

Representació de nombres en punt flotant

- La notació entera no permet representar la ampla gama de nombres que deriven de qualsevol càlcul científic.

Punt flotant

exponent E i mantissa M

$$N = M \cdot B^E$$

- On B és la **base**, normalment $B=2$.
- L'exponent i la mantissa són enters de longitud determinada
- L'exponent sol tenir entre 5 i 20 bits i la mantissa, de 8 a 100.
- Nombres negatius. Representació signe-mòdul

Signe (S)	Exponent (E)	Mantissa (M)
--------------	-----------------	-----------------

Normalització

Normalització, per evitar que un mateix número tingui diverses representacions

$$0,001 \cdot 2^3 = 0.01 \cdot 2^2$$

Dues regles de Normalització, que es poden usar indistintament:

- (i) El dígit més significatiu diferent de 0 ha d'estar immediatament **després** del punt: **0,10110**
- (ii) El dígit més significatiu diferent de 0 ha d'estar immediatament **abans** del punt: **1,0110**

Un número es normalitza desplaçant els bits de la mantissa tants llocs a l'esquerra (o a la dreta) fins que es compleixi una de les condicions anteriors, i després restant (o sumant) a l'exponent el valor del número de llocs desplaçats.

Per exemple amb la regla (i):

$$0,001101 \cdot 2^5 = 0,1101 \cdot 2^3$$

$$101,110 \cdot 2^4 = 0,101110 \cdot 2^7$$

Per exemple amb la regla (ii):

$$0,001101 \cdot 2^5 = 1,101 \cdot 2^2$$

$$1100,001 \cdot 2^6 = 1,100001 \cdot 2^9$$

Polarització

- **La representació del zero encara provoca problemes.** Quan la mantissa és 0, l'exponent pot tenir qualsevol valor ($0 \cdot 2^E = 0$),
- Un exponent negatiu gran representa els números que s'aproximen gradualment a zero
- També és aconsellable designar el grup de bits format tot per 0's com el "zero", així les representacions del zero en punt fix i flotant són idèntiques. Per això cal que l'exponent 00...00 correspongui a l'exponent més negatiu, $-E_M$.
- L'exponent es calcula restant-li E_M al camp exponent del número en punt flotant. Es diu que l'exponent està **polaritzat** E_M

El valor de l'exponent pot considerar-se com el nombre de desplaçaments que cal fer per normalitzar la mantissa (dins del marge $1/2 \leq |M| < 1$).

Si $|M| > 1$ s'ha de desplaçar a la dreta i sumar a l'exponent el mateix nombre de posicions, per exemple: $N = 101.101$ pot escriure's com $N = 2^3 (.101101)$.

En canvi, si $|M| \leq 1/2$, la mantissa es desplaça cap a l'esquerra, de manera que el bit més significatiu és 1 i ara l'exponent es resta el nombre de posicions desplaçades.

Format estàndard IEEE 754

- 1985 IEEE estableix estàndard - màxima compatibilitat entre sistemes.
- Tres tipus de representació de nombres reals en punt flotant, segons la precisió:

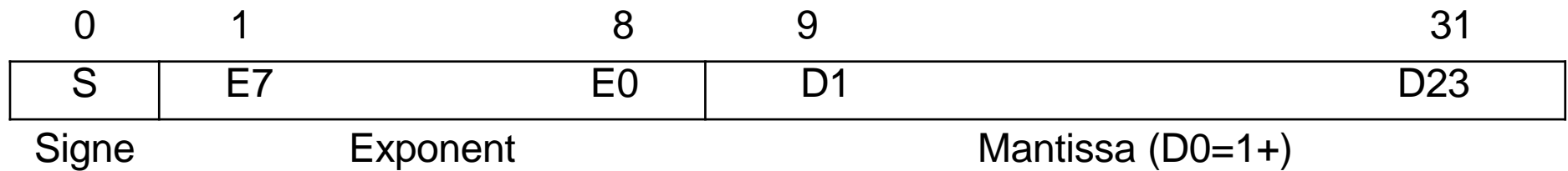
real curt (32 bits)

real llarg (64 bits)

real temporal (80 bits)

El Real-curt 32 bits (4 bytes):

1 bit per al signe, 8 bits per a l'exponent i 23 per a la mantissa (regla ii):



$$N = (-1)^S 2^{\text{EXP} - 127} \left(1 + \sum_{i=1}^{23} \frac{D_i}{2^i} \right)$$

Per exemple:

0	00000000	000000000000000000000000	= 0
1	01111111	000000000000000000000000	= -1
0	11111111	000000000000000000000000	= $+\infty$
1	10000010	011101100000000000000000	= -11,625
0	11111111	01010111100110101000011	= NaN

Bit fantasma

Aritmètica amb punt flotant

SUMA I RESTA

- És directa si ambdós números tenen el mateix exponent
- Si l'exponent és diferent cal **ALINEAMENT**:
- Desplaçar mantissa del núm. de menor exp. El desplaçament ve determinat per la resta dels dos exps.

X = 0 1000101 101000000000000000000000 (2⁶ · 1,101) (bit fantasma implícit és 1)
Y = 0 10001010 101101000000000000000000 (2¹¹ · 1,101101)

Cal desplaçar la mantissa d'X 5 llocs a la dreta i posar l'exponent a 11.

Aleshores:

$$\begin{array}{r} 0,00001101 \cdot 2^{11} \\ + 1,101101 \cdot 2^{11} \\ \hline 1,11000001 \cdot 2^{11} \end{array}$$

MULTIPLICACIÓ I DIVISIÓ

Multiplicació (divisió) es multipliquen (divideixen) les mantisses i se sumen (resten) els exponents.

Per exemple,
si $X = 2^{E_x} \cdot M_x$; $Y = 2^{E_y} \cdot M_y$, tenim que:

$$X \cdot Y = (2^{E_x + E_y}) (M_x \cdot M_y) = 2^{E_v} \cdot M^v$$

$$X/Y = (2^{E_x - E_y})(M_x/M_y) = 2^{E_w} \cdot M_w$$

El resultat ha de normalitzar-se per tal de complir les regles inicials.