



# Firmware Division, UC & Floating pointer adder

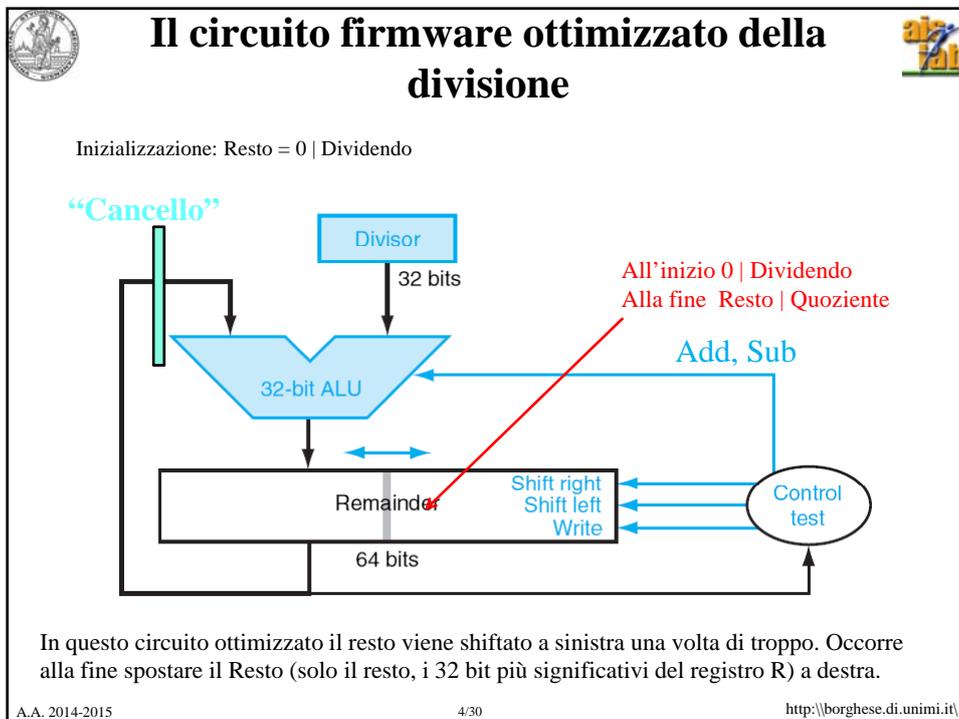
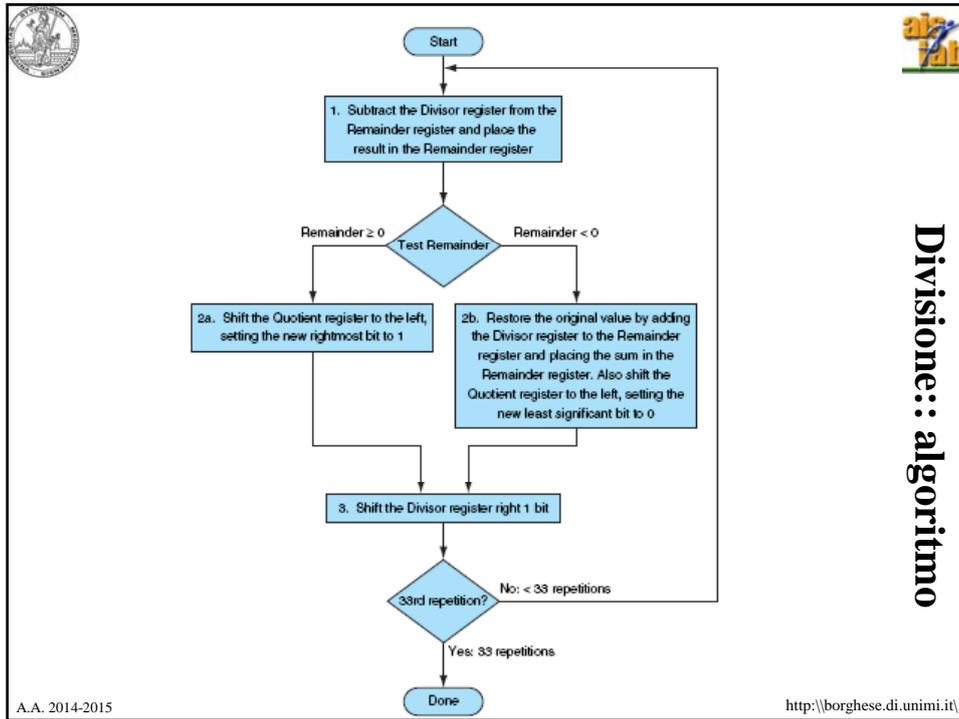
Prof. Alberto Borghese  
Dipartimento di Scienze dell'Informazione  
[borgnese@di.unimi.it](mailto:borgnese@di.unimi.it)

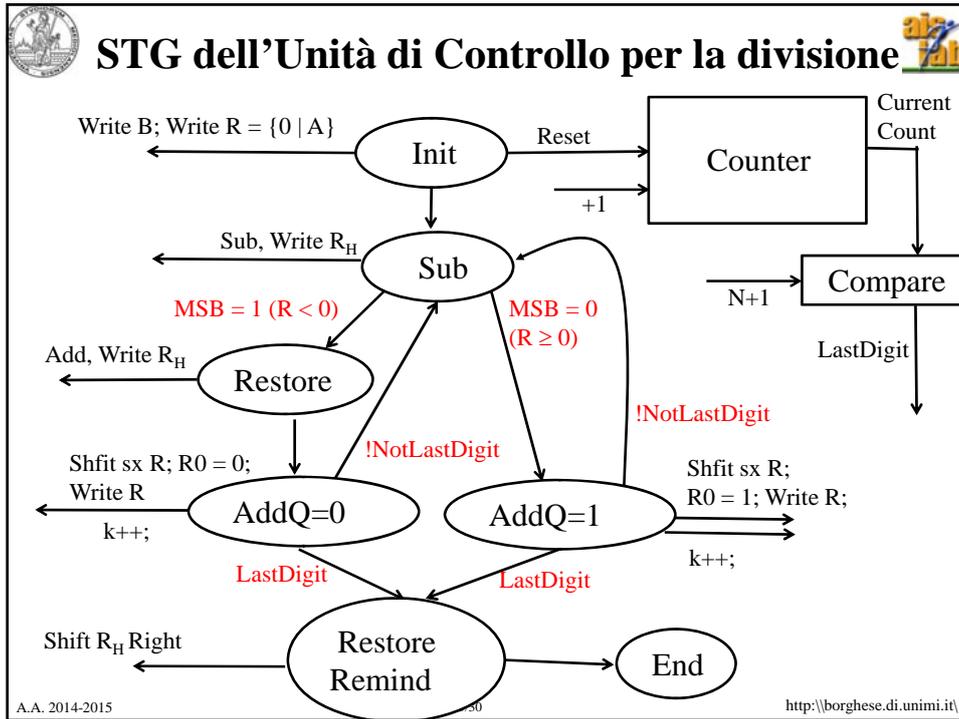
Università degli Studi di Milano  
Riferimenti sul Patterson, 5a Ed.: 3.4, 3.5



## Sommario

- UC Divisione intera
- Somma in virgola mobile







## Definizione del funzionamento della FSM: la STT



S \ I	MSB=0 LastDigit	MSB = 1 LastDigit	MSB=0 !LastDigit	MSB = 1 !LastDigit
Init	Sub	Sub	Sub	Sub
Sub	AddQ=1	Restore	AddQ=1	Restore
Restore	AddQ=0	AddQ=0	AddQ=0	AddQ=0
Add Q=0	Restore R <sub>H</sub>	Restore R <sub>H</sub>	Sub	Sub
Add Q=1	Restore R <sub>H</sub>	Restore R <sub>H</sub>	Sub	Sub
Restore R <sub>H</sub>	End	End	End	End
End	End	End	End	End

$$\text{Next State } S_{t+1} = f(S_t, I_t)$$



## Definizione del funzionamento della FSM: la STT

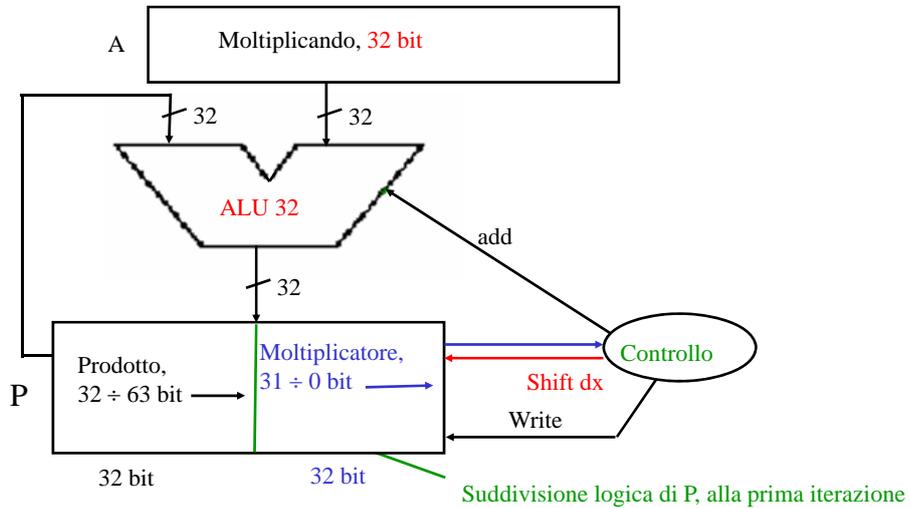


S \ I	Write B	Write R={0   A}	Reset count	Sub	Add	Shift R left	Write R0=0	Write R0=1	Write R <sub>H</sub>	Shift R <sub>H</sub> right	k++
Init	1	1	1	0	0	0	0	0	0	0	0
Sub	0	0	0	1	0	0	0	0	1	0	0
Restore	0	0	0	0	1	0	0	0	1	0	0
Add Q=0	0	0	0	0	0	1	1	0	0	0	1
Add Q=1	0	0	0	0	0	1	0	1	0	0	1
Restore R <sub>H</sub>	0	0	0	0	0	0	0	0	0	1	1
End	0	0	0	0	0	0	0	0	0	0	0

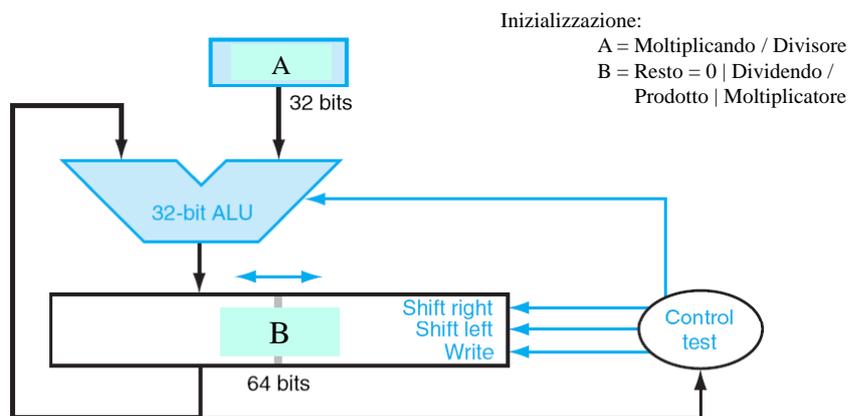
$$\text{Uscita } Y = g(S)$$



## Circuito ottimizzato della moltiplicazione



## Un unico circuito per moltiplicazione e divisione





## Definizione della struttura della FSM



$\langle S_0, S, I, Y, f(S,I), g(S) \rangle$

$S = \{ \text{Init, Sub, Restore, AddQ=1, AddQ=0, Restore } R_H, \text{End} \}$

$S_0 = \text{Init}$

$I = \{ \text{MSB, LastDigit} \}$

$Y = \{ \text{Write B, Write } R = \{0 | A\}, \text{Sub, Add, Write } R_H, R_0 = 0; R_0 = 1; k++; \text{Shift R left, Shift } R_H \text{ right} \}$

$f(S,I) = ?$

$g(S) = ?$



## Unità di controllo per divisione / moltiplicazione



$\langle S_0, S, I, Y, f(S,I), g(S) \rangle$

$S = \{ \text{Init, Sub, Restore, AddQ=1, AddQ=0, Restore } R_H, \text{Shift dx, End} \}$

$S_0 = \text{Init}$

$I = \{ \text{MSB, LastDigit, Multiply/Divide} \}$

$Y = \{ \text{Write B, Write A, Write } R = \{0 | A\}, \text{Write } R = \{0 | B\}, \text{Sub, Add, None, Write } R_H, R_0 = 0; R_0 = 1; \text{Shift R left, Shift } R_H \text{ right, } k++ \}$

$f(S,I) = ?$

$g(S) = ?$



## Sommario



- UC Divisione intera
- **Somma in virgola mobile**



## Codifica in virgola mobile Standard IEEE 754 (1980)

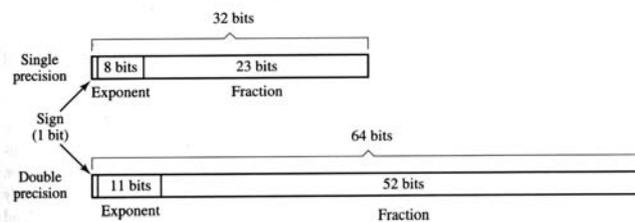


Figure 2-10 Single-precision and double-precision IEEE 754 floating point formats.

Rappresentazione polarizzata dell'esponente:

Polarizzazione pari a 127 per singola precisione =>  
1 viene codificato come 1000 0000.

Polarizzazione pari a 1023 in doppia precisione.  
1 viene codificato come 1000 0000 000.



## Esempio di somma in virgola mobile



$$a = 9,999 \times 10^1 \quad b = 1,61 \times 10^{-1} \quad a + b = ?$$

NB I numeri decimali sono normalizzati.

Una possibilità è:

$$\begin{array}{r} 99,99 \quad + \\ 0,161 \quad = \\ \hline \end{array}$$

100,151

100,151  $\rightarrow$  1,0015  $\times 10^2$  in forma normalizzata



## Algoritmo di somma in virgola mobile - I



1) Trasformare i due numeri in modo che le due rappresentazioni abbiano la stessa base: allineamento della virgola. Quale si allinea?

2) Effettuare la somma delle mantisse.

**Se il numero risultante è normalizzato termino qui. Altrimenti:**

3) Normalizzare il risultato.



## Algoritmo di somma in virgola mobile - II



- 1) Trasformare **uno dei due numeri** in modo che le due rappresentazioni abbiano la stessa base: allineamento della virgola. Si allinea all'esponente più alto (denormalizzo il numero più piccolo).
- 2) Effettuare la somma delle mantisse.

**Se il numero risultante è normalizzato termino qui. Altrimenti:**

- 3) Normalizzare il risultato.



## Esempio di somma in virgola mobile - II



$$a = 9,999 \times 10^1 \quad b = 1,61 \times 10^{-1} \quad a + b = ?$$

Supponiamo di avere a disposizione 4 cifre per la mantissa e due per l'esponente.

Devo trasformare uno dei numeri, una possibilità è:

$$\begin{array}{r} 9,999 \times 10^1 + \\ 0,016 \times 10^1 = \end{array} \quad \text{Perdo un bit perchè non rientra nella capacità della mantissa}$$


---


$$10,015 \times 10^1 \quad \text{Il risultato non è più normalizzato, anche se i due addendi sono normalizzati.}$$

Normalizzazione per ottenere il risultato finale:

$$10,015 \times 10^1 = 1,001 \times 10^2 \text{ in forma normalizzata.}$$

**NB:** In questa fase si può generare la necessità di rinormalizzare il numero (passo 3):

$$\begin{array}{l} \text{Esempio: } 9,99999 \times 10^2 \Rightarrow \text{Arrotondo alla cifra più vicina} \Rightarrow 10,00 \times 10^3 \\ \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \Rightarrow \text{Normalizzazione} \qquad \qquad \qquad \Rightarrow 1,0 \times 10^4 \end{array}$$



## Approssimazione



Interi -> risultato esatto (o overflow)

Numeri decimali -> Spesso occorrono delle approssimazioni

- Troncamento (floor):  $1,001 \times 10^2$  (cf. Slide precedente)
- Arrotondamento alla cifra superiore (ceil):  $1,002 \times 10^2$
- Arrotondamento alla cifra più vicina:  $1,002 \times 10^2$

IEEE754 prevede 2 bit aggiuntivi nei calcoli per mantenere l'accuratezza.

bit di guardia (guard)

bit di arrotondamento (round)



## Esempio: aritmetica in floating point accurata



$$a = 2,34 \quad b = 0,0256$$

$$a + b = ?$$

Codifica su 3 cifre decimali totali.

Approssimazione mediante arrotondamento.

Senza cifre di arrotondamento devo scrivere:

$$2,34 +$$

$$0,02 =$$

-----

$$2,36$$

Con il bit di guardia e di arrotondamento posso scrivere:

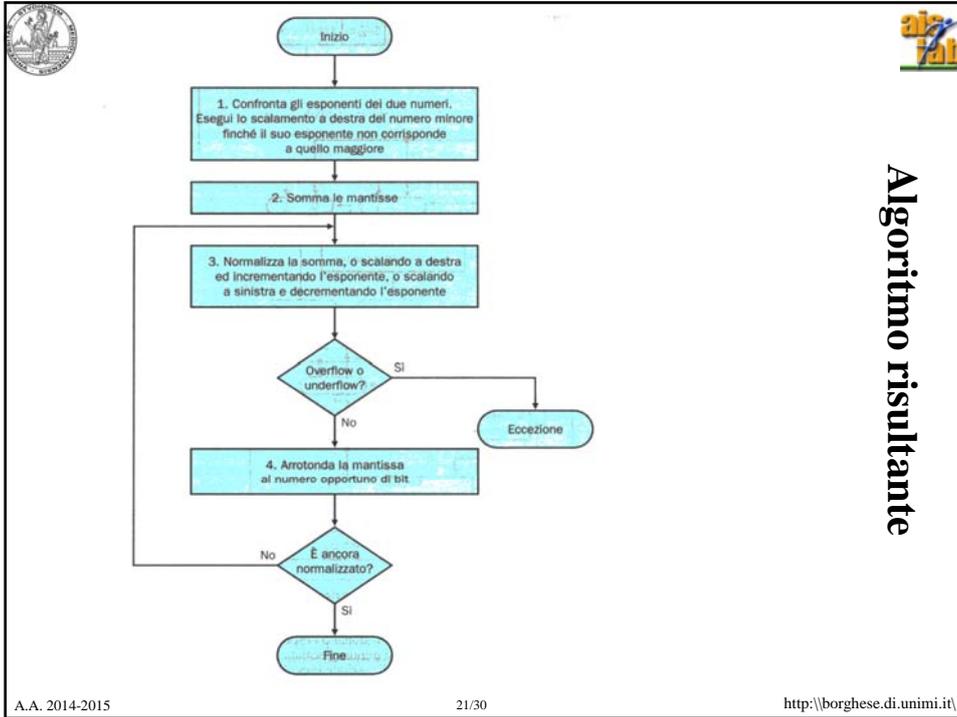
$$2,3400 +$$

$$0,0256 =$$

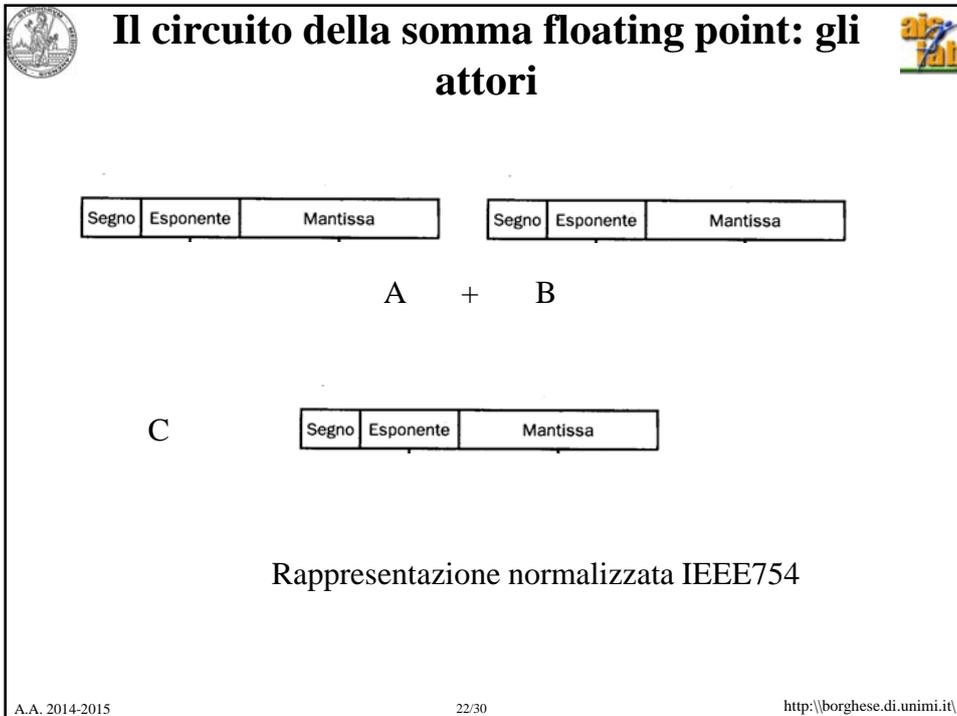
-----

$$2,3656$$

L'arrotondamento finale fornisce per rientrare in 3 cifre decimali fornisce: **2,37**

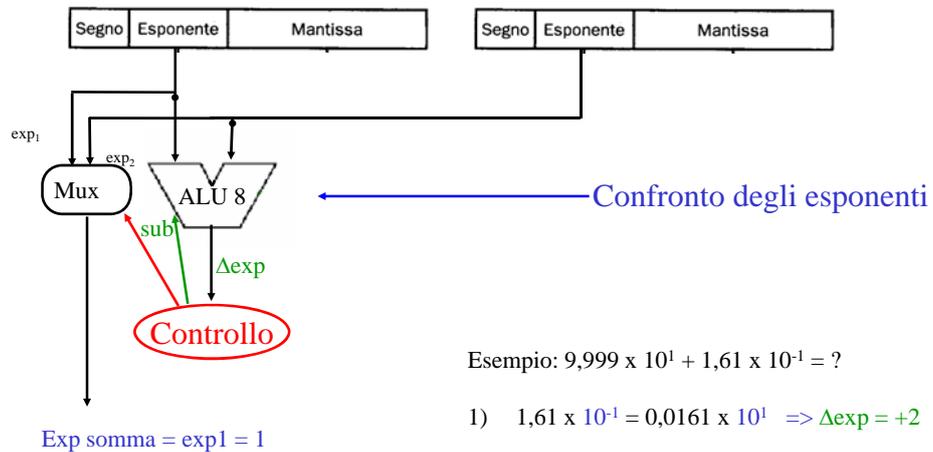


Algoritmo risultante

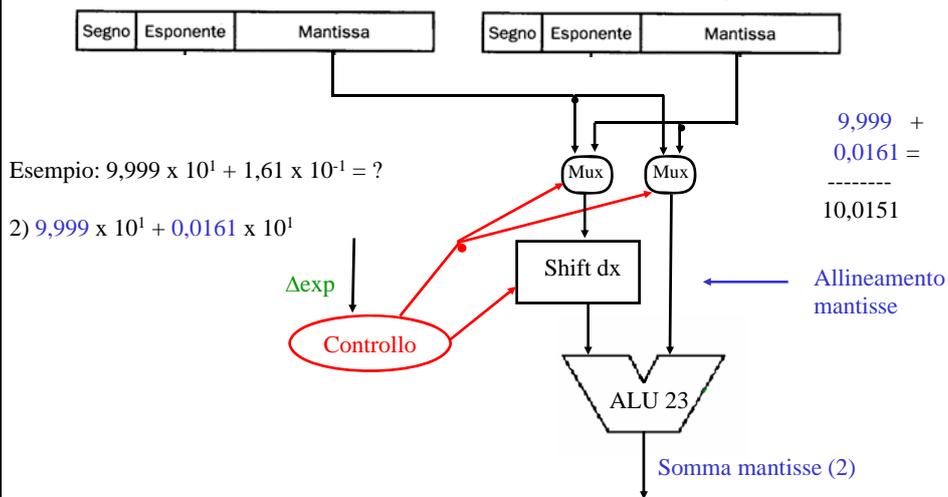


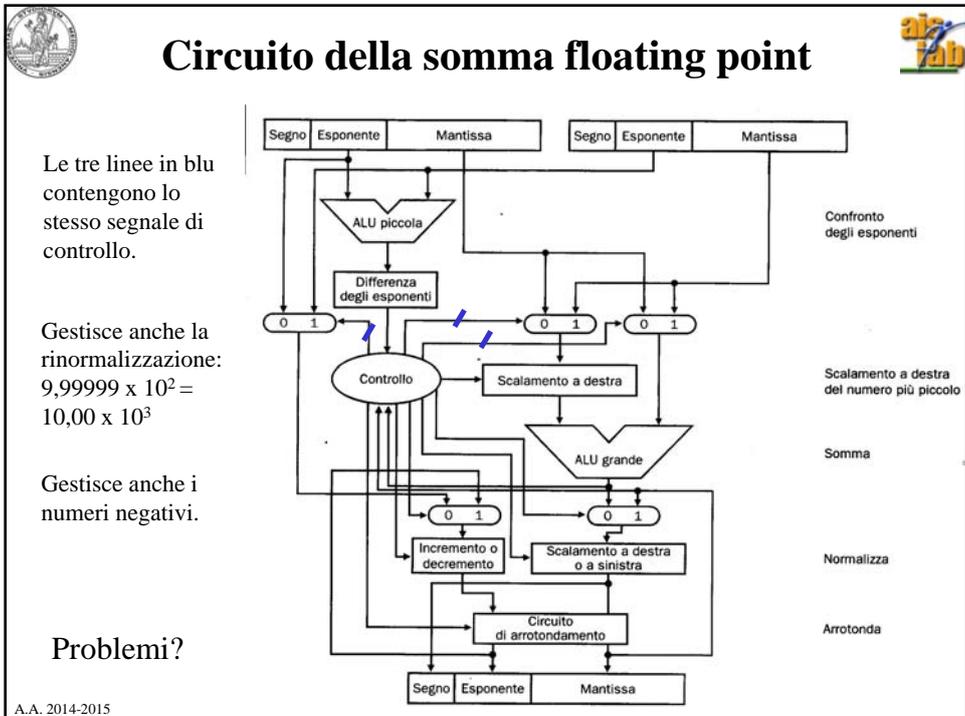
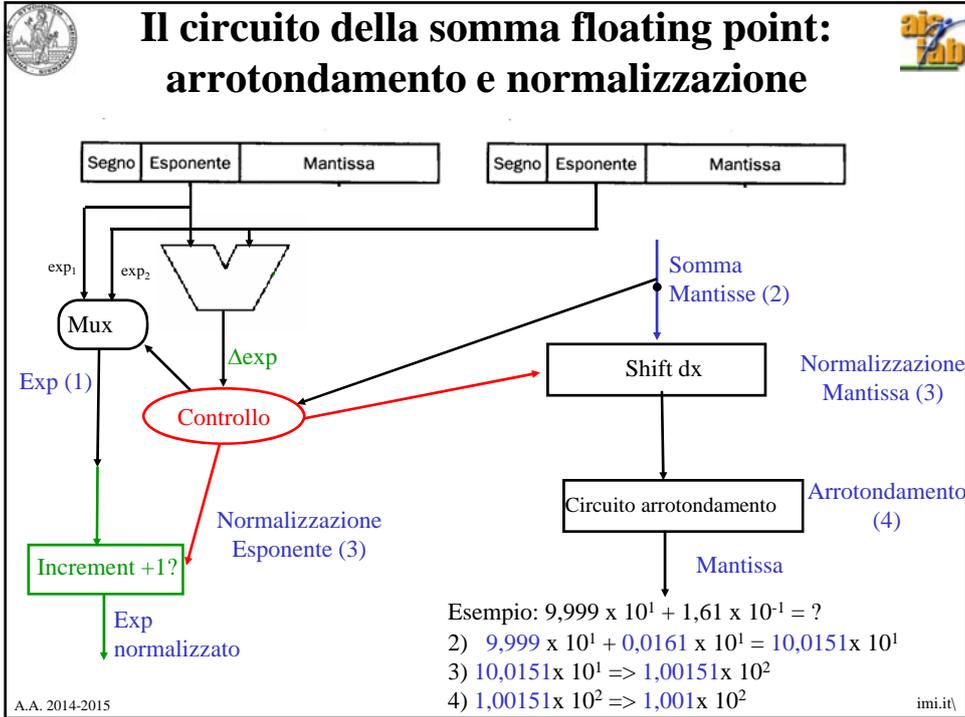


## Il circuito della somma floating point: determinazione dell'esponente comune



## Il circuito della somma floating point: allineamento delle mantisse e somma







## Esempio



$P = 0,5 - 0,4375$  somma binaria con precisione di 4 bit.

$A = 1,000 \times 2^{-1}$   
 $B = -1,110 \times 2^{-2}$  Espressione in forma normalizzata.

- 1) Allineamento di A e B. Trasformo B:  $-1,110 \times 2^{-2} = -0,1110 \times 2^{-1}$
- 2) Somma delle mantisse:  $1,000 + (-0,111) = 0,001 \times 2^{-1}$
- 3) Normalizzazione della somma:  $0,001 \times 2^{-1} = 1,000 \times 2^{-4}$   
Controllo di overflow o underflow:  $-126 \leq -4 \leq 127$ .
- 4) Arrotondamento della somma: 1,000 non lo richiede.

$1,000 \times 2^{-4} = 1/2^4 = 1/16 = 0,0625$  c.v.d.



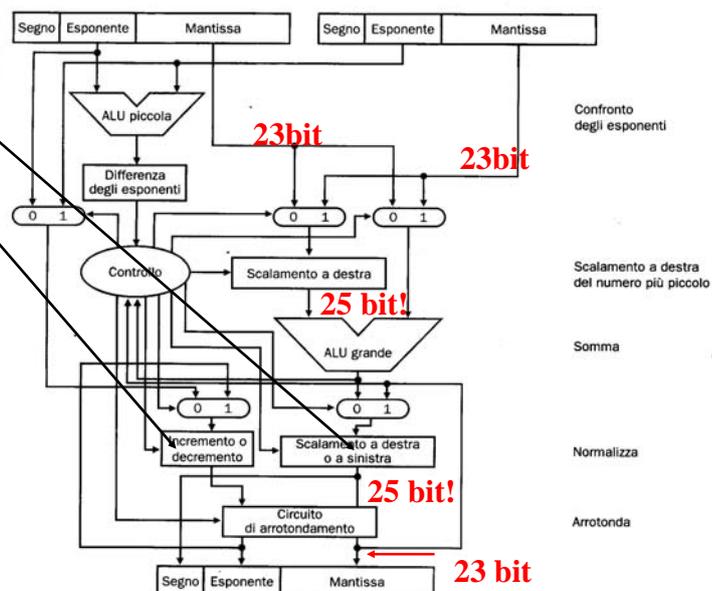
## Circuito della somma floating point con bit di arrotondamento



In quale caso la mantissa viene scalata a sx?

In quale caso l'esponente viene decrementato?

La rappresentazione interna, secondo IEEE 754, prevede 2 bit aggiuntivi: **bit di guardia** e **bit di arrotondamento**.





## Prodotto e divisione in virgola mobile



- Prodotto delle mantisse
- Somma degli esponenti
- Normalizzazione
  
- Divisione in virgola mobile = Prodotto di un numero per il suo inverso.



## Sommario



- UC Divisione intera
- Somma in virgola mobile