#### Grupo ARCOS

### uc3m Universidad Carlos III de Madrid

#### Tema 2 Coma flotante

Estructura de Computadores Grado en Ingeniería Informática



#### Contenidos

#### I. Introducción

- Motivación y objetivos
- 2. Sistemas posicionales

#### 2. Representaciones

- Alfanuméricas
  - Caracteres
  - 2. Cadenas de caracteres
- 2. Numéricas
  - Naturales y enteras
  - 2. Coma fija
  - 3. Coma flotante (estándar IEEE 754)

### Otras necesidades de representación

#### ¿Cómo representar?

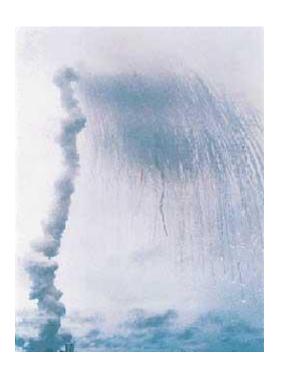
- Números muy grandes: 30.556.926.000<sub>(10</sub>
- Números muy pequeños: 0.000000000529177<sub>(10</sub>
- Números con decimales: 1,58567

### Ejemplo de fallo...

- Explosión del Ariane 5 (primer viaje)
  - Enviado por ESA en junio de 1996
  - Coste del desarrollo:
     10 años y 7000 millones de dólares
  - Explotó 40 segundos después de despegar,
     a 3700 metros de altura.



El software del sistema de referencia inercial realizó la conversión de un valor real en coma flotante de 64 bits a un valor entero de 16 bits. El número a almacenar era mayor de 32767 (el mayor entero con signo de 16 bits) y se produjo un fallo de conversión y una excepción.



### Coma fija [racionales]

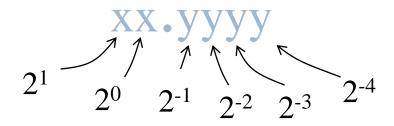
 Se fija la posición de la coma binaria y se utilizan los pesos asociados a las posiciones decimales

Ejemplo:

$$|00|.|0|0 = 2^4 + 2^0 + 2^{-1} + 2^{-3} = 9,625$$

# Representación de fracciones con representación binaria en coma fija

Ejemplo de representación con 6 bits:



- Ejemplo de número:  $10,1010_{(2} = 1 \times 2^{1} + 1 \times 2^{-1} + 1 \times 2^{-3} = 2.62510$
- Asumiendo esta coma fija, el rango sería:
  - □ [0 a 3.9375 (casi 4)]

## Potencias negativas

i	2-i	
0	1.0	1
1	0.5	1/2
2	0.251/4	
3	0.125	1/8
4	0.0625	1/16
5	0.03125	1/32
6	0.015625	
7	0.0078125	
8	0.00390625	
9	0.001953125	
10	0.0009765625	

#### Contenidos

#### I. Introducción

- Motivación y objetivos
- 2. Sistemas posicionales

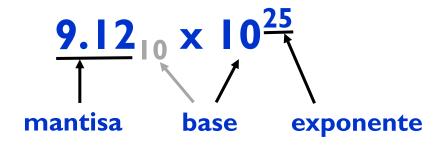
#### 2. Representaciones

- Alfanuméricas
  - Caracteres
  - 2. Cadenas de caracteres

#### 2. Numéricas

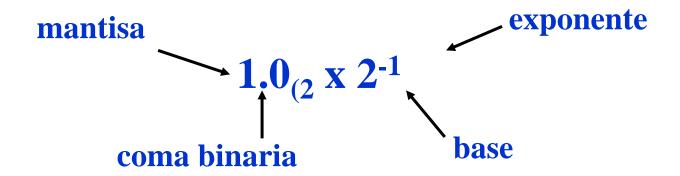
- Naturales y enteras
- 2. Coma fija
- 3. Coma flotante (estándar IEEE 754)

#### Notación científica decimal



- Cada número lleva asociado una mantisa y un exponente
- Notación científica decimal usada: notación normalizada
  - Solo un dígito distinto de 0 a la izquierda del punto
- Se adapta el número al orden de magnitud del valor a representar, trasladando la coma decimal mediante el exponente

#### Notación científica en binario



- Forma normalizada: Un I (solo un dígito) a la izq. de la coma
  - Normalizada:  $1.0001 \times 2^{-9}$
  - No normalizada:  $0.0011 \times 2^{-8}$ ,  $10.0 \times 2^{-10}$

# Estándar IEEE 754 [racionales]



- Estándar para coma flotante usado en la mayoría de los ordenadores.
- Características (salvo casos especiales):
  - Exponente: en exceso con sesgo 2 num\_bits\_exponente I I
  - Mantisa: signo-magnitud, normalizada, con bit implícito
- Diferentes formatos:
  - Precisión simple: 32 bits (signo: I, exponente: 8 y mantisa: 23)
  - **Doble precisión**: 64 bits (signo: I, exponente: I I y mantisa: 52)
  - Cuádruple precisión: 128 bits (signo: 1, exponente: 15 y mantisa: 112)

### Normalización y bit implícito

#### Normalización

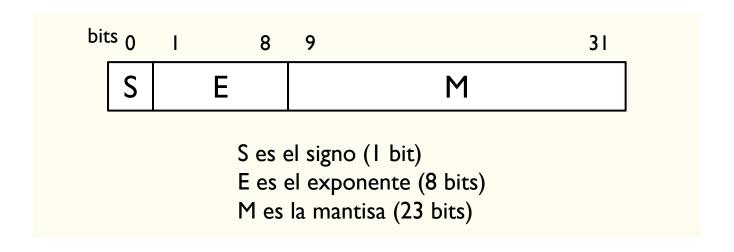
Para normalizar la mantisa se ajusta el exponente para que el bit más significativo de la mantisa sea I

- Ejemplo:  $00010000000010101 \times 2^3$  (no lo está)  $100000000010101000 \times 2^0$  (ahora sí)

#### Bit implícito

Una vez normalizado, dado que el bit más significativo es 1, **no** se almacena para dejar espacio para un bit más (aumenta la precisión)

Así se puede representar mantisas con un bit más



▶ El valor se calcula con la siguiente expresión (salvo casos especiales):

$$N = (-1)^{S} \times 2^{E-127} \times 1.M$$

#### donde:

S = 0 indica número positivo, S = I indica número negativo

0 < E < 255 (E=0 y E=255 indican casos especiales)

Existencia de casos especiales:

(-1)s \* 0.mantisa \* 2-126

Exponente	Mantisa	Valor especial
0 (0000 0000)	0	+/- 0 (según signo)
0 (0000 0000)	No cero	Número NO normalizado
255 (1111 1111)	No cero	NaN (0/0,)
255 (1111 1111)	0	+/-infinito (según signo)
I-254	Cualquiera	Número normalizado (no
		especial)

(-I)s \* I.mantisa \* 2exponente-127

## Ejemplos (incluyen casos especiales)

S	E	M	N
I	00000000	000000000000000000000000000000000000000	-0 (Excepción 0) E=0 y M=0.
I	01111111	000000000000000000000000000000000000000	$-2^{0} \times 1.0_{2} = -1$
0	10000001	111000000000000000000000000000000000000	$+2^2 \times 1.111_2 = +2^2 \times (2^0 + 2^{-1} + 2^{-2} + 2^{-3}) = +7.5$
0	111111111	000000000000000000000000000000000000000	∞ (Excepción ∞) E=255 y M=0
0	111111111	100000000000000000000000000000000000000	NaN (Not a Number) E=255 y M≠0.

### Ejercicio

a) Calcular el valor correspondiente al número
 0 10000011 1100000000000000000
 dado en coma flotante según norma 754 de simple precisión

### Ejercicio (solución)

- a) Calcular el valor correspondiente al número
   0 10000011 11000000000000000000
   dado en coma flotante según norma 754 de simple precisión
  - a) Bit de signo:  $0 \Rightarrow (-1)^0 = +1$
  - Exponente:  $10000011_2 = 131_{10} \Rightarrow E 127 = 131 127 = 4$

Por tanto el valor decimal del n° es  $+1 \times 2^4 \times 1,75 = +28$ 

## Ejercicio

b) Expresar según norma IEEE 754 de simple precisión el n°-9

### Ejercicio (solución)

b) Expresar según norma IEEE 754 de simple precisión el n°-9

$$-9_{10} = -1001_2 = -1001_2 \times 2^0 = -1,001_2 \times 2^3$$
 (mantisa normalizada)

- a) Bit de signo: negativo  $\implies$  S=I
- Exponente: 3+127 (exceso) =  $130 \implies 10000010$

- Rango de magnitudes representables (sin considerar el signo):
  - Menor normalizado:
  - Mayor normalizado:

- Menor no normalizado:
- Mayor no normalizado:

 $(-1)^s * 0.mantisa * 2^{-126}$ 

Exponente	Mantisa	Valor especial
0	<b>≠ 0</b>	No normalizado
1-254	cualquiera	normalizado

(-I)<sup>s</sup> \* I.mantisa \* 2<sup>exponente-127</sup>

- Rango de magnitudes representables (sin considerar el signo):
  - Menor normalizado:

Mayor normalizado:

- Menor no normalizado:
- Mayor no normalizado:

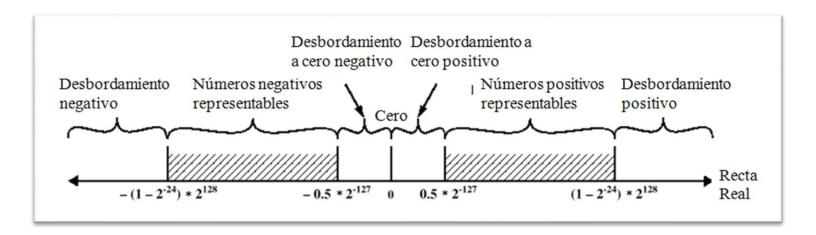
#### Truco:

$$X = 2 - 2^{-23}$$

- Rango de magnitudes representables (sin considerar el signo):
  - Menor normalizado:

Mayor normalizado:

- Menor no normalizado:
- Mayor no normalizado:



## Ejercicio

¿Cuántos números de floats (coma flotante de simple precisión) hay entre el 1 y el 2 (no incluido)?

¿Cuántos números de floats (coma flotante de simple precisión) hay entre el 2 y el 3 (no incluido)?

### Ejercicio (solución)

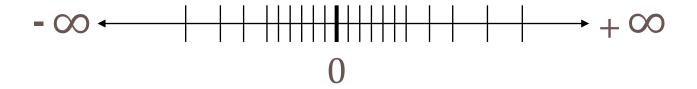
- ¿Cuántos números de floats (coma flotante de simple precisión) hay entre el 1 y el 2 (no incluido)?

  - ▶ Entre I y 2 hay 2<sup>23</sup> números
- ¿Cuántos números de floats (coma flotante de simple precisión) hay entre el 2 y el 3 (no incluido)?

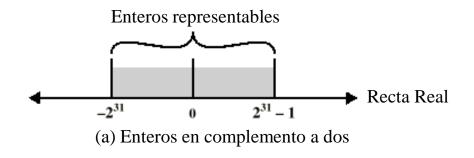
  - ▶ Entre 2 y 3 hay 2<sup>22</sup> números

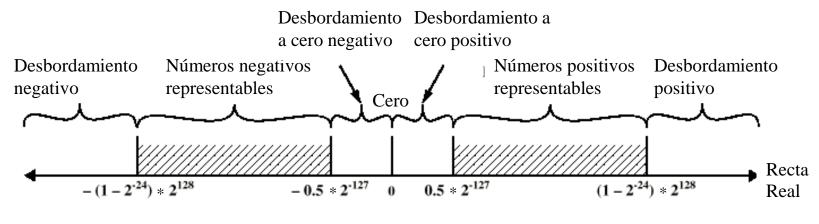
### Números representables

Resolución variable:
 Más denso cerca de cero, menos hacia el infinito



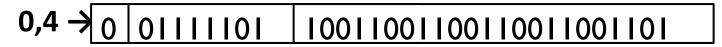
#### Números representables





(b) Números en coma flotante

# Ejemplo 1 imprecisión





3.9999998 x 10<sup>-1</sup>





9.9999994 x 10<sup>-2</sup>

# Ejemplo 2 imprecisión

¿Cómo realiza C una división?

```
t2.c
#include <stdio.h>
int main ()
 float a;
 a = 3.0/7.0;
 if (a == 3.0/7.0)
      printf("Igual\n");
 else printf("No Igual\n");
 return (0);
```

# Ejemplo 2 imprecisión

¿Cómo realiza C una división?

```
t2.c
#include <stdio.h>
int main ()
 float a;
  a = 3.0/7.0;
  if (a == 3.0/7.0)
      printf("Igual\n");
  else printf("No Igual\n");
  return (0);
```

```
$ gcc -o t2 t2.c
$ ./t2
No Igual
```

# Ejemplo 2 imprecisión

¿Cómo realiza C una división?

```
t2.c
         #include <stdio.h>
         int main ()
           float a;
                             double
float
           a = 3.0/7.0;
           if (a == 3.0/7.0)
               printf("Igual\n");
           else printf("No Igual\n");
           return (0);
```

```
$ gcc -o t2 t2.c
$ ./t2
No Igual
```

## Ejemplo 3 imprecisión

La propiedad asociativa no siempre se cumple ¿ a + (b + c) = (a + b) + c ?

```
#include <stdio.h>

int main ()
{
    float x, y, z;

    x = 10e30; y = -10e30; z = 1;
    printf("(x+y)+z = %f\n",(x+y)+z);
    printf("x+(y+z) = %f\n",x+(y+z));

    return (0);
}
```

# Ejemplo 3 imprecisión

La propiedad asociativa no siempre se cumple a + (b + c) = (a + b) + c?

```
#include <stdio.h>

int main ()
{
    float x, y, z;

    x = 10e30; y = -10e30; z = 1;
    printf("(x+y)+z = %f\n",(x+y)+z);
    printf("x+(y+z) = %f\n",x+(y+z));

    return (0);
}
```

```
$ gcc -o t1 t1.c
$ ./t1
(x+y)+z = 1.000000
x+(y+z) = 0.000000
```

#### **Asociatividad**

La coma flotante no es asociativa

$$x = -1.5 \times 10^{38}, y = 1.5 \times 10^{38}, y z = 1.0$$

$$(x + y) + z = (-1.5 \times 10^{38} + 1.5 \times 10^{38}) + 1.0$$

$$= (0.0) + 1.0 = 1.0$$

- Las operaciones coma flotante no son asociativas
  - Los resultados son aproximados
  - $ightharpoonup 1.5 imes 10^{38}$  es mucho más grande que 1.0
  - $1.5 \times 10^{38}$  + 1.0 en la representación en coma flotante sigue siendo  $1.5 \times 10^{38}$

#### Conversión int $\rightarrow$ float $\rightarrow$ int

```
if (i == (int)((float) i)) {
    printf("true");
}
```

- ▶ No siempre es cierto
- Muchos valores enteros grandes no tienen una representación exacta en coma flotante
- ¿Qué ocurre con double?

- ▶ El número 133000405 en binario es:
- Se normaliza

  - S = 0 (positivo)
  - $\rightarrow$  e = 26  $\rightarrow$  E = 26 + 127 = 153
- El número realmente almacenado es
  - $\downarrow$  1, 11111011010110110011010  $\times$  2<sup>26</sup> =

#### Conversión float $\rightarrow$ int $\rightarrow$ float

```
if (f == (float)((int) f)) {
    printf("true");
}
```

- No siempre es cierto
- Los números con decimales no tienen representación entera

#### Redondeo

- El redondeo elimina cifras menos significativas de un número para obtener un valor aproximado.
- ▶ Tipos de redondeo:
  - ▶ Redondeo hacia + ∞
    - ▶ Redondeo "hacia arriba":  $2.001 \rightarrow 3$ ,  $-2.001 \rightarrow -2$
  - ▶ Redondeo hacia ∞
    - ▶ Redondea "hacia abajo":  $1.999 \rightarrow 1$ ,  $-1.999 \rightarrow -2$
  - ▶ Truncar
    - $\blacktriangleright$  Descarta los últimos bits: 1.299  $\rightarrow$  1.2
  - Redondeo al más cercano
    - ightharpoonup 2.4 ightharpoonup 2.6 ightharpoonup 3, -1.4 ightharpoonup -1

#### Redondeo

- El redondeo supone ir perdiendo precisión.
- ▶ El redondeo ocurre:
  - > Al pasar a una representación con menos representables:
    - Ej.: Un valor de doble a simple precisión
    - ▶ Ej.: Un valor en coma flotante a entero
  - Al realizar operaciones aritméticas:
    - Ej.: Después de sumar dos números en coma flotante (al usar dígitos de guarda)

# Dígitos de guarda

- Se utilizan dígitos de guarda para mejorar la precisión: internamente se usan dígitos adicionales para operar.
- $\blacktriangleright$  Ejemplo: 2,65 x 10<sup>0</sup> + 2.34 x 10<sup>2</sup>

	SIN dígitos de guarda	CON dígitos de guarda
I igualar exponentes	$0.02 \times 10^2 + 2.34 \times 10^2$	$0.0265 \times 10^{2}$ + $2.3400 \times 10^{2}$
2 sumar	$2,36 \times 10^{2}$	$2,3665 \times 10^2$
3 redondear	$2,36 \times 10^{2}$	$2,37 \times 10^2$

## Operaciones en coma flotante

#### Sumar

#### Restar

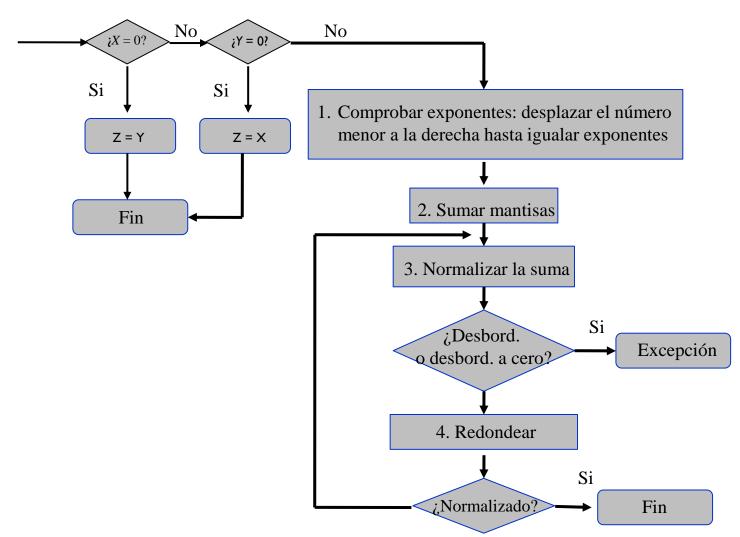
- 1. Comprobar valores cero.
- 2. Igualar exponentes (desplazar número menor a la derecha).
- 3. Sumar/restar las mantisas.
- 4. Normalizar el resultado.

#### Multiplicar

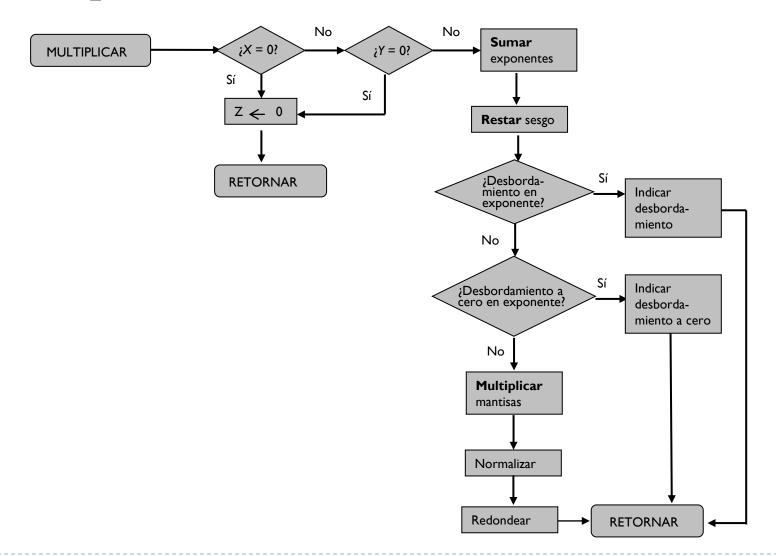
#### Dividir

- 1. Comprobar valores cero.
- 2. Sumar/restar exponentes.
- 3. Multiplicar/dividir mantisas (teniendo en cuenta el signo).
- 4. Normalizar el resultado.
- 5. Redondear el resultado.

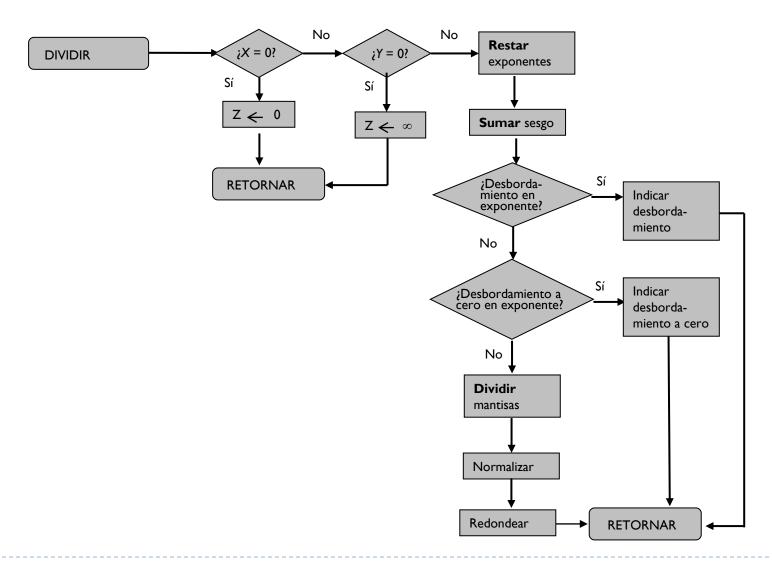
## Suma y resta: Z=X+Y y Z=X-Y



# Multiplicación: Z=X\*Y



### División: Z=X/Y



# Ejercicio

Usando el formato IEEE 754, sumar 7,5 y 1,5 paso a paso

Pasar a binario

- 1) 7,5 + 1,5 =
- 2)  $1,111*2^{2} + 1,1*2^{0} =$
- 3)  $1,111*2^2 + 0,011*2^2 =$
- 4)  $10,010*2^2 =$

5)  $1,0010*2^3$ 

exponentes

Igualar

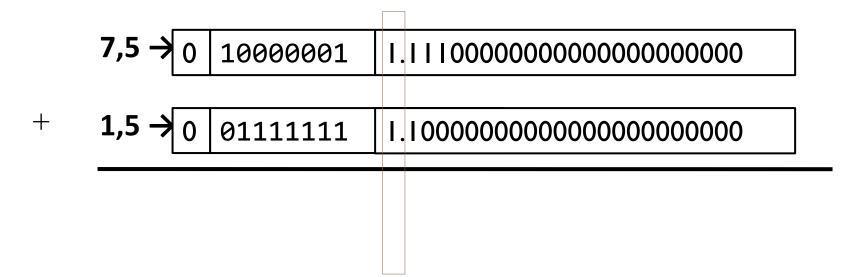
Sumar

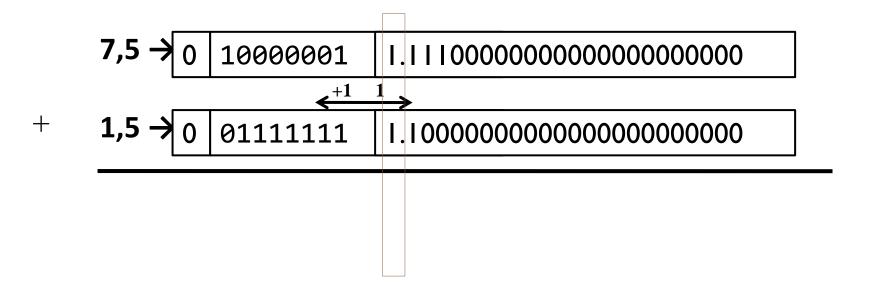
Ajustar exponentes

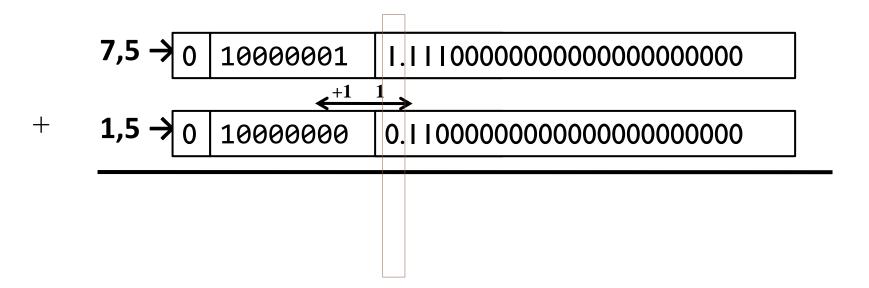
Representación de los números

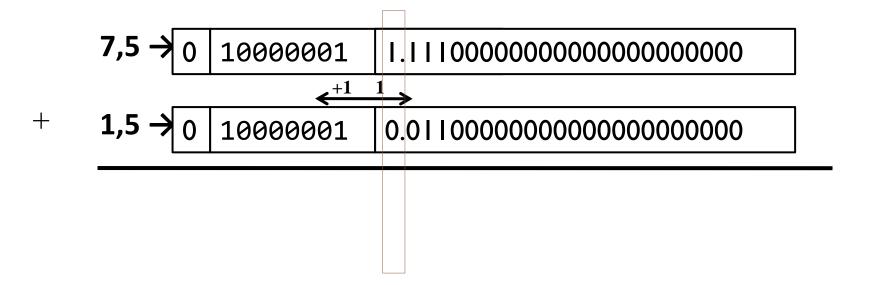
Se separa exponentes y mantisas y se añade el bit implícito

$7,5 \rightarrow 0$	10000001	1.	111000000000000000000000000000000000000
<b>1,5</b> → 0	01111111	1.	100000000000000000000000000000000000000

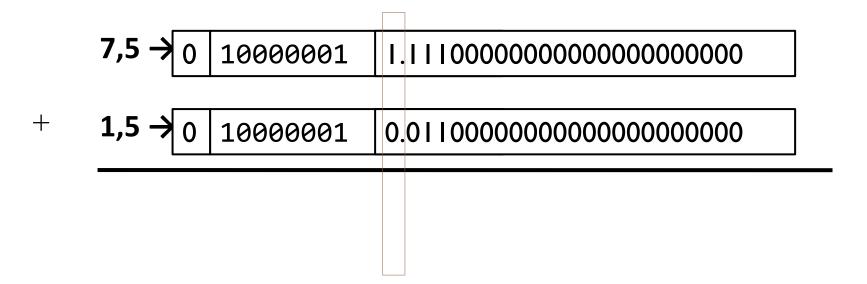




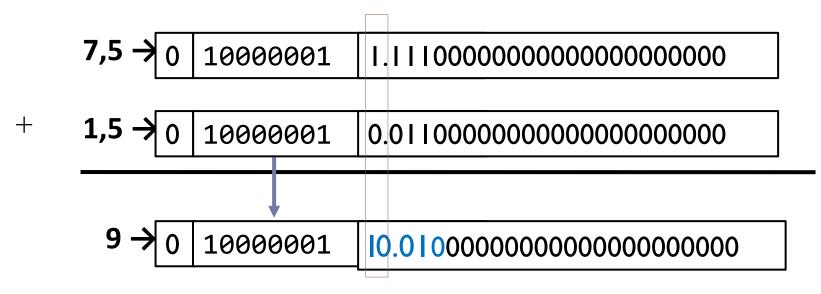




#### Sumar mantisas

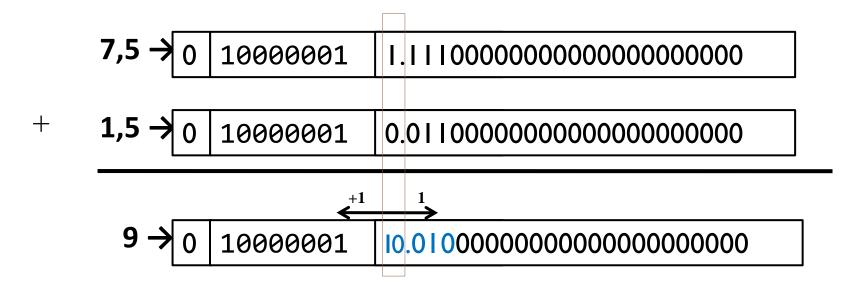


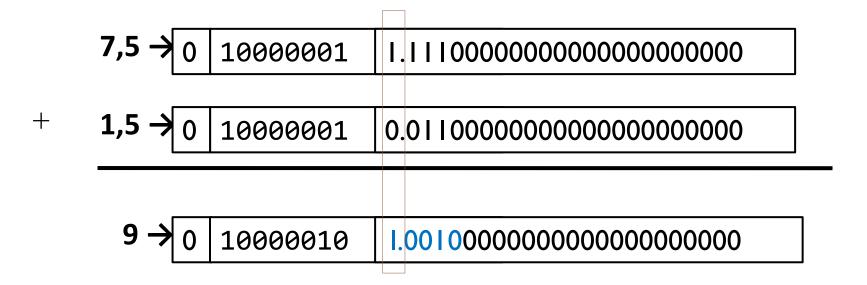
#### Normalizar el resultado



Se produce un acarreo, mantisa no normalizada

#### Normalizar el resultado





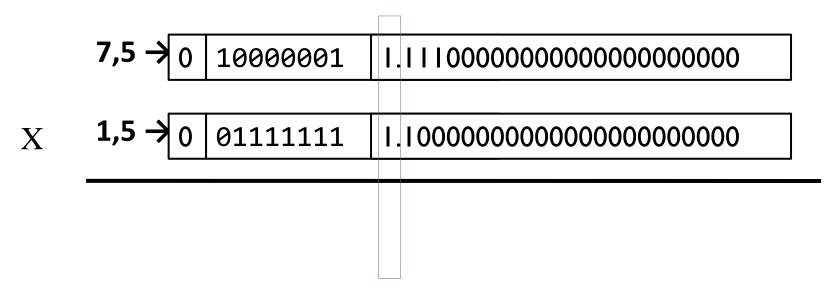
Se almacena el resultado eliminando el bit implícito

# Ejercicio

Usando el formato IEEE 754, multiplicar 7,5 y 1,5 paso a paso

Representación de los números

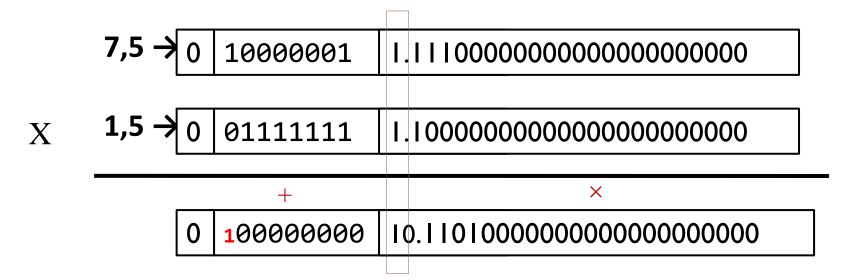
Se separan exponentes y mantisas y se añade bit implícito



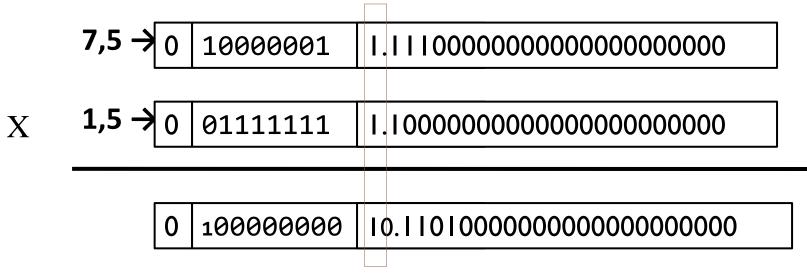
Se añade el bit implícito para operar

60

Multiplicar: sumar exponentes y multiplicar mantisas

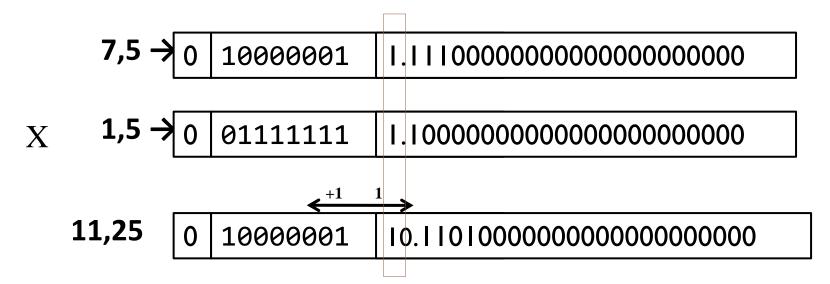


Multiplicar: quitar el sesgo al exponente (hay dos)



- 01111111

Multiplicar: normalizar el resultado



▶ Resultado normalizado...

	7,5 <del>→</del>	0	10000001		.11100000000000000000000000000000000000
X	1,5 →	0	01111111	1	.10000000000000000000000000000000000000
	11,25	0	10000010		.01101000000000000000000000000000000000

Se almacena el resultado eliminando el bit implícito

#### Evolución de IEEE 754

- ▶ 1985 IEEE 754
- ▶ 2008 IEEE 754-2008 (754+854)
- ▶ 2011 ISO/IEC/IEEE 60559:2011 (754-2008)

Name	Common name	Base	Digits	E min	E max	Notes	Decimal digits	Decimal E max
binary16	Half precision	2	10+1	-14	+15	storage, not basic	3.31	4.51
binary32	Single precision	2	23+1	-126	+127		7.22	38.23
binary64	Double precision	2	52+I	-1022	+1023		15.95	307.95
binary128	Quadruple precision	2	112+1	-16382	+16383		34.02	4931.77
decimal32		10	7	-95	+96	storage, not basic	7	96
decimal64		10	16	-383	+384		16	384
decimal 128		10	34	-6143	+6144		34	6144

http://en.wikipedia.org/wiki/IEEE\_floating\_point

#### Grupo ARCOS

# uc3m Universidad Carlos III de Madrid

## Tema 2 Coma flotante

Estructura de Computadores Grado en Ingeniería Informática

