



# Firmware Division, UC & Floating pointer adder

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

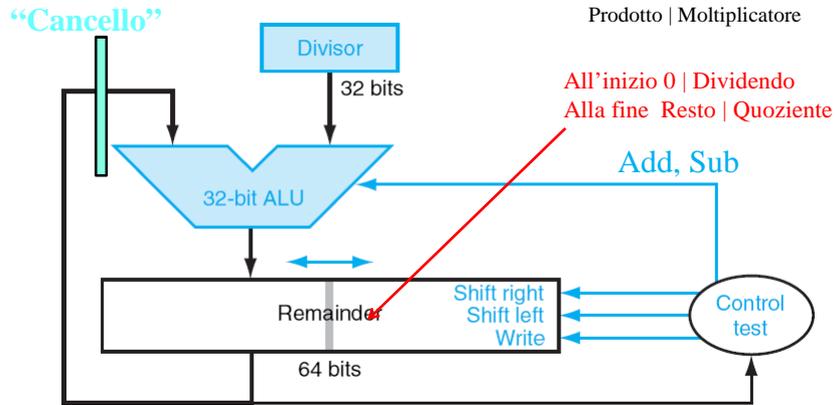


# Il circuito firmware ottimizzato della divisione

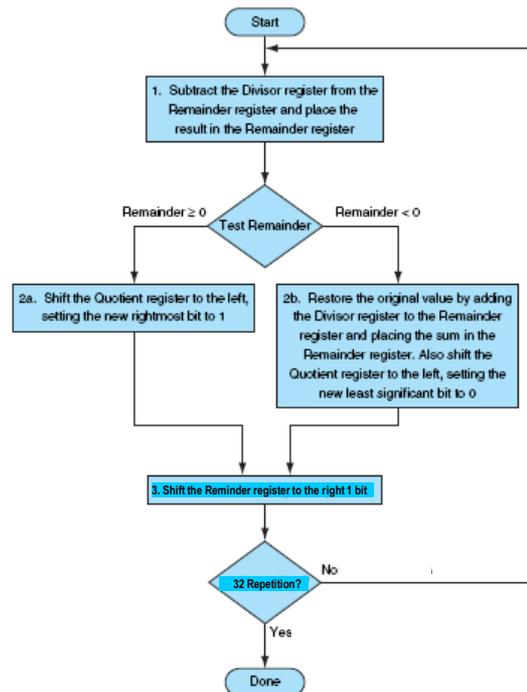


Inizializzazione:

A = Moltiplicando / Divisore  
 B = Resto = 0 | Dividendo (Quoziente)  
 Prodotto | Moltiplicatore



In questo circuito ottimizzato il resto viene shiftato a sinistra una volta di troppo. Occorre alla fine spostare il Resto (solo il resto, i 32 bit più significativi del registro R) a destra.



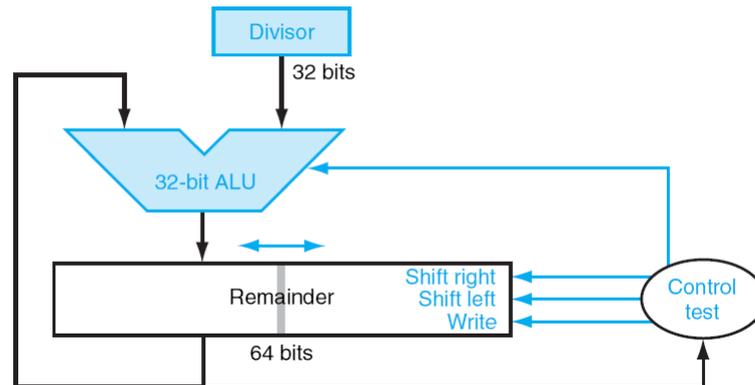
Divisione:: algoritmo



## Esempio su 4 bit

$9 : 2 = 4$  con resto di 1

$2 * 4 = 8$



## Il segno della divisione

Dividendo = Quoziente x Divisore + Resto

|   |   |   |
|---|---|---|
| + | + | + |
| + | - | - |
| - | + | - |
| - | - | + |

Qual'è il segno del resto?

Esempio:

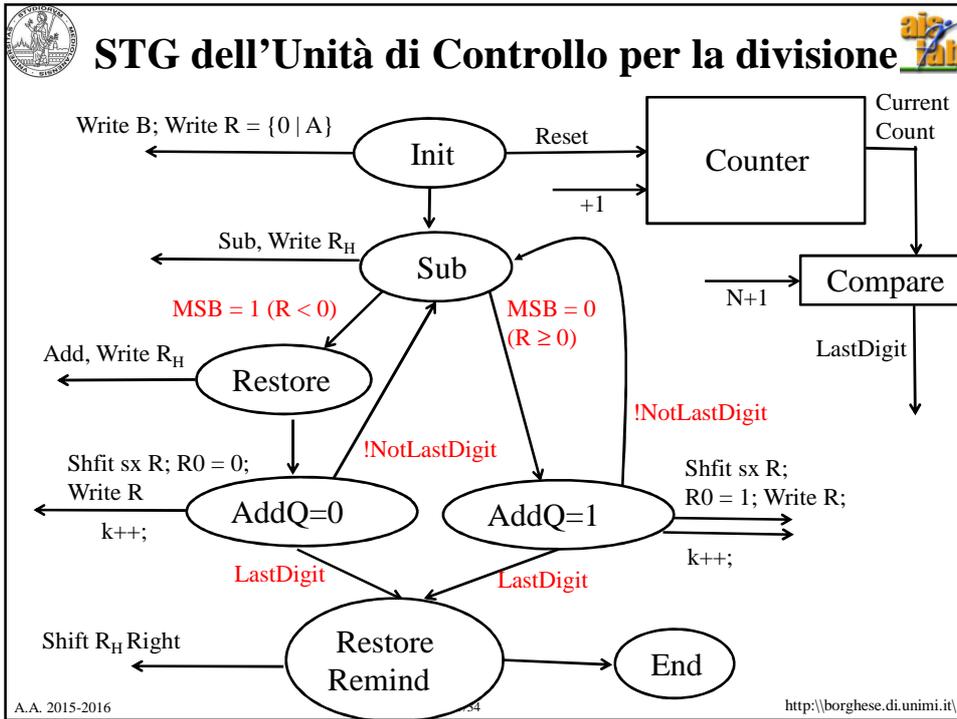
$$+7 = +2 \times +3 + 1$$

$$-7 = -2 \times +3 - 1 \quad (\text{sarebbe uguale anche a } -2 \times +4 + 1!!, \text{ ma cambierebbe il quoziente in funzione del segno del resto!})$$

$$+7 = -2 \times -3 + 1$$

$$-7 = +2 \times -3 - 1$$

Osserviamo che il quoziente è positivo se i segni di Divisore e Dividendo sono concordi.  
Il resto ha il segno del Dividendo.





## Definizione del funzionamento della FSM: la STT



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

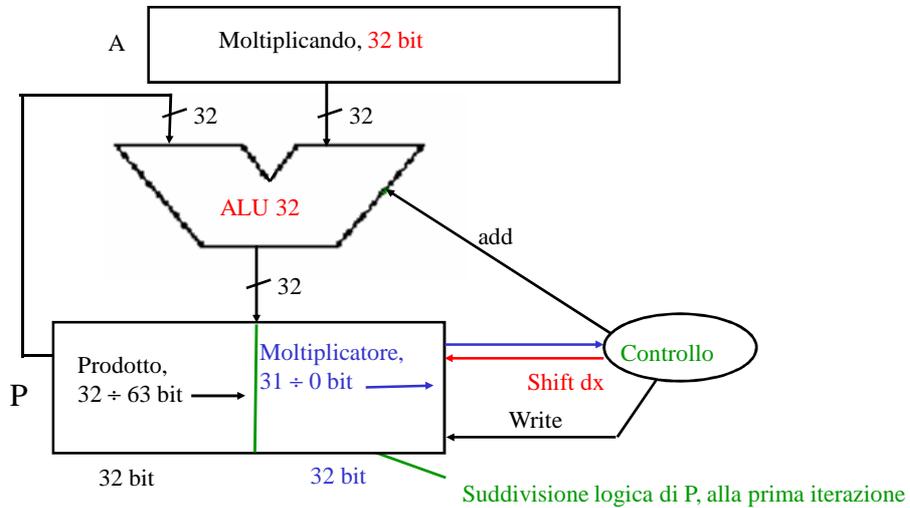


| S \ I                     | Write<br>B | Write<br>R={0   A} | Reset<br>count | Sub | Add | Shift<br>R left | Write<br>R0=0 | Write<br>R0=1 | Write<br>R <sub>H</sub> | Shift<br>R <sub>H</sub> right | k++ |
|---------------------------|------------|--------------------|----------------|-----|-----|-----------------|---------------|---------------|---------------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

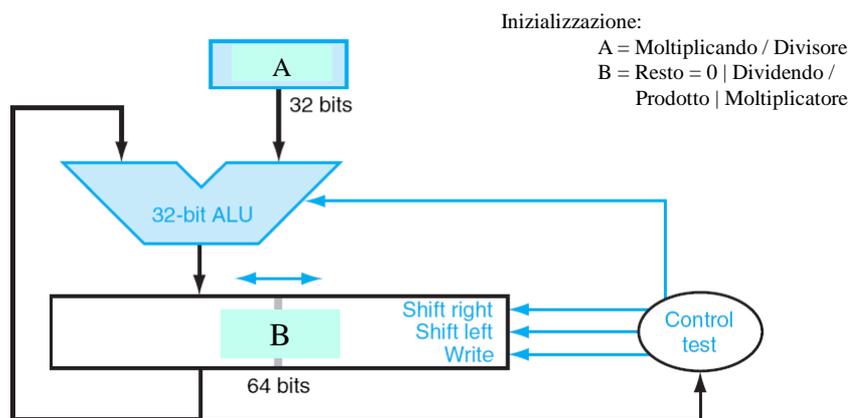
$$\text{Uscita } Y = g(S)$$



## 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)

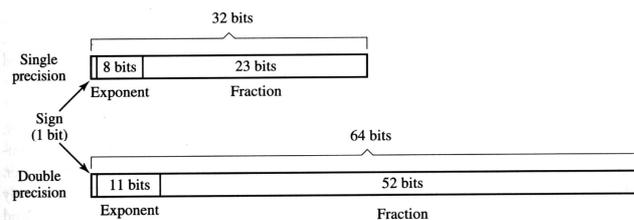


Figure 2-10 Single-precision and double-precision IEEE 754 floating point formats.

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



$$a = 7,999 \times 10^1 \quad b = 1,61 \times 10^{-1} \quad a + b = ?$$

NB I numeri decimali sono normalizzati.

Una possibilità è:

$$\begin{array}{r} 79,99 \quad + \\ 0,161 \quad = \\ \hline \end{array}$$

$$80,151 \times 10^0 = 80,151 \rightarrow 8,0151 \times 10^1 \text{ in forma normalizzata}$$

Altre possibilità sono:

$$\begin{array}{r} 799,9 \quad + \\ 1,61 \quad = \\ \hline \end{array}$$

$$801,51 \times 10^{-1} = 8,0151 \times 10^1 \text{ in forma normalizzata}$$

$$\begin{array}{r} 7,999 \quad + \\ 0,0161 \quad = \\ \hline \end{array}$$

$$8,0151 \times 10^1$$



## Quale forma conviene utilizzare?



Supponiamo di avere 4 cifre in tutto: 1 per la parte intera e 3 per la parte decimale

$$\begin{array}{r} 79,99 \quad + \\ 0,161 \quad = \\ \hline \end{array}$$

$$80,151 \times 10^0$$

$$\begin{array}{r} 799,9 \quad + \\ 1,61 \quad = \\ \hline \end{array}$$

$$801,51 \times 10^{-1}$$

$$\begin{array}{r} 7,999 \quad + \\ 0,0161 \quad = \\ \hline \end{array}$$

$$8,0151 \times 10^1$$

La rappresentazione migliore è:

$$\begin{array}{r} 7,999 \quad + \\ 0,0161 \quad = \\ \hline \end{array}$$

$$8,0151 \times 10^1$$

Con la quale posso scrivere: 1 cifra prima della virgola (8) e 3 cifre dopo la virgola (015), 1 va perso.

Con la rappresentazione più a sinistra, perdo le decine, con quella in mezzo decine e centinaia commettendo un errore grande sulla rappresentazione.

Allineo al numero con esponente maggiore.



## 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 (al più vicino) fornisce per rintrare in 3 cifre decimali fornisce: **2,37**



## L'effetto perverso del troncamento



$C = A + B$   
if ( $C \leq A$ ) then  
    (a)...  
else  
    (b)....

$$A = 7,999 \times 10^1 \quad B = 1,61 \times 10^{-1} \quad C = A + B = 8,0151$$

Passando alla codifica su 4 bit ottengo:

$$A = 7,999 \times 10^1 \quad B = 1,61 \times 10^{-1} \quad C = A + B = 8,015 \quad \Rightarrow \quad C > A \text{ correttamente}$$

$$A = 7,999 \times 10^1 \quad B = 1,61 \times 10^{-4} \quad C = A + B = 7,999161$$

Passando alla codifica su 4 bit ottengo:

$$A = 7,999 \times 10^1 \quad B = 1,61 \times 10^{-4} \quad C = A + B = 7,999 \quad \Rightarrow \quad C = A \text{ errore!!!}$$

Questo è un errore molto comune quando si considera l'aritmetica con i numeri decimali



## 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.

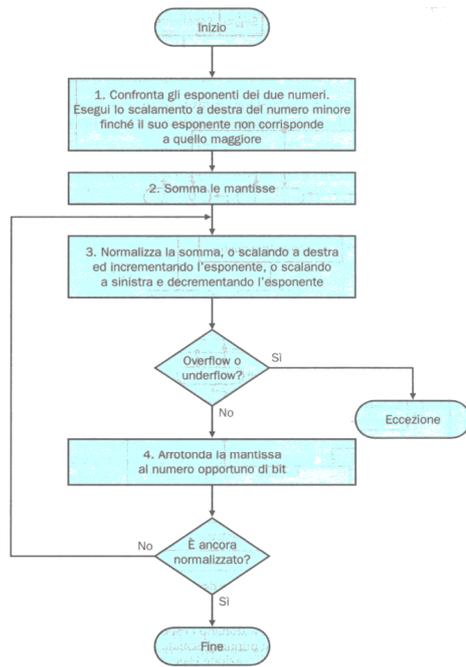
**Se il numero risultante è normalizzato termino qui. Altrimenti:**

3) Normalizzare il risultato.





# Algoritmo risultante



## Il circuito della somma floating point: gli attori



A + B

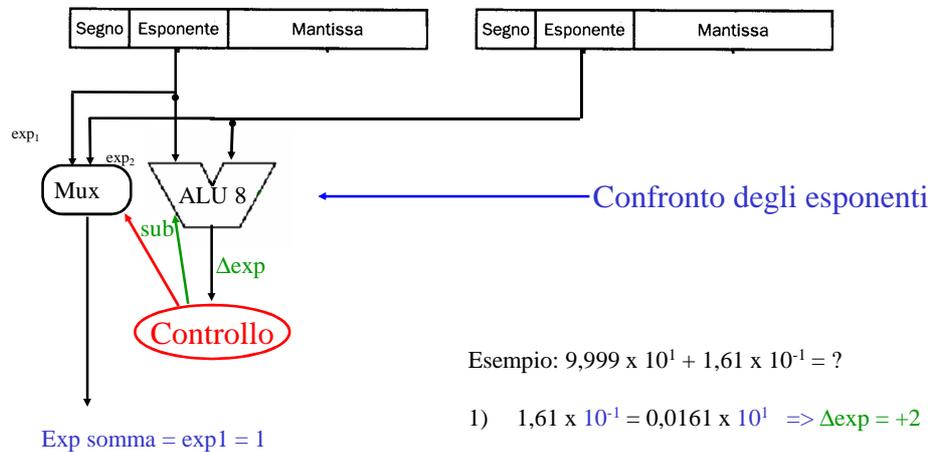
C



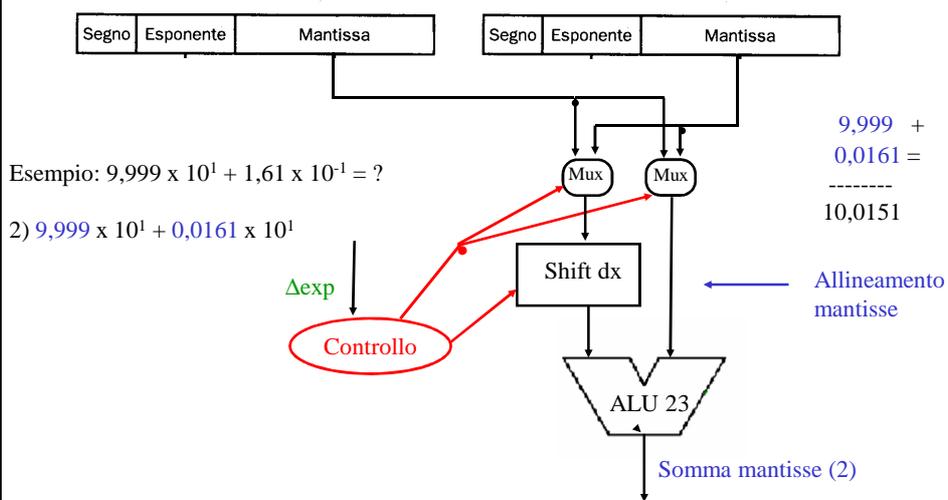
Rappresentazione normalizzata IEEE754

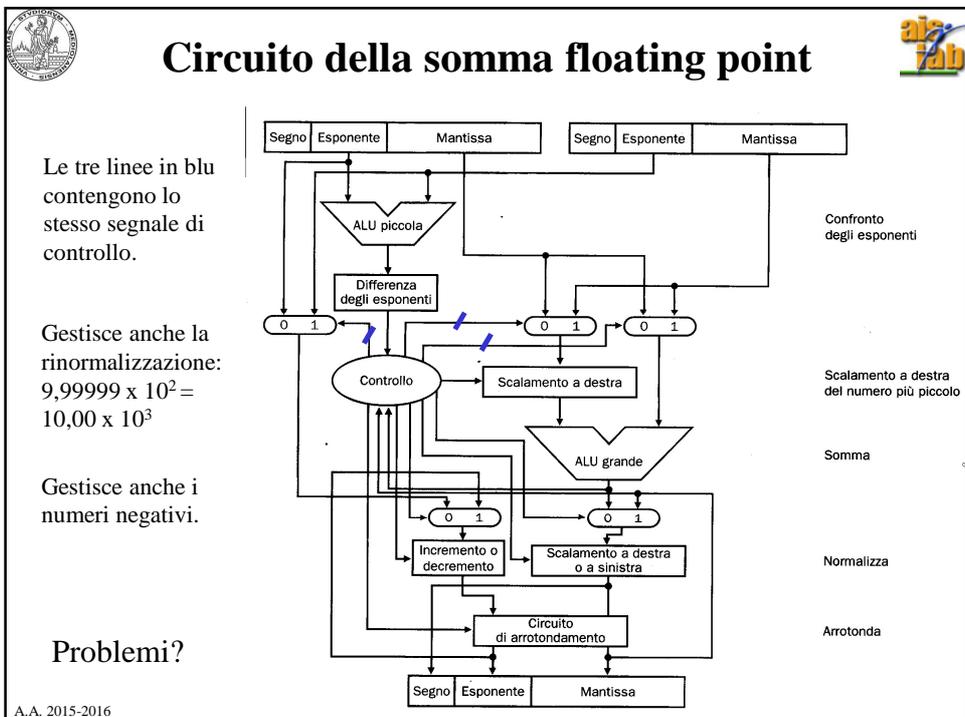
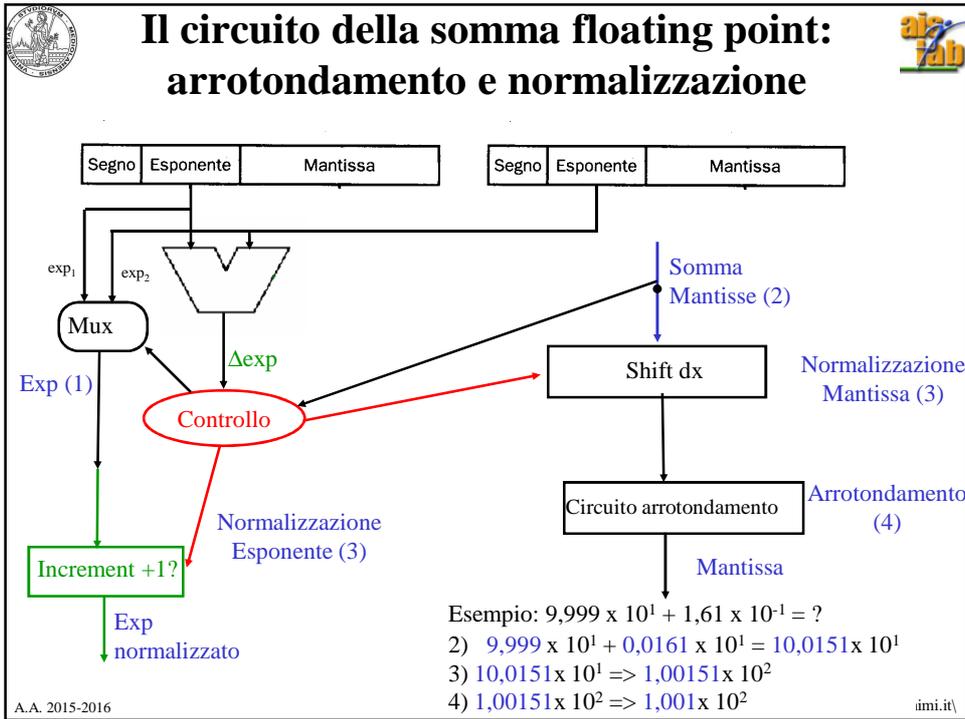


## Il circuito della somma floating point: determinazione dell'esponente comune



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







# Sottrazione



$P = 0,5 - 0,4375$  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) Sottrazione 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.



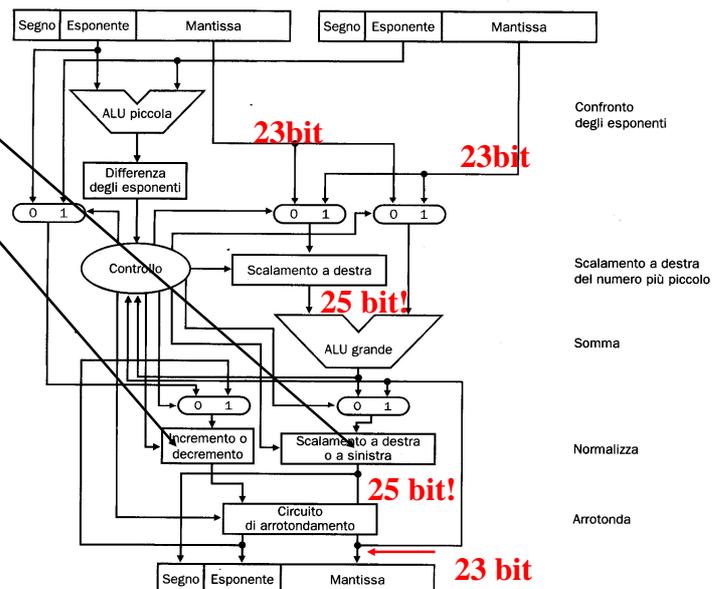
# Circuito della somma floating point con bit di arrotondamento



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