

## 1.2 IEEE standard 754 : FLOATING POINT

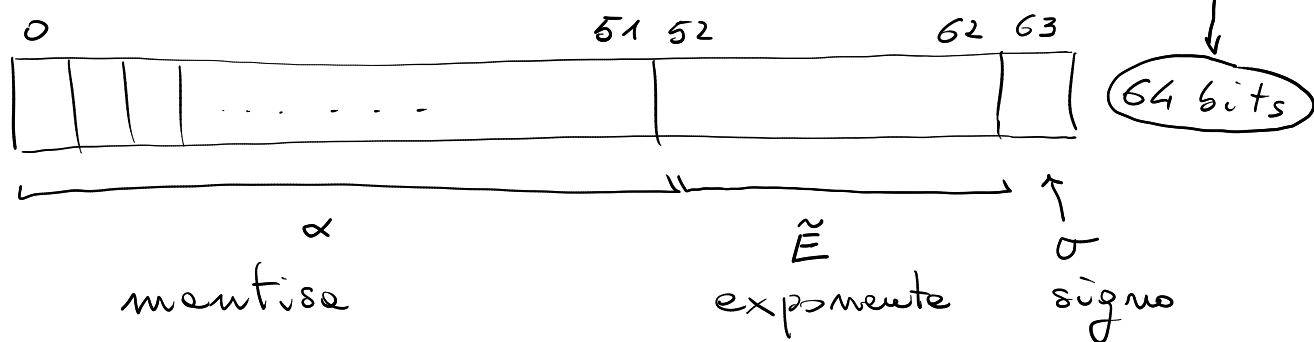
↘ Institute of Electrical & Electronics Engineers

los números reales se representan en la máquina como

$$X = (-1)^{\sigma} (1.\alpha_1\alpha_2\dots\alpha_m)_2 2^E \quad \left[ \begin{array}{l} \text{obs: } (1.011)_2 \cdot 2 = 10.11 \\ (11.001)_2 \cdot 2^{-1} = 1.1001 \end{array} \right]$$

$\sigma \in \{0, 1\}$ ,  $\alpha_j \in \{0, 1\}$ ,  $j = 1 \dots m$ ,  $m = 52$  double precision

$E \in \{-1022, \dots, 1023\}$



$\tilde{E}$  : 11 bits  $\Rightarrow$  se pueden representar todos los enteros entre 0 y  $2^{11}-1 = 2047$

$\tilde{E} = 0$  reservado para representar 0  
(y para tratar los problemas de underflow)

$\tilde{E} = 2047$  reservado para representar  $\infty$

$\hookrightarrow$  quedan 2046 valores disponibles para definir  $E$

$$E = \tilde{E} - 1023$$

- $\text{realmin}$  = número más pequeño que se pueda representar (en valor absoluto)

$$(1.000\dots 0)_2 \cdot 2^{-1022} = 2^{-1022} \approx 2.225 \cdot 10^{-308}$$

- $\text{realmax}$  = número más grande que se pueda representar (en valor absoluto)

$$(1.11\dots 1)_2 \cdot 2^{1023} = (2 - 2^{-52}) \cdot 2^{1023} \approx 1.8 \cdot 10^{308}$$

- $\text{flintmax}$  = número entero más grande hasta el que se pueden representar todos los enteros

es  $2^{53}$ . entre  $2^{53}$  y  $2^{54}$  tenemos 1 entero cada 2  
entre  $2^{54}$  y  $2^{55}$  " " 4

- $\text{eps}$  :  $\epsilon$  - máquina / machine precision  
precisión con la que se pueden representar los números cerca de 1

- el número IEEE 754 64 bits  $x > 1$  más pequeño es

$$x = (1.00\dots 01)_2 \cdot 2^0 = \underbrace{1 + 2^{-52}}$$

$\uparrow$   
 $2^{-52}$

observar que  $1 + 2^{-53} = (1.000\dots 01)_2 \cdot 2^0 = 1$   
52 X ↑ se necesitaría un bit 53

- el número IEEE 754 64 bits  $x < 1$  más grande es

$$x = (1.11\dots 1)_2 \cdot 2^{-1} = 1 - 2^{-53}$$

observación: la aritmética en floating point  
no es asociativa

$$1 - (\underbrace{1 + 2^{-53}}_1) = 0$$

$$(\underbrace{1 - 1}_0) - 2^{-53} = -2^{-53}$$

|                              |   |  |
|------------------------------|---|--|
| $1 + 2^{-53} - 1 = 0$        | } | <u>ORDEN</u> : de<br>izquierda a derecha |
| $-1 + 1 + 2^{-53} = 2^{-53}$ |   |  |