples

x = 1 + 1?

Un ordinateur ne sait pas vraiment calculer

$$x = 1 + 1$$

$$\bullet$$
 x = 1 + 100000000?

Un ordinateur ne sait pas vraiment calculer

$$x = 1 + 1$$

$$\bullet$$
 x = 1 + 100000000

Un ordinateur ne sait pas vraiment calculer

$$x = 1 + 1$$

$$\bullet$$
 x = 1 + 100000000

$$\bullet$$
 x = 1.0 + 1000000000.0 ?

Un ordinateur ne sait pas vraiment calculer

$$x = 1 + 1$$

$$\bullet$$
 x = 1 + 100000000

$$\bullet$$
 x = 1.0 + 100000000.0

Un autre exemple

Un incrément

```
float compteur = 0.0f;
float increment = 1.0f;
for (int i = 0 ; i < 1000000000 ; i++){
   compteur = compteur + increment;
}
print(compteur);</pre>
```

Un autre exemple

Un incrément

```
float compteur = 0.0f;
float increment = 1.0f;
for (int i = 0 ; i < 1000000000 ; i++){
   compteur = compteur + increment;
}
print(compteur);</pre>
```

\$> 16777216.0

16 Millions au lieu de 100 Millions!

Calculer ou ne pas calculer

- Le système de numération utilisé
- Machine finie vs. infinité de nombre

- Représentation binaire sur ordinateur
- Format = notation scientifique

$$s \times m \times 2^e$$

$$s = \{+, -\}$$

m = mantisse

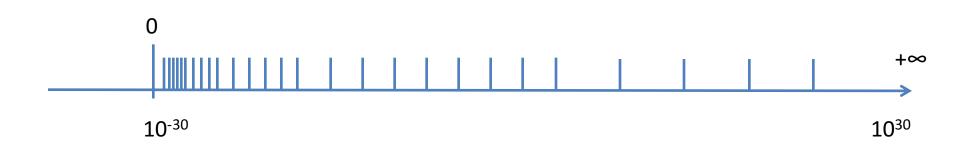
e = exposant

- Les nombres flottants
 - base 2
 - 32 bits (float), 64 bits (double)

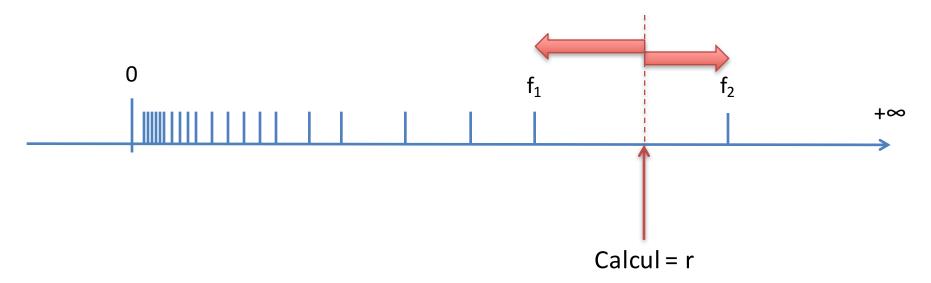
64 bits			
signe	mantisse	exposant	
1 bit	52 bits	11 bits	

32 bits			
signe	mantisse	exposant	
1 bit	23 bits	8 bits	

- Idée sous-jacente
 - Infinité de nombre à représenter
 - Étaler les flottants sur l'ensemble du spectre
- Répartition
 - ◆ Très petit nombre, ex: 10⁻³⁰
 - Très grand nombre, ex: 10³⁰



Arrondis



Mode d'arrondi

- Plus proche
- Vers zéro
- Vers +∞
- ♦ Vers -∞

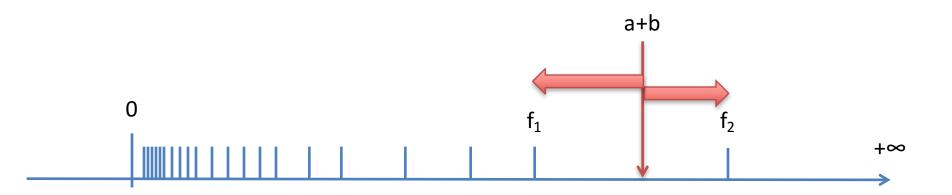
$$arrondi(r) = f_2$$

$$arrondi(r) = f_1$$

$$arrondi(r) = f_2$$

$$arrondi(r) = f_1$$

Notion d'arrondi correct

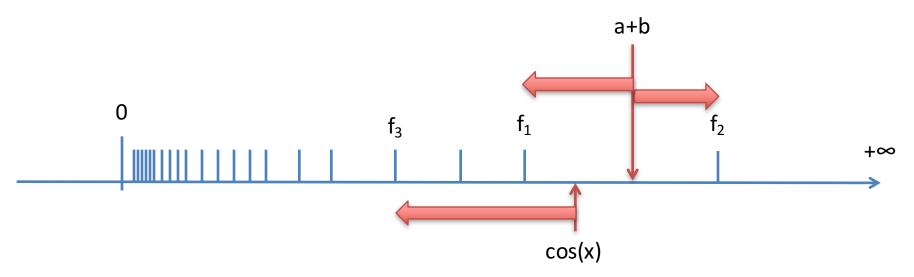


Si mon calcul est à l'arrondi correct alors arrondi (calcul) vaut soit f₁ soit f₂

Les opérateurs à arrondis correctes sont :

- **+**
- **-**
- ×
- **♦** ÷
- **♦** √

Notion d'arrondi pas toujours correct



Si mon calcul est à l'arrondi correct alors arrondi (calcul) vaut soit f₁ soit f₂

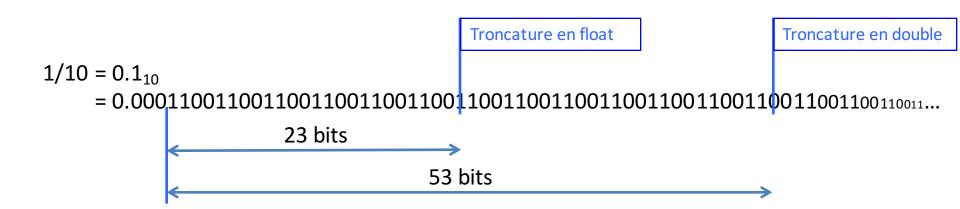
Les opérateurs à arrondis correctes sont :

- +
- **-**
- ***** ×
- **+** :
- 1

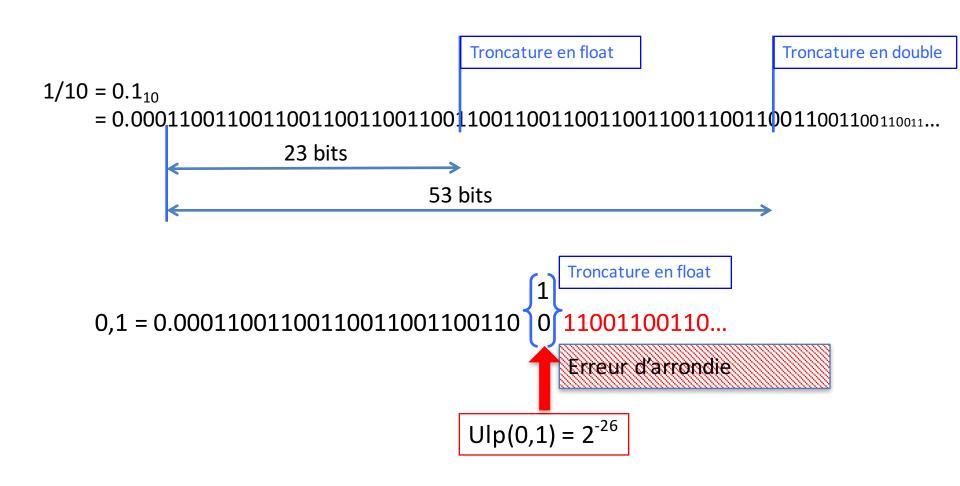
Ne sont **pas** arrondis correctement :

- cos
- ♦ sin
- tan
- log
- exp
- •

Exemples typiques d'arrondi qui tuent



Exemples typiques d'arrondi qui tuent



 $0.1 \approx 0.10000001490116119384765625$ en arrondi au plus près

La valeur 0,1

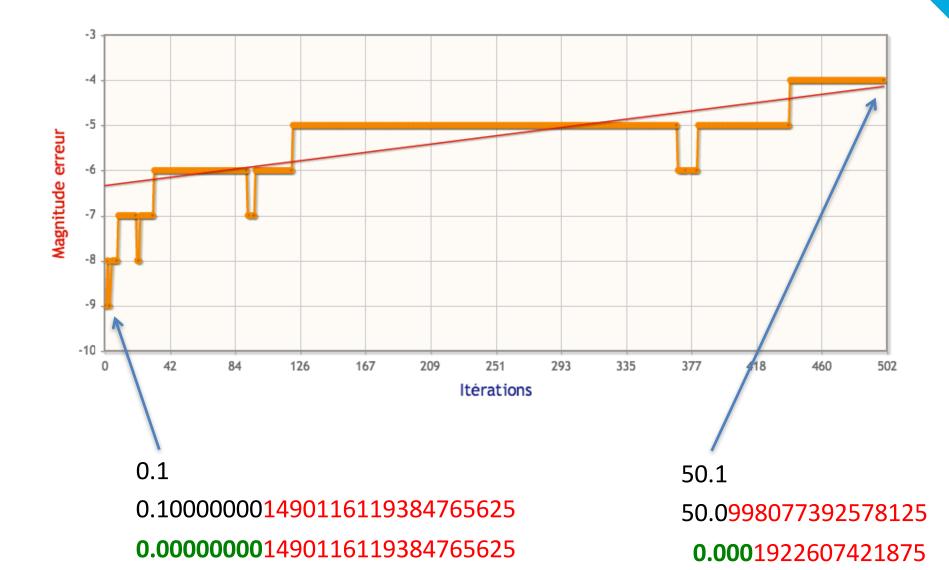
Missile Patriot

Le missile a une horloge en 10^{ème} de seconde

- + 0.1 en binaire c'est 0.10000001490116119384765625
- + 100h de fonctionnement où l'écart s'accumule
- = 600 mètres de déviation

```
float compteur = 0.0f;
float increment = 0.1f;
while (missile is on) {
  compteur = compteur + increment;
}
```





- Exemples typiques d'arrondi qui tuent
 - L'absorption

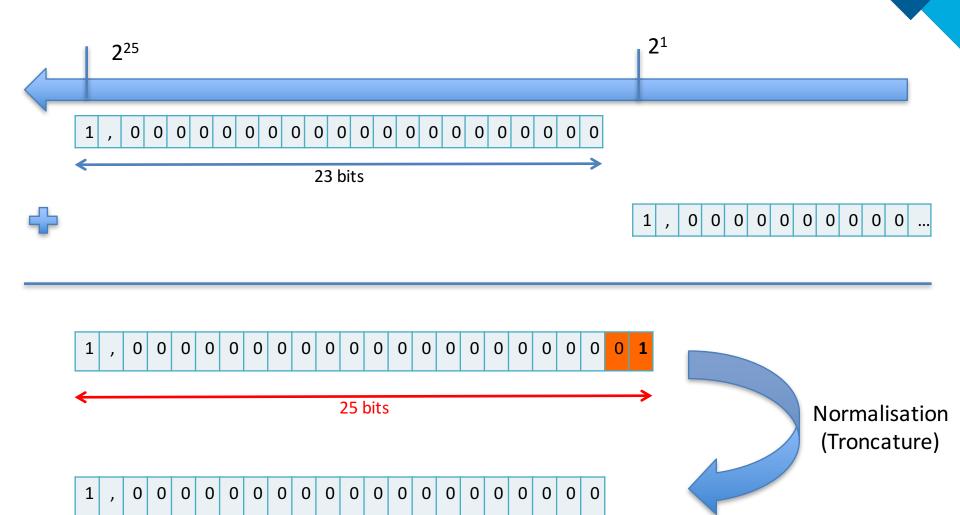
Pourquoi 1.0 + 100 000 000.0 ≠ 100 000 001.0 ?

- Exemples typiques d'arrondi qui tuent
 - L'absorption

Pourquoi 1.0 + 100 000 000.0 ≠ 100 000 001.0 ?

Les petits et les grands s'additionnent mal

◆ Pourquoi?



23 bits

Et si j'utilise le type double?

- Exemples typiques d'arrondi qui tuent
 - L'annulation catastrophique

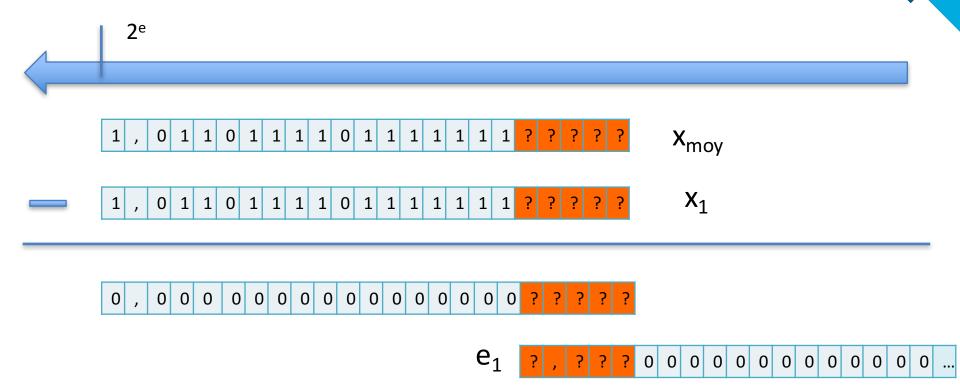
Comment deux calculs un peu faux font un résultat aléatoire

- Cas concret
 - 2 capteurs redondants
 - Calcul de la moyenne
 - Calcul de l'écart à la moyenne

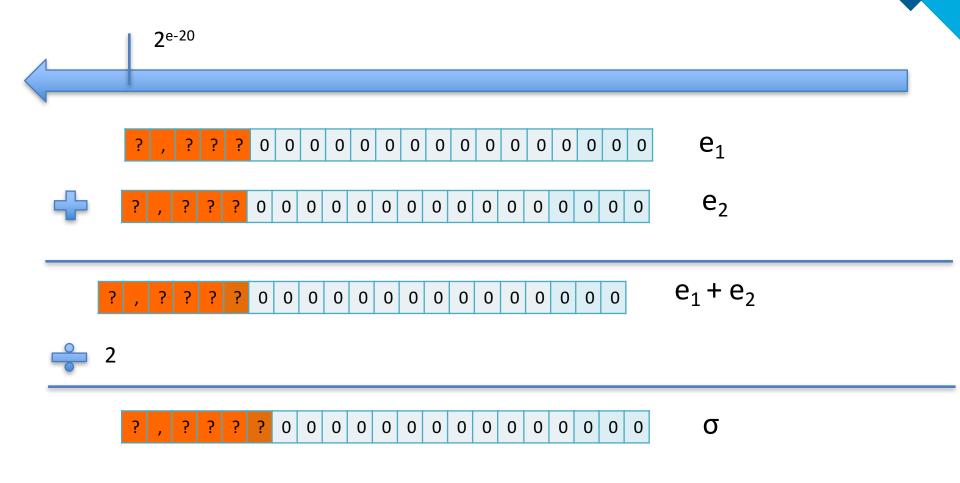
- Valeurs bruités mais proches entre elles
 - Les valeurs (x₁ et x₂) sont entachés d'erreurs
 - Mais leurs valeurs sont proches

$$x_{moy} = (x_1 + x_2)/2 \approx x_1 \approx x_2$$

- Écart type
 - \bullet e₁ = $(x_{mov} x_1)^2$
 - \bullet e₂ = $(x_{mov} x_2)^2$
 - \bullet $\sigma = (e_1 e_2) / 2$
- Pourquoi ce calcul risque de mal se terminer ?

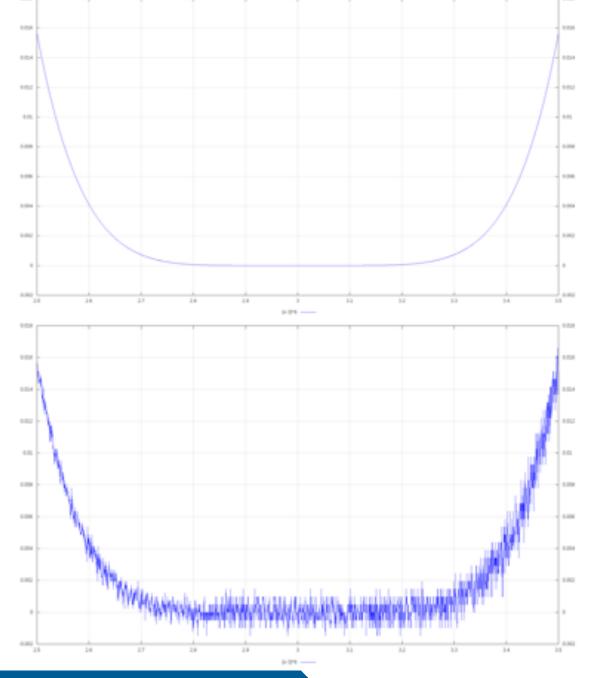


On a la même chose pour **e**₂



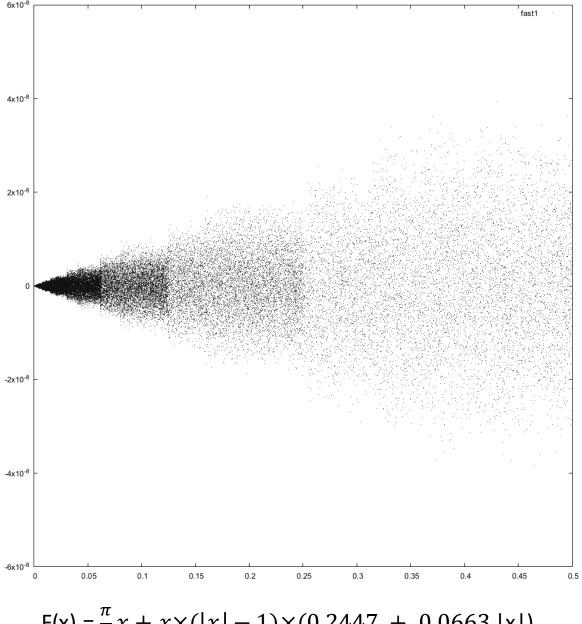
Et si je multiplie σ par un <u>très gros</u> coefficient?

 $F(x) = (x-3)^6$

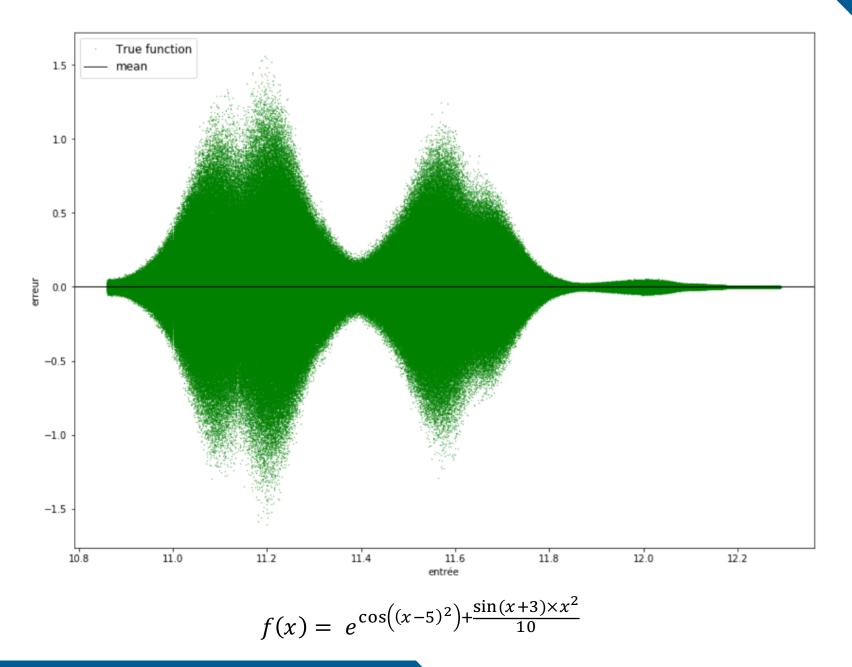


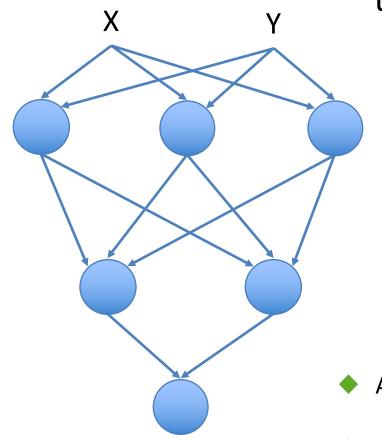
Flottant

Réel



$$F(x) = \frac{\pi}{4}x + x \times (|x| - 1) \times (0.2447 + 0.0663 |x|)$$





```
tanh( p_{13} * tanh(p_7 * tanh(p_1 * X + p_2 * Y) + p_8 * tanh(p_3 * X + p_4 * Y) + p_9 * tanh(p_5 * X + p_6 * Y))
```

```
+ p_{14} * tanh (p_{10} * tanh(p_1 * X + p_2 * Y) + p_{11} * tanh(p_3 * X + p_4 * Y) + p_{12} * tanh(p_5 * X + p_6 * Y))
```

- Architecte non connu (GPU, CPU, PFGA ? Mix ?)
- Précision non connue (16, 32, 64 bits ? Mix ?)
- Poids non connus (online learning?)

- Architecture variable
- Librairie avec performances différentes
- Arrondis des valeurs dû aux calculs

Reproduire des résultats n'est pas simple!

- L'industrie commence à regarder tout ça
 - Aéronautique, aérospatiale, énergie, défense
 - Finance
 - ◆ IA
- Notion de précision garantie
 - Type des variables et manière d'écrire le code
 - Entrées plus précises
- Equilibre entre
 - Performance (temps, consommation, occupation)
 - Précision (stabilité, durabilité)

Alors que faire ?

- Utiliser des types les plus grands ?
 - Perte importante de performances
- Faire des simulations?
 - Nombre exponentiel de cas à traiter
 - Très peu de cas pathologiques
- Avis d'expert ?
 - Très rares
- Outils de preuves de programmes ?
 - Pas forcément très simples et user-friendly

numalis

```
float s = 0.20, t = 1.00;
float risk = 0.05, divd = 0,01;
float S = 60.00, K = initK();
assert(K >= 65.00 \&\& K <= 70.00);
float d1 = (\log(S/K) + (risk-divd+(s*s)/2.0)*t) / (s*sqrt(t));
float d2 = d1 - s * sqrt(t);
float call = S*exp(-divd*t) * N(d1) - K*exp(-risk*t) * N(d2);
    Write it more like this:
    \exp(-\text{divd}*t) * (S*N(d1) - \exp((\text{divd}-\text{risk})*t)*N(d2)*K)
```