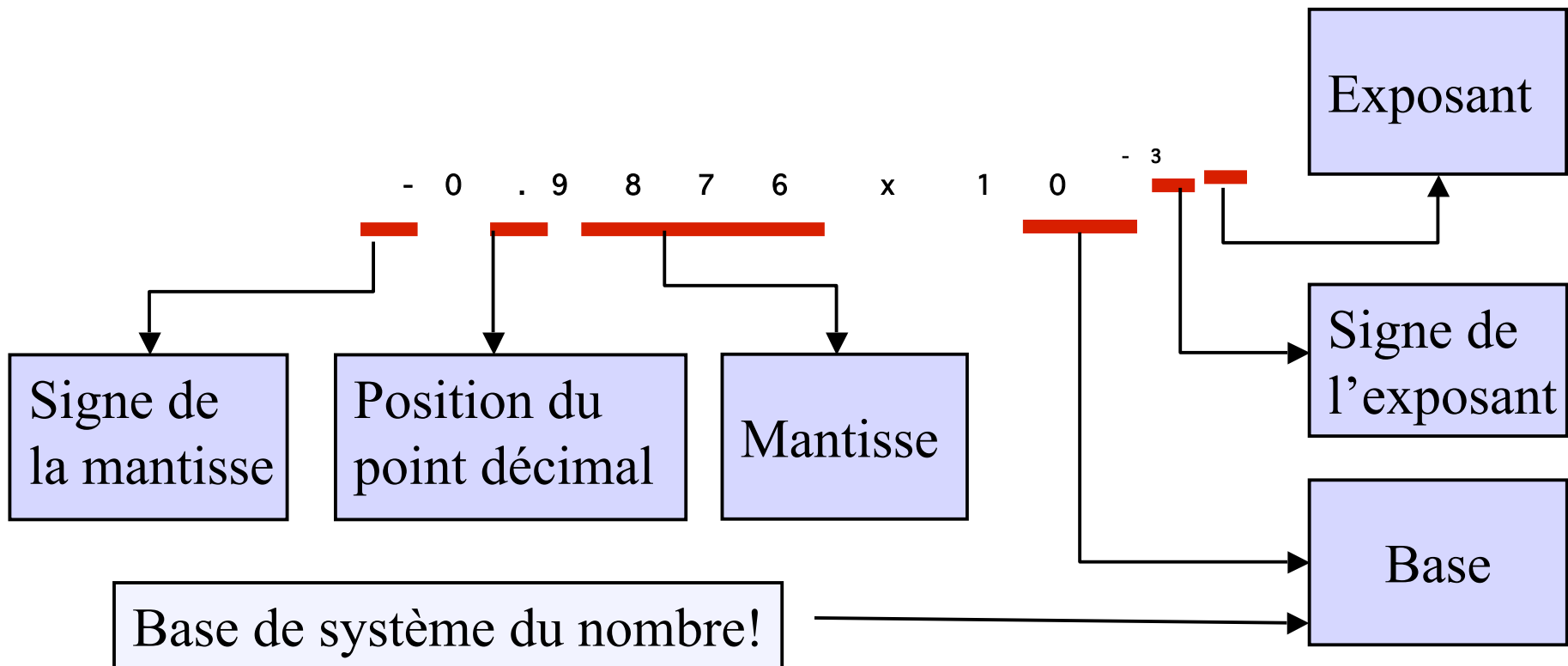


# Représentation des nombres flottants



# Éléments de la notation exponentielle



# Représentation normalisée

- Un nombre représenté en virgule flottante est normalisé s'il est sous la forme:
  - $\pm 0, M * X^{\pm c}$
  - $M$  – un nombre dont le premier chiffre est non nul
- Exemple:
  - $+ 59,4151 * 10^{-5} \Rightarrow$

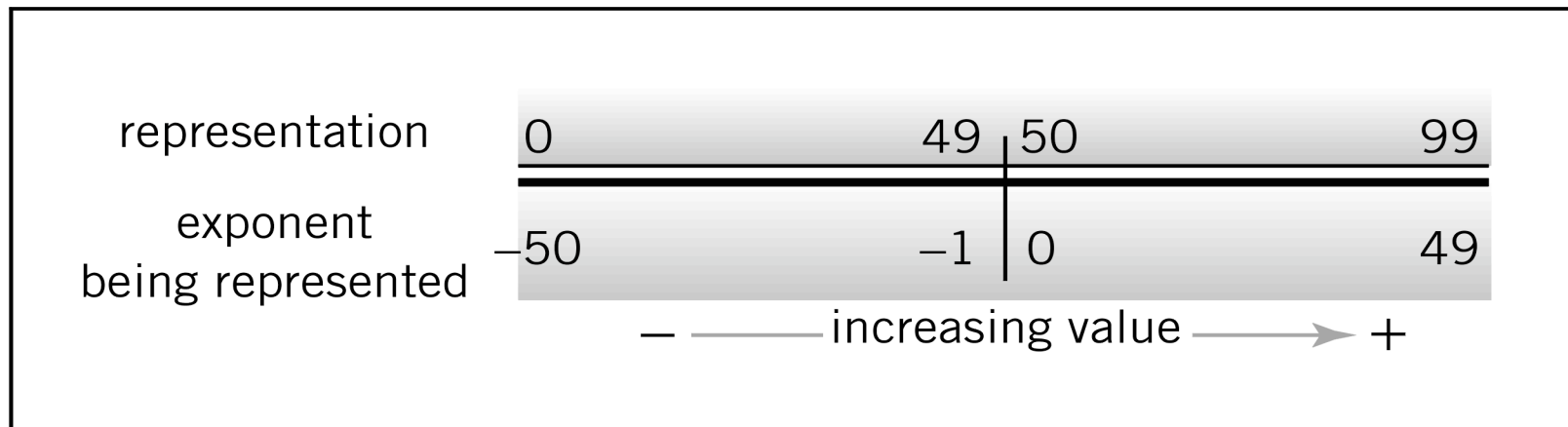
Normalisé:  $+0,594151 * 10^{-3}$

# Représentation de l'exposant et de son signe

- L'exposant est translatée de manière à toujours coder en interne une valeur positive
- Avec 2 digits réservés au codage de l'exposant
  - Les valeurs positives:  $[+0, +99]$
  - En appliquant une translation  $k=50$ :
    - Les exposants représentables  $\Rightarrow [-50,49]$
- La constante  $k$  est appelée constante d'excentrement

# Représentation en virgule flottante

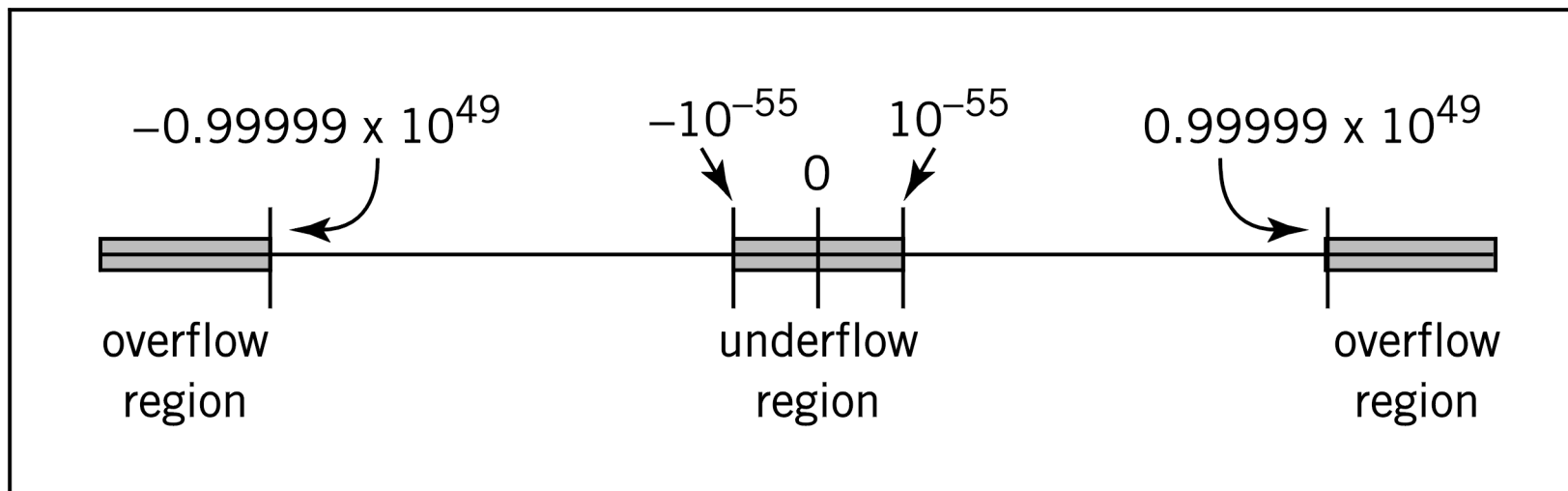
- Avec 2 digits réservés au codage de l'exposant avec un excentrement égal à  $50_{10}$  et 5 digits pour la mantisse on peut représenter
  - de  $.00001 \times 10^{-50}$  à  $.99999 \times 10^{49}$



Englander: The Architecture of Computer  
Hardware and Systems Software, 2nd edition  
Chapter 5, Figure 05-01

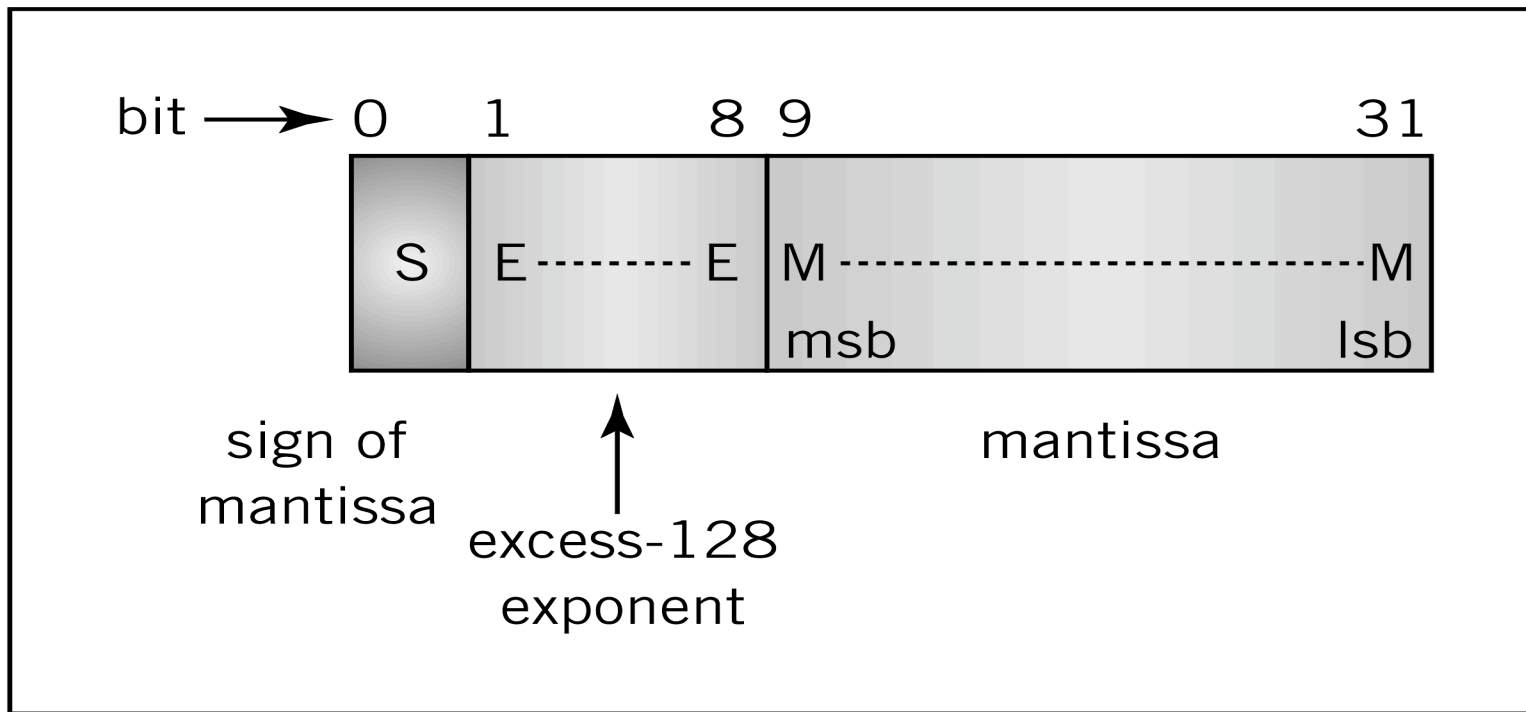
# Overflows / Underflows

- De  $.00001 \times 10^{-50}$  à  $.99999 \times 10^{49}$   
 $1 \times 10^{-55}$  à  $.99999 \times 10^{49}$



Englander: The Architecture of Computer  
Hardware and Systems Software, 2nd edition  
Chapter 5, Figure 05-02

# Format typique



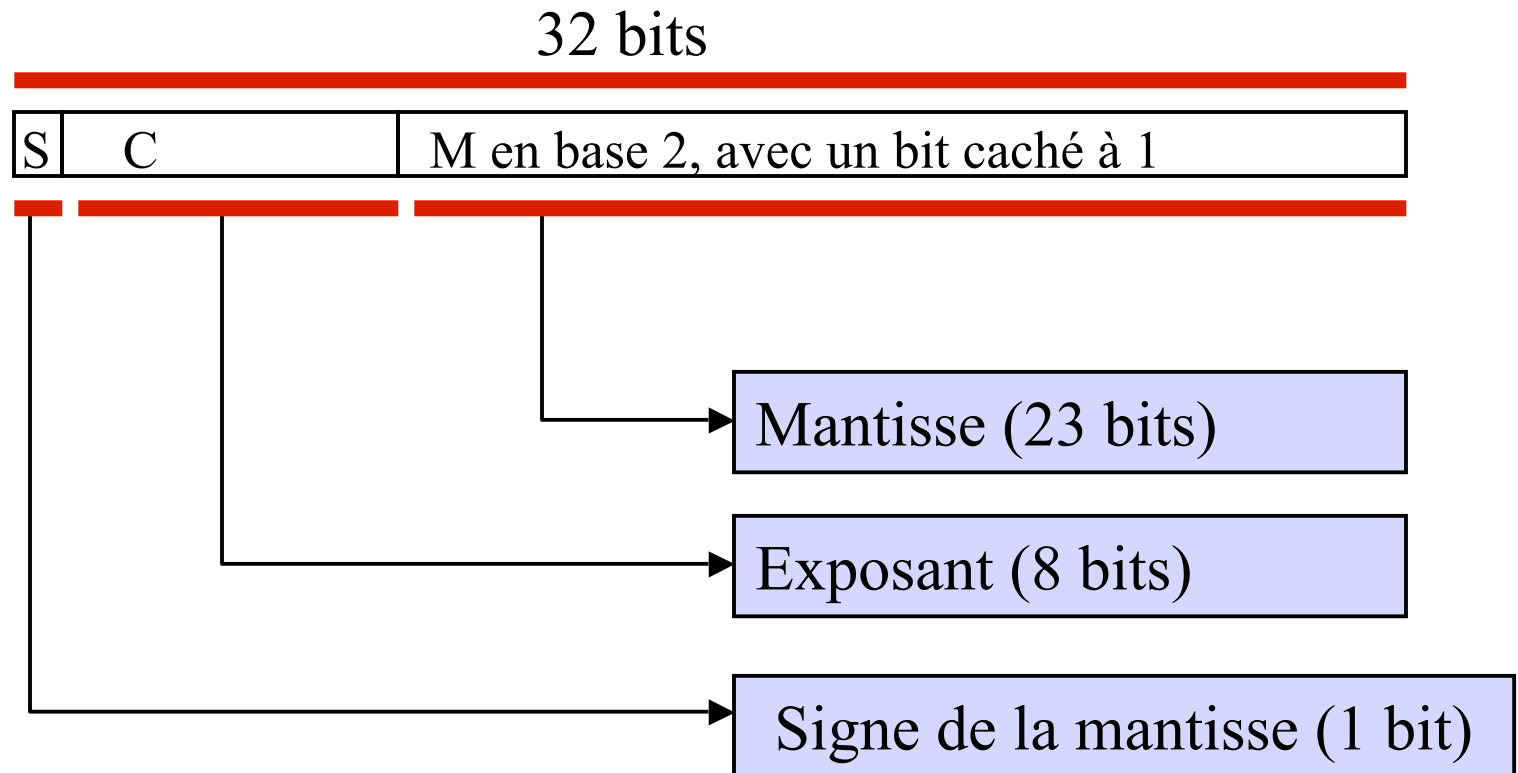
Englander: The Architecture of Computer Hardware and Systems Software, 2nd edition  
Chapter 5, Figure 05-04



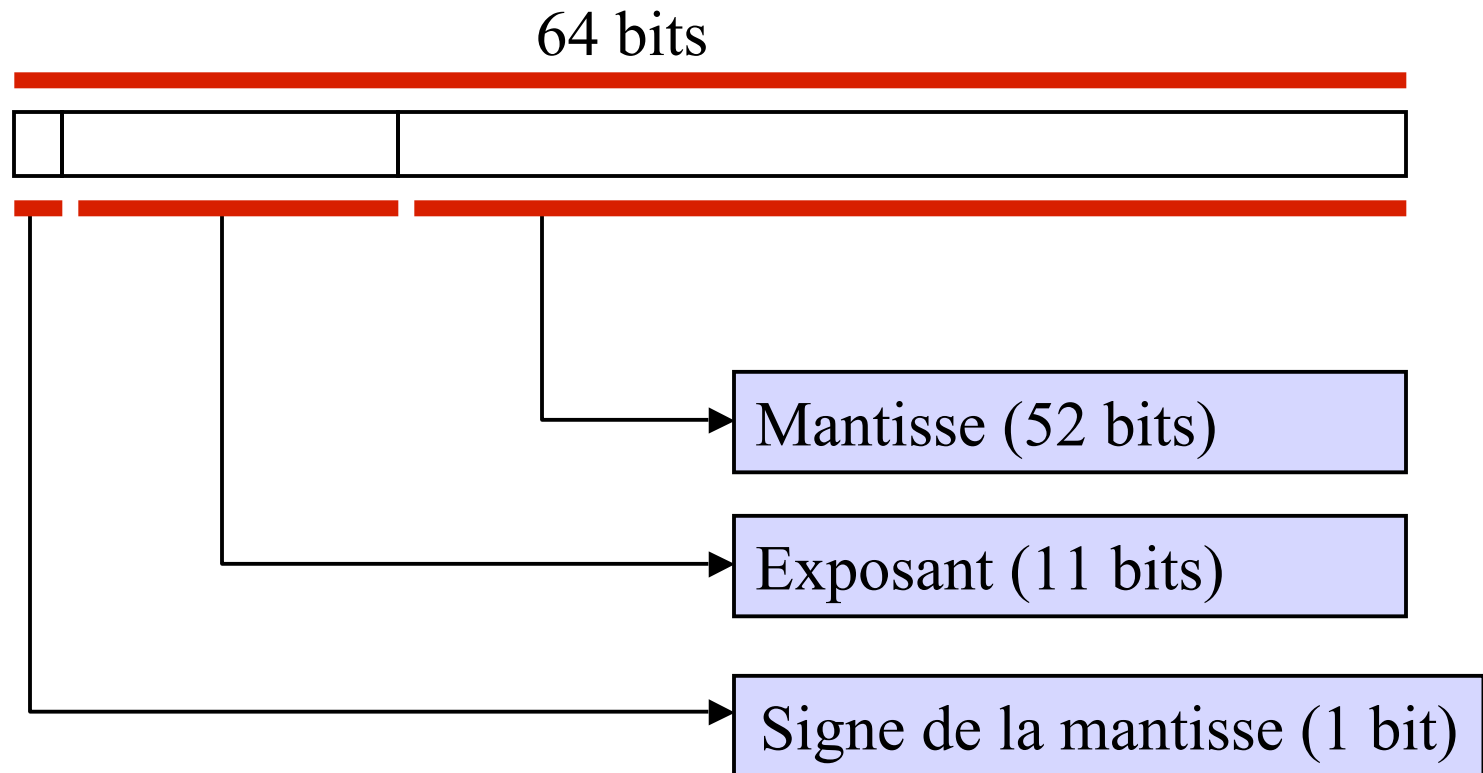
# La norme IEEE 754

- Un format standardisé
- **Format simple précision**: 32 bits
  - Bit du signe (1 bit)
  - Exposant (8 bits)
  - Mantisse (23 bits)
- **Format double précision**: 64 bits
  - Bit du signe (1 bit)
  - Exposant (11 bits)
  - Mantisse (52 bits)

# Format simple précision



# Format Double Précision



# Normalisation dans le format IEEE 754

- La mantisse est normalisé sous la forme
  - $\pm 1, M * 2^{\pm c}$
  - Pseudo mantisse
  - Le 1 précédant la virgule n'est pas codé en machine et est appelé bit caché
- Exemple:
  - Mantisse:
  - Représentation:  $1\ 0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$

$$1.101_2 = 1.625_{10}$$

# IEEE 754, Représentation de l'exposant

- Constante  $k$  d'excentrement appliquée à l'exposant
  - Simple précision:  $+127_{10}$
  - Double précision:  $+1023_{10}$
- L'exposant  $c$  codé en interne
  - $\pm c + 127_{10}$
  - $\pm c + 1023_{10}$
- Ex.,  $-k = 127_{10}$ ,
  - Exposant: 1 0 0 0 0 1 1 1 2
  - Représentation: 1 3 5 <sub>10</sub> - 1 2 7 <sub>10</sub> - 8 <sub>10</sub> ( v a l e u r )

# Représentation de l'exposant et de son signe

- Exemple -

Représentez l'exposant  $14_{10}$  avec un excentrement 127:

$$127_{10} = + 01111111_2$$

$$14_{10} = + \underline{00001110}_2$$

$$\text{Représentation} = 10001101_2$$

# Représentation de l'exposant et de son signe

- Exemple -

Représentez l'exposant  $-8_{10}$  avec un excentrement 127:

$$\begin{array}{rcl} 127_{10} & = & + \mathbf{01111111}_2 \\ - 8_{10} & = & - \mathbf{\underline{00001000}}_2 \\ \text{Représentation} & = & \mathbf{01110111}_2 \end{array}$$





# Exercice – Conversion en virgule flottante IEEE 754

- Quelle est la valeur décimale des représentations internes suivantes?

1 1 0 0 0 0 0 1 0  
1 1 1 1 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

- Réponse: -

# Exercice – Conversion en virgule flottante IEEE 754

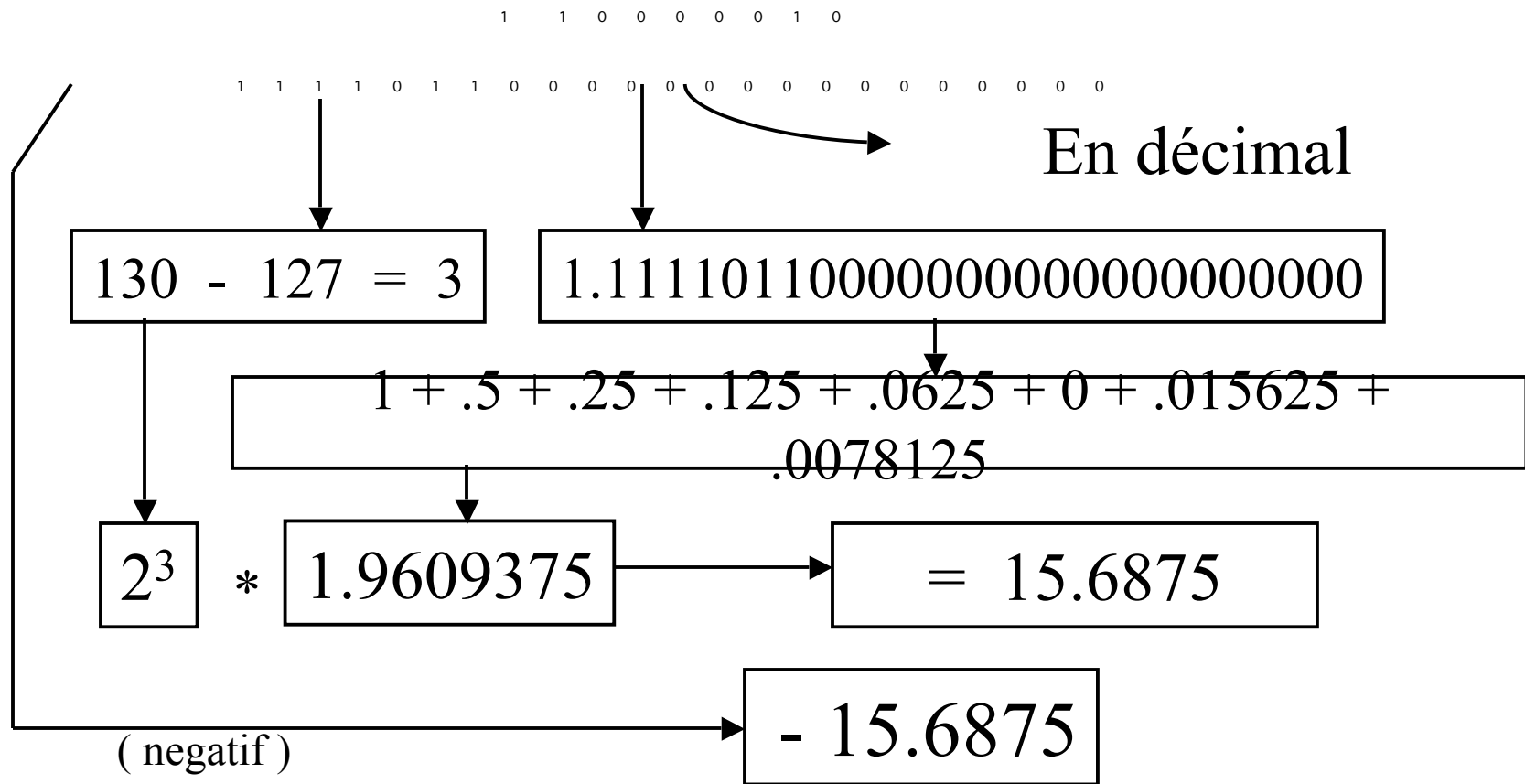
Réponse

- Quelle est la valeur décimale des représentations internes suivantes?

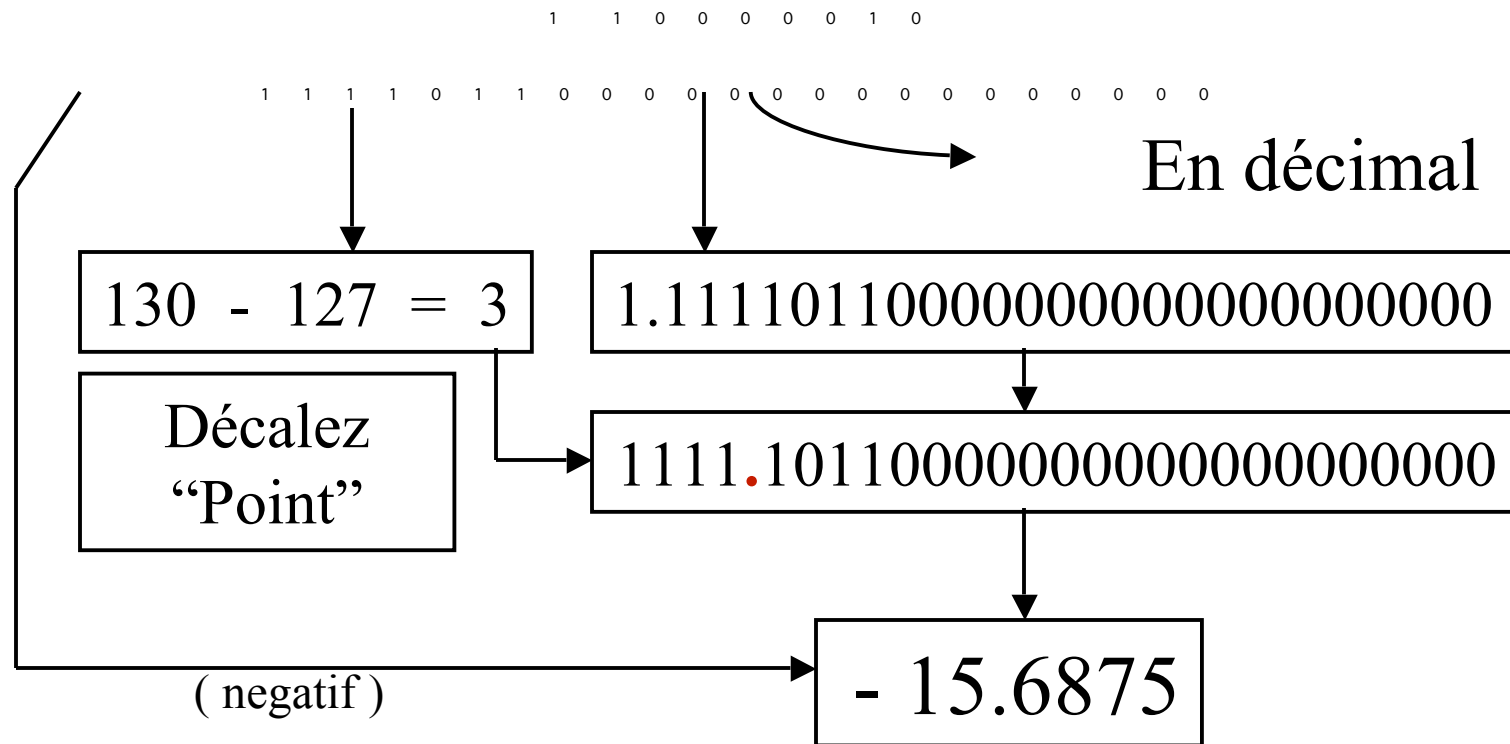
1 1 0 0 0 0 0 1 0  
1 1 1 1 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

- Réponse: -15.6875

# Solution



## Solution : Méthode Alternative



# Exercice – Conversion en virgule flottante IEEE 754

- Quelle est la représentation interne du nombre  $3.14_{10}$ ?
- Remarque: utiliser seulement les 10 chiffres significatifs pour la mantisse
- Réponse: -



## Solution : 3.14 en IEEE Simple Précision

---

3.14 En Binaire (approx):

- Normalisez ( $2^1$ )
- Enlevez le bit caché

11.001000111101

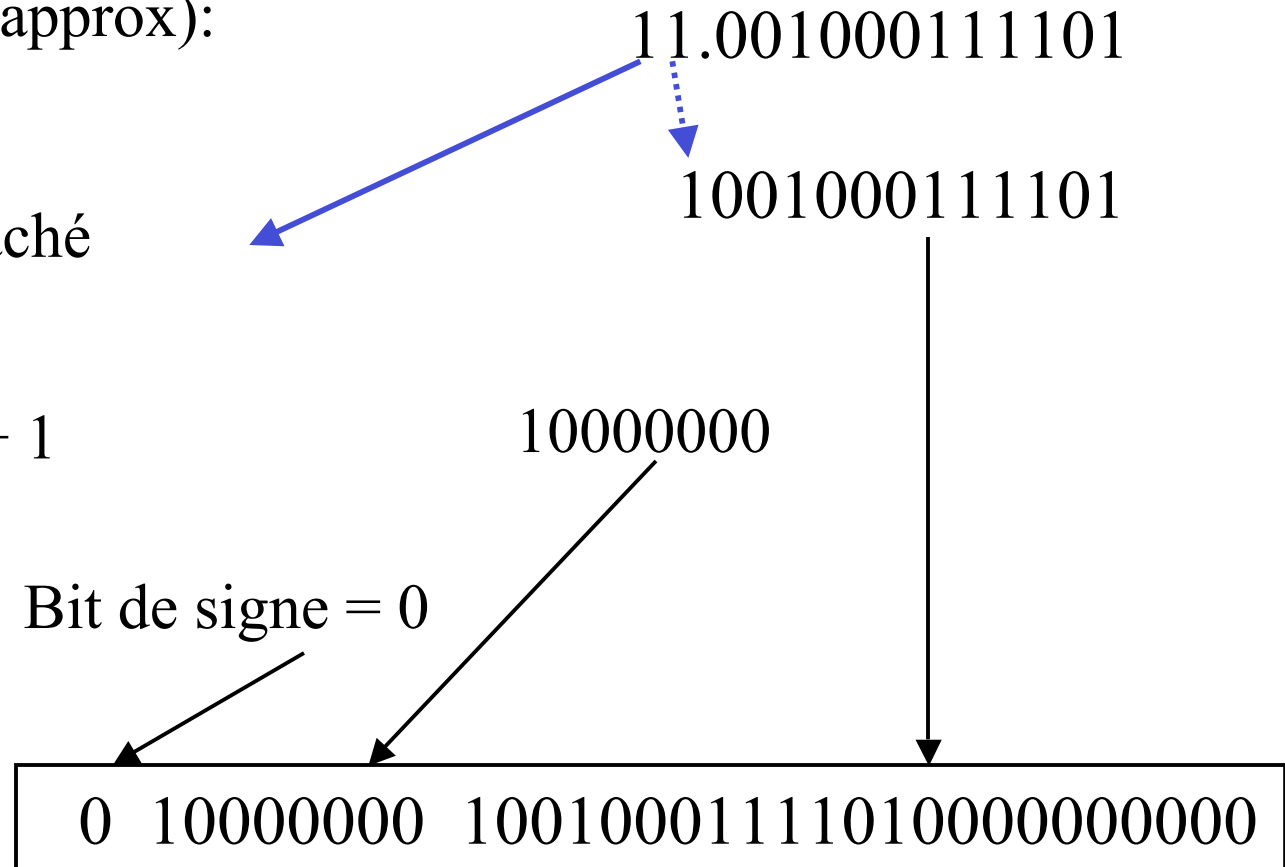
1001000111101

Exposant =  $127 + 1$

10000000

Valeur est positive: Bit de signe = 0

0 10000000 100100011110100000000000







# Représentation du zéro, des infinis, représentations dénormalisées

- Le norme IEEE admet des codages spéciaux pour la représentation
  - 0
  - $+\infty$
  - $-\infty$
  - Représentations dénormalisées

# Représentation du zéro, des infinis, représentations dénormalisées

Exposant	Mantisse	Valeur
0	$\pm 0$	0
0	Non 0	$\pm 2^{-126} * 0.M$
-126 - +127	Tout	$\pm 2^{E+127} * 1.M$
$\pm 128$	$\pm 0$	$\pm \infty$
$\pm 128$	Non 0	Conditions spéciales

# Addition et soustraction de deux nombres décimales en virgule flottante

Opérandes	Alignement	Normaliser et arrondir
$6.144 \times 10^2$	$0.06144 \times 10^4$	$1.003644 \times 10^5$
$+9.975 \times 10^4$	$+9.975 \times 10^4$	$+ .0005 \times 10^5$
	$10.03644 \times 10^4$	$1.004 \times 10^5$

Opérandes	Alignement	Normaliser et arrondir
$1.076 \times 10^{-7}$	$1.076 \times 10^{-7}$	$7.7300 \times 10^{-9}$
$-9.987 \times 10^{-8}$	$-0.9987 \times 10^{-7}$	$+ .0005 \times 10^{-9}$
	$0.0773 \times 10^{-7}$	$7.730 \times 10^{-9}$

# Calcul en virgule flottante: Addition

- Nombres doivent être alignés : avoir les mêmes exposants (le plus élevé pour protéger la précision)
- Additionner mantisses. Si overflow, ajuster l'exposant
- Ex. 0 51 99718 (e = 1) et 0 49 67000 (e = -1)

- Aligner les nombres: 
$$\begin{array}{r} 0\ 51\ 99718 \\ 0\ 51\ 00670 \end{array}$$

- Additionner: 
$$\begin{array}{r} 99718 \\ + 00670 \\ \hline \underline{1\ 00388} \end{array} \quad \leftarrow \text{Overflow}$$

- Arrondir le nombre et ajuster l'exposant: 0 52 10039

# Calcul en virgule flottante: Multiplication

- $(a * 10^e) * (b * 10^f) = a * b * 10^{e+f}$
- Règle: multiplier les mantisses; additionner les exposants

But: Codage en excédent,  $(n + e) + (n + f) = 2 * n + e + f$

→ Besoin soustraire constante d'excentrement  $n$  a partir du résultat

- Ex. 0 51 99718 ( $e = 1$ ) and 0 49 67000 ( $e = -1$ )  
Mantisses:  $.99718 * .67000 = 0.6681106$   
Exposants:  $51 + 49 = 100$  and  $100 - 50 = 50$   
Normaliser:  $.6681106 \rightarrow .66811$   
Résultat:  $.66811 * 10^0$  (50 signifie  $e = 0$ )