

Modulo: Approfondimenti sui Sistemi Aritmetici di un computer: tipo reale [P2_03]

Unità didattica: Sistema Aritmetico Reale Floating-Point [1-AT]

Titolo: Rappresentazione in memoria dei numeri reali

Argomenti trattati:

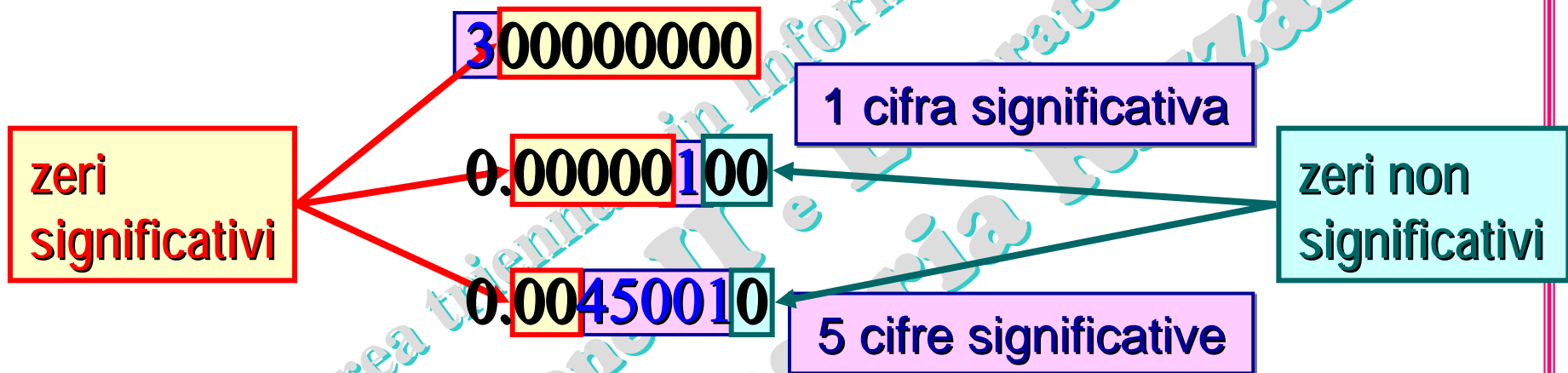
- ✓ Notazione scientifica dei numeri reali (segno, mantissa, esponente)
- ✓ Rappresentazione binaria della mantissa a bit implicito
- ✓ Sistema Aritmetico Floating Point IEEE Standard 754
- ✓ Oggetti del S.A. Standard e loro caratterizzazione

Prerequisiti richiesti: aritmetica binaria, sistema arit-metico intero

Notazione scientifica

La notazione più compatta per rappresentare i **numeri reali** è quella **scientifica** in cui si rappresentano solo le *cifre significative*.

Esempi



Nella **notazione scientifica** un numero reale x si rappresenta tramite una **mantissa** m , contenente le cifre significative, ed un **esponente** p , indicazione degli zeri significativi, tali che $x = m \cdot \beta^p$ dove β è la **base** del sistema di numerazione.

			notazione scientifica
Esempi:	300000000	$= 3 \cdot 10^8$	$= 3e+8$
	0.00000100	$= 1 \cdot 10^{-6}$	$= 1e-6$
	0.00450010	$= 4.5001 \cdot 10^{-3}$	$= 4.5001e-3$

Per individuare **univocamente** la notazione scientifica di un numero, si opera la **normalizzazione** della mantissa da cui discende il valore dell'esponente:

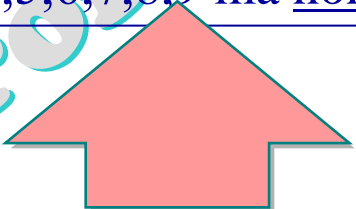
Esempio: 2 schemi di normalizzazione (base $\beta=10$)

356000000	=	3.5600 · 10 ⁺⁸	=	0.35600 · 10 ⁺⁹
0.00000123	=	1.2300 · 10 ⁻⁶	=	0.12300 · 10 ⁻⁵
0.00450010	=	4.5001 · 10 ⁻³	=	0.45001 · 10 ⁻²

1

mantissa m tale che
 $1 \leq m < \beta=10$

la cifra delle unità può essere
 1,2,3,4,5,6,7,8,9 ma non 0



...oppure...

2

mantissa m tale che
 $0.1=\beta^{-1} \leq m < 1$

la cifra dei decimi può essere
 1,2,3,4,5,6,7,8,9 ma non 0

Esempio: base $\beta=2$

bit esplicito

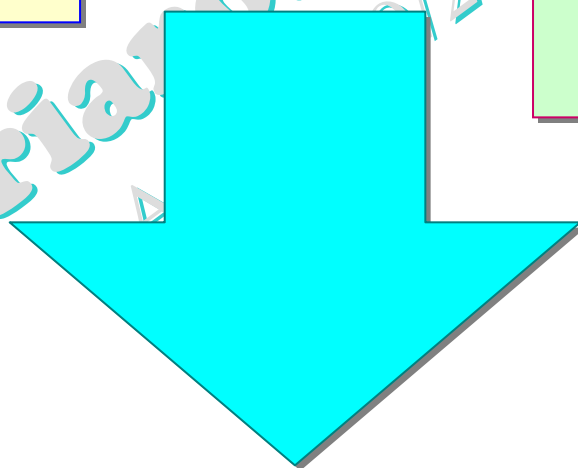
1 10110000	=	1.1011 · 2 ⁺¹⁰⁰⁰	=	0.11011 · 2 ⁺¹⁰⁰¹
0.000001011	=	1.0110 · 2 ⁻⁰¹¹⁰	=	0.10110 · 2 ⁻⁰¹⁰¹
0.0010111	=	1.0111 · 2 ⁻⁰⁰¹¹	=	0.10111 · 2 ⁻⁰⁰¹⁰

1

mantissa m tale che
 $1 \leq m < \beta=2$

2

mantissa m tale che
 $1/2=\beta^{-1} \leq m < 1$



la prima cifra della mantissa può essere solo 1

La prima cifra (significativa) della mantissa è sempre 1
(allora può **non** essere memorizzata!)



bit implicito

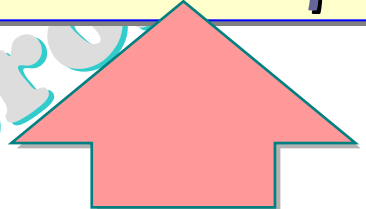
110110000	=	1 1011 · 2 ⁺¹⁰⁰⁰	=	0.1 1011 · 2 ⁺¹⁰⁰¹
0.000001011	=	1 .0110 · 2 ⁻⁰¹¹⁰	=	0.1 0110 · 2 ⁻⁰¹⁰¹
0.0010111	=	1 .0111 · 2 ⁻⁰⁰¹¹	=	0.1 0111 · 2 ⁻⁰⁰¹⁰

1

mantissa m tale che
 $1 \leq m < \beta$

2

mantissa m tale che
 $\beta^{-1} \leq m < 1$



SISTEMA ARITMETICO REALE FLOATING-POINT

Sistema Aritmetico Floating Point

$$F(\beta, t, E_{\min}, E_{\max})$$

Denota l'insieme dei numeri reali rappresentati in un computer e le operazioni definite su di essi.

Parametri del sistema aritmetico $F(\beta, t, E_{\min}, E_{\max})$:

β = base del sistema di numerazione

t = numero di cifre β per la mantissa

$E_{\min} < E_{\max}$ = limitazioni per il campo esponente

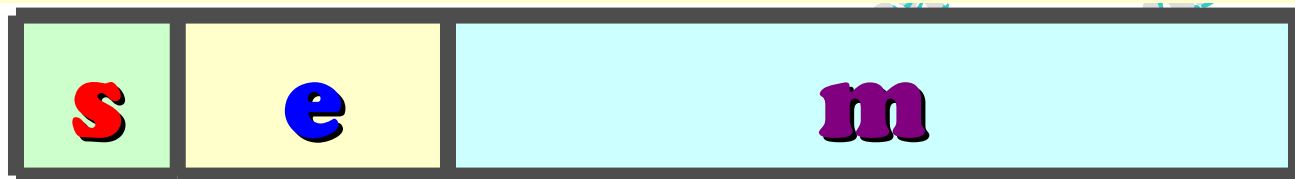
Sistema Aritmetico Floating Point IEEE Standard 754 (1985)

$$\beta = 2$$

Formato	Tipo	num. bit	ling. C
Basic	<i>single</i>	32 (4 byte)	<i>float</i>
	<i>double</i>	64 (8 byte)	<i>double</i>
Extended	<i>double</i>	80 (10 byte)	<i>long double</i>

opzionale

Nel **S.A. IEEE Std. F(2, t, E_{min}, E_{max})** un numero (normalizzato) del **formato Basic** è rappresentato in memoria come:



dove

segno

esponente

mantissa

● **s** denota il *segno della mantissa*;

● l'*esponente e* è rappresentato come "intero biased";

● la *mantissa m* è rappresentata con **s** per segno e modulo su **t** bit a *bit implicito* ℓ^\dagger (*precisione* = **t**+1) ed è generata con lo schema del *round to nearest*;

ed il suo **valore** è pertanto $x = (-1)^s [\ell.m] \times 2^{e - \text{Bias}}$

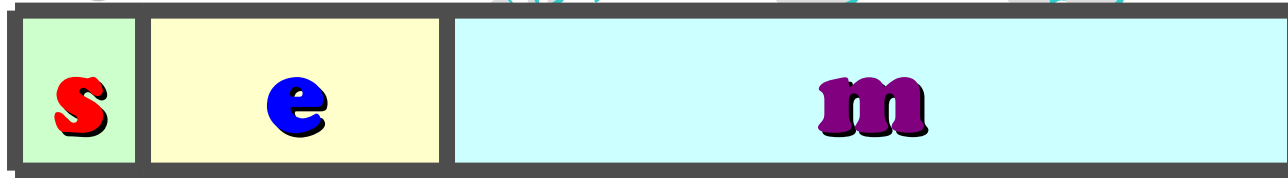
\dagger il primo bit ℓ non è rappresentato esplicitamente, ma assume un valore convenzionale: pertanto è come se la mantissa fosse rappresentata su **t**+1 bit.

Perché i campi si susseguono in questo modo?

È così assicurato l'ordinamento dei numeri!!!

(considerando tutti i bit insieme come se fosse un intero)

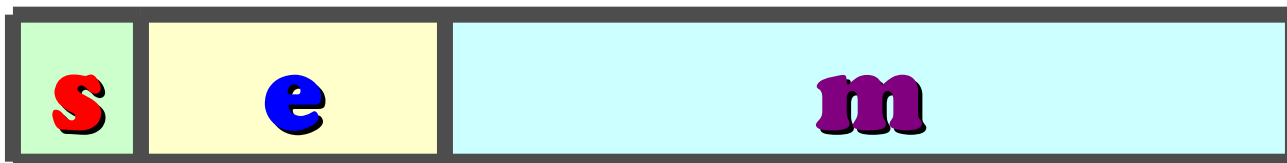
bit +signif. ← → bit -signif.



... a parità di segno ...

l'ordine di grandezza predomina

I negativi precedono i positivi



IEEE Standard 754

Tipo	numero bit esponente	Bias	t = numero bit mantissa
<i>single</i>	8 $e \in \{0..255\}$	127	23
<i>double</i>	11 $e \in \{0..2047\}$	1023	52
<i>double EXTEND</i>	15 $e \in \{0..32767\}$	16383	64

E_{min}

E_{max}

senza bit implicito

opzionale

II **S.A. IEEE Standard 754** stabilisce:

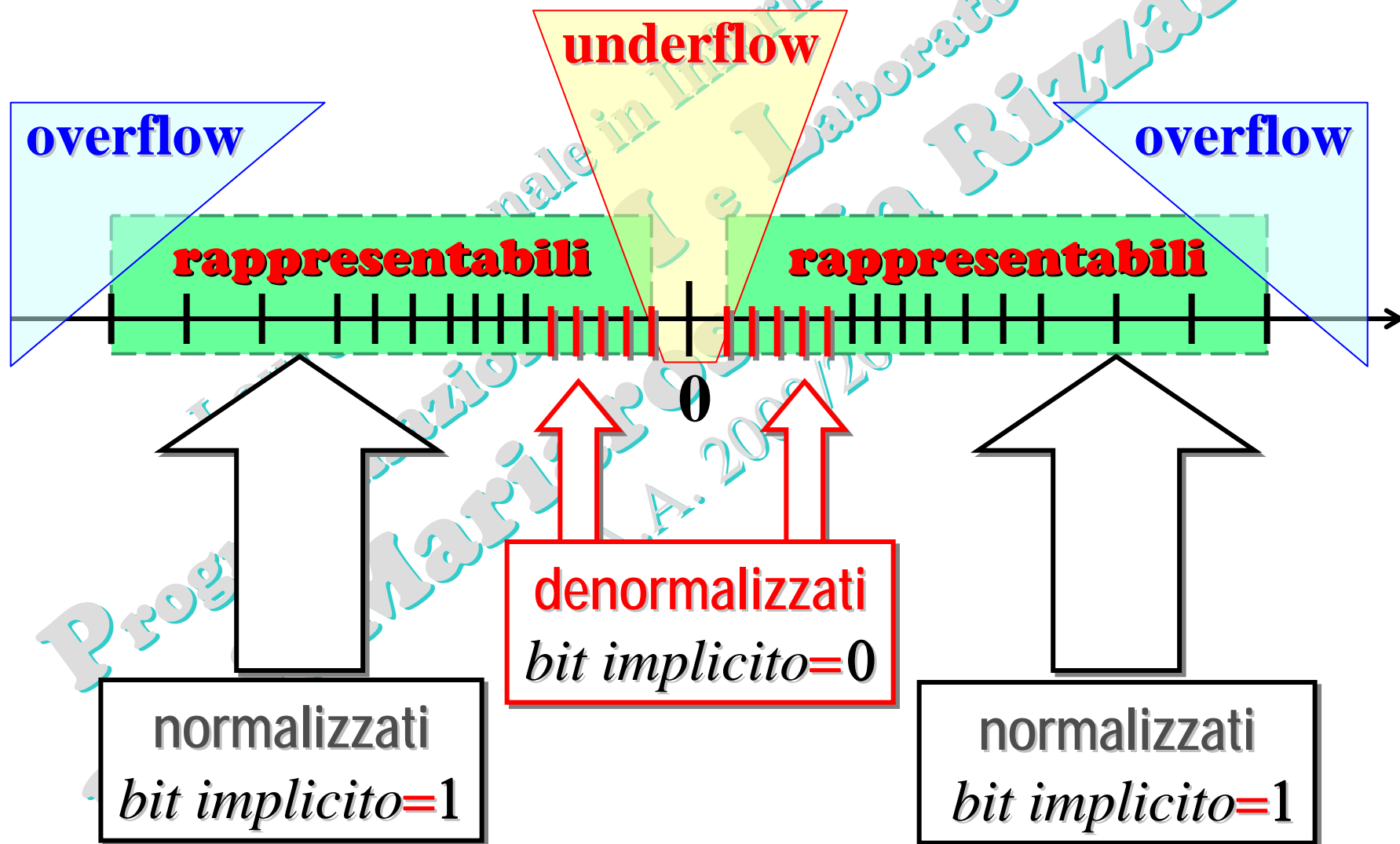
gli oggetti del *Sistema Aritmetico Floating-Point*;
le operazioni definite su ciascun oggetto.

valore bit
implicito

Oggetto	Caratterizzazione		ℓ
	esponente e	mantissa m	
Numeri normalizzati valore = $(-1)^s [\ell.m] \times 2^{e-\text{Bias}}$	$E_{\min} < e < E_{\max}$	$m \geq 0$	1
Infinito con segno	$e = E_{\max}$	$m = 0$	-
NaN (Not A Number)	$e = E_{\max}$	$m \neq 0$	-
Zero con segno	$e = E_{\min}$	$m = 0$	-
Numeri denormalizzati valore = $(-1)^s [\ell.m] \times 2^{e-\text{Bias}+1}$	$e = E_{\min}$	$m \neq 0$	0

I numeri Floating Point non sono uniformemente distribuiti sull'asse reale.

Esempio: numeri del **Formato Basic**



Esempio

$F(\beta=2, t=2, E_{\min}=0, E_{\max}=3)$

2 bit per la mantissa

n = 2 bit per l'esponente
(Bias = 1)

s	e		valore espon.	m		ℓ	valore numero
0	0	0	-	0	0	0	+ zero
0	0	0	0	0	1	0	0.01 2 ⁰ = +1/4
0	0	0	0	1	0	0	0.10 2 ⁰ = +1/2
0	0	0	0	1	1	0	0.11 2 ⁰ = +3/4
0	0	1	0	0	0	1	1.00 2 ⁰ = +1
0	0	1	0	0	1	1	1.01 2 ⁰ = +5/4
0	0	1	0	1	0	1	1.10 2 ⁰ = +3/2
0	0	1	0	1	1	1	1.11 2 ⁰ = +7/4
...							

valore = $(-1)^s [\ell.m] \times 2^{E_{\min} - \text{Bias} + 1}$
denormalizzati

realmin

valore = $(-1)^s [\ell.m] \times 2^{\text{Bias}}$
normalizzati

...								normalizzati	
s	e		valore espon.	m		<i>l</i>	valore numero	$\text{valore}=(-1)^{\textcolor{red}{s}}[\textcolor{blue}{l}.\textcolor{violet}{m}]\times 2^{\textcolor{blue}{e}-\text{Bias}}$	
0	1	0	+1	0	0	1	1.00 $2^{+1}=\textcolor{red}{+2}$	<div><i>realmax</i></div>	
0	1	0	+1	0	1	1	1.01 $2^{+1}=\textcolor{red}{+5/2}$		
0	1	0	+1	1	0	1	1.10 $2^{+1}=\textcolor{red}{+3}$		
0	1	0	+1	1	1	1	1.11 $2^{+1}=\textcolor{red}{+7/2}$		
0	1	1		0	0		Inf		
0	1	1		0	1		} NaN		
0	1	1		1	0				
0	1	1		1	1				
1	0	0	-	0	0	0	- zero	<div>NEGATIVI</div>	
1	0	0	0	0	1				
1				

$$F(\beta=2, t=2, E_{\min}=0, E_{\max}=3)$$

2 bit per la
mantissa

2 bit per
l'esponente
(bias = 1)

realmin

realmax

underflow

overflow

overflow

rappresentabili

rappresentabili

-3.5 -3 -2 -1 0 +1 +2 +3 +3.5

denormaliz.

$$F(\beta=2, t=2, E_{\min}=0, E_{\max}=3)$$