

## Rappresentazione dei numeri razionali

$$\text{Valore di } N = \sum_{i=-m}^{n-1} c_i b^i$$

$n \equiv$  numero cifre parte intera

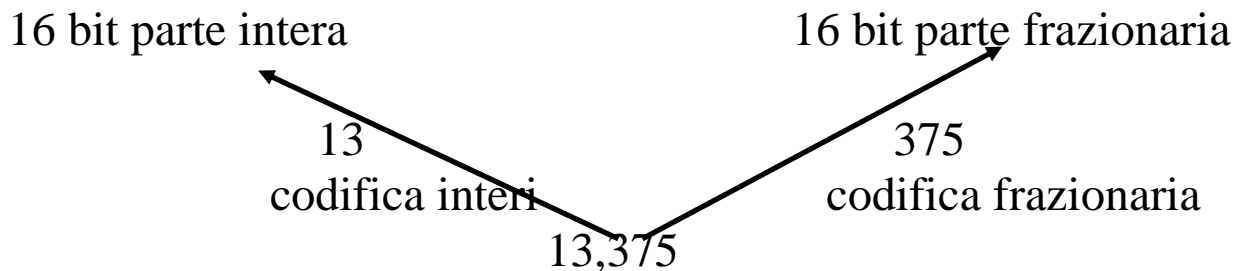
$m \equiv$  numero cifre parte decimale

Es.

$$13,375_{10} = 5 \cdot 10^{-3} + 7 \cdot 10^{-2} + 3 \cdot 10^{-1} + 3 \cdot 10^0 + 1 \cdot 10^1$$

Rappresentazione in virgola fissa

- Precisione costante della parte frazionaria
- $E_a$  costante



codifica parte frazionaria:

moltiplicazione per 2 sino a quando vale 0 e i bit corrispondono ai riporti nell'ordine prodotto:

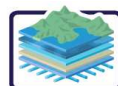
Es.:

	x 2	riporto
0,375	0,75	0 $2^{-1}$
0,75	1,5	1 $2^{-2}$
0,5	1,0	1 $2^{-3}$
0		

$$0,375_{10} = 011\ 0000000\ldots$$



POLITECNICO  
DI MILANO



SpatialDBgroup

Politecnico di Milano – DEI –  
Prof. Mauro Negri

- L'algoritmo può non convergere

0,1	0,2	0	
0,2	0,4	0	0,1 <sub>10</sub> =0.0001100011... (∞ bit)
0,4	0,8	0	
0,8	1,6	1	
0,6	1,2	1	
0,2	0,4 ....		si ripete in modo periodico

e se non converge si introduce un'approssimazione (imprecisione)

- Errore assoluto costante  $2^{-16} = 0,000015$
- Inoltre la codifica dimostra che non esiste relazione tra il numero di cifre della parte frazionaria nel numero decimale e in quello corrispondente binario
- Rigidità della pre-divisione dei bit

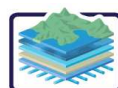
8.750.000.000      non rappresentabile nei 16 bit della parte intera  
16 bit frazionaria inutilizzati

0,00000000000875    16 bit parte intera inutilizzati  
16 bit parte frazionaria a 0 (approssimazione)



POLITECNICO  
DI MILANO

Politecnico di Milano – DEI –  
Prof. Mauro Negri



SpatialDBgroup

## Real numbers in finite representation “a large grey area”

Before 1985

Non esiste un accordo su un format per real numbers (fighting)

1985 Agreement

Standard IEEE 754-1985 for binary FP arithmetic

“What Every Computer Scientist Should Know About Floating-Point Arithmetic”, David Goldberg, ACM Computing Surveys, Vol 23, No 1, March 1991, pp. 5-48

- non evita tutti i problemi
- stabilisce vincoli sull'entità dei “rounding error” per le operazioni aritmetiche
- un'implementazione hw è IEEE compliant se produce risultati uguali a quelli degli algoritmi IEEE



Program1 (CPU1) = Program1 (CPU2)

Tutto risolto?

New York Times, nov. 1994 “Intel’s Pentium problem persists”  
(300 milioni di dollari per il ritiro)

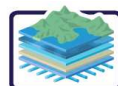
2008 IEEE-754-2008 standard per decimal FP arithmetic

Altri problemi più avanti



POLITECNICO  
DI MILANO

Politecnico di Milano – DEI –  
Prof. Mauro Negri



**SpatialDBgroup**

## Virgola mobile (standard ANSI/IEEE 754-1985 binary FP arithmetics)

Obiettivi della rappresentazione:

- autoadattabilità tra parte intera e frazionaria
- errore relativo costante

Rappresentazione normalizzata  $N = s M \cdot 2^e$   
 $s = \pm$        $1 \leq M < 2$        $b=2$        $e = \pm \text{numero intero}$

Codifiche dedicate per 0, NaN ( $\sqrt{-2}$ ,  $0/0$ ,  $\infty+/-/\infty$ )  $\infty$

Es.  $11_{10} = 1,375 \cdot 2^3$        $0.25_{10} = 1 \cdot 2^{-2}$



Modelli di precisione

	hidden	$n_s$	$n_M$	$n_e$	$e + (2^{(n_e-1)} - 1)$	precisione
singola precis. FP32	0	1	23	8 -126 – 127	$e+127$	24
doppia precis. FP64	0	1	52t	11 -1022 – 1023	$e+1023$	53
Intel 80 bit	1	1	63	15		64

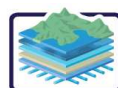
s    0=positivo 1=negativo

M    codifica della sola parte frazionaria

e    codifica del valore positivo es. FP32  $\Rightarrow e+127$



POLITECNICO  
DI MILANO



**SpatialDBgroup**

Politecnico di Milano – DEI –  
Prof. Mauro Negri

(1,) sottinteso in FP

Es.:  $11_{10} = 1,375 \cdot 2^3$

$s=0$

$M = 011000000000000000000000$

$e(3) + 127 = 130 \Rightarrow 10000010$



0 10000010 011000000000000000000000

s

e

M

Esempi di intervalli

Tipi C	Bit	Intervallo possibile (float.h)
<i>float</i>	FP32	$-3,4 \cdot 10^{38} \dots -1,1 \cdot 10^{-38}$ (FLT_MIN) $1,1 \cdot 10^{-38} \dots 3,4 \cdot 10^{38}$ (FLT_MAX) FLT_MIN = $(1+0,5+0,25+\dots) 2^{127} = 2^{128}$ FLT_MIN = $(1+0) 2^{-126}$
<i>double</i>	FP64	$-1,7 \cdot 10^{308} \dots -2,2 \cdot 10^{-308}$ (DBL_MIN) $2,2 \cdot 10^{-308} a 1,7 \cdot 10^{308}$ (DBL_MAX)

**Osservazione 1** FP vs v.fissa

8.750.000.000

v. fissa: non rappresentabile

FP32:  $0,875 \cdot 10^{10} = 1,018634065 \dots \cdot 2^{33}$

$= 0\ 10100000\ 00000100110001010011001$

0,00000000000875

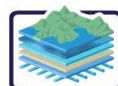
v.fissa: rappresentato come 0

FP32:  $0,875 \cdot 10^{-10} = 1,5032385 \dots \cdot 2^{-34}$

$= 0\ 01011101\ 10000000110101000011110$



POLITECNICO  
DI MILANO



SpatialDBgroup

Politecnico di Milano – DEI –  
Prof. Mauro Negri

## Osservazione 2 Come si esegue una somma?

$$S \cdot 2^E + T \cdot 2^F \text{ con } E > F$$



- denormalization of T: shifting right T (+hidden bit) di (E-F) bits
- sum, normalize and rounding

Es. Mantissa da 7 bit

$$\begin{array}{rcl} 3 + & 1.5 \cdot 2^1 & 1000000 \\ 0.75 = & 1.5 \cdot 2^{-1} & 1]1000000 \text{ shift 2bit sx } 0110000 \\ & & \hline & 1.875 \cdot 2 & \leftarrow 1110000 \end{array}$$

## Osservazione 3 NON corrispondenza precisione decimale e FP

numero

1)  $2,1 = 1,05 \cdot 2^1 \Rightarrow 00001100110011\dots\dots$

2)  $1,5 = 1,5 \cdot 2^0 \Rightarrow 1(0..0) = 1,5$

3)  $1,8750 = 1,875 \cdot 2^0 \Rightarrow 111(0\dots0) = 1,875$

4)  $16,5625 = 1,035 \cdot 2^4 \Rightarrow 00001001$

Corrispondenza precisa

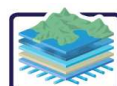
## Osservazione 4. Corrispondenza biunivoca tra decimale e FP

Si afferma che esiste in FP32/64 con decimali composti da 6/15 cifre decimali.



POLITECNICO  
DI MILANO

Politecnico di Milano – DEI –  
Prof. Mauro Negri



SpatialDBgroup

## Come si deriva

Detta

- p la precisione binaria del sistema
- q una precisione decimale

Esiste la corrispondenza se

$$\begin{array}{ccc} \text{configurazioni binarie} & & \text{configurazioni decimali} \\ 2^{(p-1)} & \geq & 10^q \end{array}$$

Da cui  $q = \lfloor (p-1) \log_{10} 2 \rfloor$  che per  $p-1=23/52$  ottiene 6/15

## Cosa significa

Numeri decimali composti da 6/15 cifre complessive sono in corrispondenza biunivoca con configurazioni FP32/64  
NON che i decimali con 6/15 cifre della parte frazionaria .....

## Per capire meglio

$\epsilon$  (machine epsilon – FLT\_EPSILON): distanza tra il numero 1 e quello immediatamente successivo in FP

$$1.0000000...0001 - 1.0000000...0000$$

$$\begin{array}{ccc} & \text{FP32} & \text{FP64} \\ \epsilon = 2^{-(p-1)} & \epsilon = 2^{-23} & 2^{-52} \end{array}$$

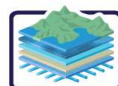
Distanza tra un valore x e il successivo in FP (**u**nits in **l**ast **p**lace – ulp(x)) ?

$$\text{ulp}(x) = \epsilon * 2^e = 2^{-(p-1) + e} \quad \text{deriva da } (1+2^{-(p-1)}) * 2^e - 1 * 2^e$$



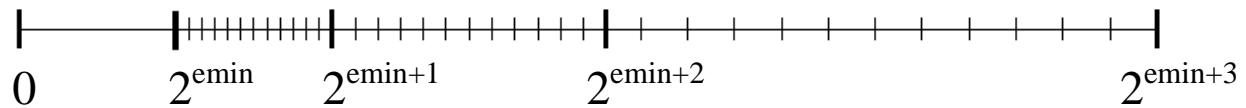
POLITECNICO  
DI MILANO

Politecnico di Milano – DEI –  
Prof. Mauro Negri



SpatialDBgroup

- cresce al crescere di e
- è costante nell'area di valori che hanno lo stesso valore di e



Ciò significa che l'errore assoluto commesso da un'approssimazione è costante in ognuno degli intervalli precedenti e cresce al crescere del valore.

Per avere un'idea concreta

Distanza minima in prossimità di 0:

$$\text{FP32} \quad 2^{-23-126} = 10^{-45}$$

Distanza nell'intervallo  $[1, 2[$

$$\text{FP32} \quad 2^{-23+0} = 10^{-7}$$

$$\text{FP64} \quad 2^{-52+0} = (10^{-16})$$

Distanza 1 dove  $2^{-(p-1)+e} = 1$  ossia  $-(p-1) + e = 0$

$$\text{FP32 } e=23 \text{ ossia in } [2^{23}, 2^{24}[ \text{ ossia } [8.388608, 16.777216[$$

$$\text{FP64 } e=52 \text{ ossia in } [2^{52}, 2^{53}[ \text{ ossia } [4.5 \cdot 10^{15}, 9.07 \cdot 10^{15}[$$

Distanza massima

$$\text{FP32} \quad 2^{-23+127} = 10^{31}$$

$$\text{FP64} \quad 2^{-52+1023} = 10^{292}$$

Conclusione: il numero di cifre decimali varia in base all'intervallo considerato

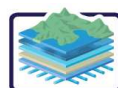


$$\text{int 32} \quad \approx 4 \cdot 10^9$$

$$\text{FP 32} \quad \approx 6.8 \cdot 10^{38}$$



POLITECNICO  
DI MILANO



SpatialDBgroup

Politecnico di Milano – DEI –  
Prof. Mauro Negri



## Osservazione 5 La flessibilità nella gestione dei bit

Nell'esempio supponiamo di usare  $N=7$  bit per la mantissa

origine	norm.	p.intera	p.fraz	mantissa	stored
2,1	$1,05 * 2^1$	10	00011001100.. $\infty$	0000110	2.093
1,5	$1,5 * 2^0$	1	100000000000	1000000	1.5
1,8750	$1,875 * 2^0$	1	111000000000	1110000	1.8750
16,5625	$1,035 * 2^4$	10000	100100000000	0000100	16.5

- si considerano gli “e” bit della parte intera da destra
- hidden bit
- si aggiungono N-e bit della parte frazionaria
- p.frazionaria decresce al crescere della p.intera
- l'arrotondamento è costante all'interno dell'intervallo  $[1,00... * 2^x \text{ e } 1,00 * 2^{(x+1)}]$  perché hanno lo stesso numero di bit della parte frazionaria.

In altri termini (consideriamo FP 32 per semplicità)

Intervallo  $[1..2[$                       ossia numeri  $1,xx * 2^0$

23 bit mantissa per la parte frazionaria    ossia 8.000.000 combinazioni  
ossia 6 cifre decimali

Intervallo  $[2..4[$                       ossia numeri  $1,xx * 2^1$

1 bit mantissa per parte intera

22 bit mantissa per la parte frazionaria    ossia 4.000.000 combinazioni  
ossia 2.000.000 a unità ossia ancora 6 cifre decimali

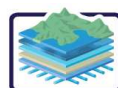
Intervallo  $[\approx 8\text{milioni} - 16777215]$                       ossia  $1,xx * 2^{23}$

23 bit mantissa per parte intera

0 bit mantissa per la parte frazionaria rappresentazione interi



POLITECNICO  
DI MILANO



SpatialDBgroup

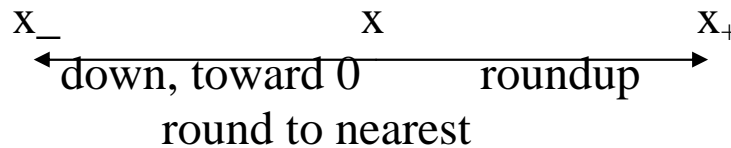
Politecnico di Milano – DEI –  
Prof. Mauro Negri

**Osservazione 6.** A proposito dell'arrotondamento e dell'errore  
Se la sequenza dei bit necessari è maggiore dei bit disponibili  
l'algoritmo approssima attraverso la funzione "round"

Errore assoluto  $Ea(x)$ : spostamento che  $x$  subisce

Errore relativo  $Er(x)$ :  $Ea(x)$  rapportato al valore di  $x$

Dato  $x > 0$



Consideriamo round to nearest

1.  $Ea(x) \quad 0 \leq |x - \text{round}(x)| < \text{ulp}(x)/2 = 2^{-p+e}$

2.  $Er(x) = |(x - \text{round}(x)) / x|$

Dalla 1) e dato che  $x > 2^e$  si deriva che:

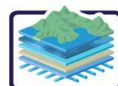
$$Er(x) < 2^{-p+e} / 2^e \text{ ossia } < 2^{-p} = \epsilon/2 \text{ (machine epsilon)}$$

$Er(x)$  è costante e va bene quindi per micro e macro numeri



POLITECNICO  
DI MILANO

Politecnico di Milano – DEI –  
Prof. Mauro Negri



**SpatialDBgroup**

Non va bene per domini applicativi nei quali il valore dipende da convenzioni

Es. Bounding Box Regione Lombardia del sistema UTM32/WGS84

- $X \in [459.973, 683.970]$   
 $\text{pi}(19\text{bit}) \text{ pf}(33) = 2^{19} * 2^{-52} = 0,1 * 10^{-6} \text{ mm}$
- $Y \in [4.949.981, 5.169.976]$   
 $\text{pi}(22\text{bit}) \text{ pf}(30) = 2^{22} * 2^{-52} = 0,9 * 10^{-6} \text{ mm}$

-----

Adesso è più chiaro perché....

**Esempio 1.** round superiore all'unità

```
#include <stdio.h>
```

```
int is=0, ix=7,i; float s=0.0, x=7.0;
```

```
void main()
```

```
{ for(i=1;i<=10000000; i++) {s=s+x; is=is+ix;}
```

```
printf("is= %d e s= %f",is,s);
```

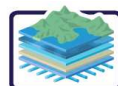
```
}
```

Risultato is= 70.000.000 e s= 77.603.248 (inizio differenza intorno a 16.777.216)



POLITECNICO  
DI MILANO

Politecnico di Milano – DEI –  
Prof. Mauro Negri



SpatialDBgroup

## Esempio 2. Denormalizzazione perde numeri piccoli

```
#include <stdio.h>
float a=0.0, b;
void main()
{   for(;;) {b=a; a=a+1.0; printf("n= %f e n+1= %f",b,a);
          if(b!=(a-1)) exit();
      }
}
```

Problema a 16777215

Oppure che: numeri distanti o molto vicini tra loro

$$x^2 - y^2 \neq (x-y) * (x+y)$$

## Esempio 3. Uguaglianza non sempre funziona

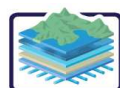
```
#include <stdio.h>
#include <math.h>
int main ()
{ float a=.1; float i; printf("\na inizio=%f\n", a);  $\Rightarrow$  0.1000001
  for (i=0.1; i!=10.0; i=i+0.1) {a=a+0.1; printf("\n%f", i);}
  ciclo  $\infty$ 
  passa da 9.90000210 a 10.00000210
```

Nota: Matlab introduce concetto di tolerance per l'uguaglianza



POLITECNICO  
DI MILANO

Politecnico di Milano – DEI –  
Prof. Mauro Negri



**SpatialDBgroup**

## Esempio 4 Cancellation problem in a-b

$$f'(x) = \frac{df(x)}{dx} = \lim_{\varepsilon \rightarrow 0} \frac{f(x + \varepsilon) - f(x)}{\varepsilon}$$

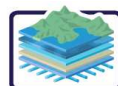
```
#include <stdio.h> ... <math.h>
int main()// derivata di sqrt(x)
{int c;double x, epsilon; printf("x?");scanf("%lf",&x);epsilon=1.;
while (epsilon > 1.e-17)
{printf("\n x= %f, epsilon= %e,derivata= %f",
      x, epsilon, (sqrt(x + epsilon)-sqrt(x))/epsilon);
  epsilon=epsilon/10.;
}
}
```

Risultato      x?1.

x= 1.000000, epsilon= 1.000000e+000,derivata=	0.414214
x= 1.000000, epsilon= 1.000000e -001,derivata=	0.488088
x= 1.000000, epsilon= 1.000000e -002,derivata=	0.498756
x= 1.000000, epsilon= 1.000000e -003,derivata=	0.499875
x= 1.000000, epsilon= 1.000000e -004,derivata=	0.499988
x= 1.000000, epsilon= 1.000000e -005,derivata=	0.499999
x= 1.000000, epsilon= 1.000000e -006,derivata=	0.500000
x= 1.000000, epsilon= 1.000000e -007,derivata=	0.500000
x= 1.000000, epsilon= 1.000000e -008,derivata=	0.500000
x= 1.000000, epsilon= 1.000000e -009,derivata=	0.500000
x= 1.000000, epsilon= 1.000000e -010,derivata=	0.500000
x= 1.000000, epsilon= 1.000000e -011,derivata=	0.500000
x= 1.000000, epsilon= 1.000000e -012,derivata=	0.500044
x= 1.000000, epsilon= 1.000000e -013,derivata=	0.499600
x= 1.000000, epsilon= 1.000000e -014,derivata=	0.488498
x= 1.000000, epsilon= 1.000000e -015,derivata=	0.444089
x= 1.000000, epsilon= 1.000000e -016,derivata=	0.000000



POLITECNICO  
DI MILANO



SpatialDBgroup

Politecnico di Milano – DEI –  
Prof. Mauro Negri

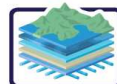
## Esempio 5

The screenshot shows the Microsoft technical support website. At the top, there's a navigation bar with 'Vai alla Gallery Add-on', 'Italia', 'Cambia', and 'Tutti i siti Microsoft'. Below this is the 'Supporto tecnico Microsoft' header. A search bar contains the text 'Cerca in Supporto tecnico Microsoft'. The main navigation menu includes 'Home page Supporto tecnico', 'Centri di supporto', 'Ricerca avanzata', and 'Acquista prodotti'. The article title is 'Operazioni aritmetiche dei valori in virgola mobile potrebbe produrre risultati non accurati in Excel'. It mentions the article ID 78113, last modified on May 13, 2013, and version 7.0. A note indicates the article is automatically translated from English. On the right, there's a sidebar with 'Altre risorse' (Other resources) including 'Altri siti di supporto', 'Community', and 'Richiedi assistenza'. Below that is 'Traduzione articoli' (Article translation) set to 'Inglese (Stati Uniti)'. Further down is 'Centri di supporto tecnico correlati' (Related technical support centers) listing Excel 2010, Excel, Excel 2007, Excel 2003, and Excel 2000. The main content area has a section titled 'Visualizza i prodotti a cui si riferisce l'articolo,' (View the products to which the article refers,) with buttons for 'In questa pagina' (On this page), 'Espandi tutto / Chiudi tutto' (Expand all / Collapse all), and 'Somma' (Summary).



POLITECNICO  
DI MILANO

Politecnico di Milano – DEI –  
Prof. Mauro Negri



SpatialDBgroup