

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

Unità didattica: Errori di roundoff [3-AT]

Titolo: Accuratezza statica e dinamica del S. A. Standard
IEEE 754

Argomenti trattati:

- ✓ Misure di errore (errore assoluto ed errore relativo) e loro significato
- ✓ Errore di roundoff di rappresentazione (accuratezza statica)
- ✓ Errore di roundoff delle operazioni aritmetiche floating-point
(accuratezza dinamica)
- ✓ Ottimalità del S.A. Standard: massima accuratezza statica e
massima accuratezza dinamica

Prerequisiti richiesti: Rappresentazione dei numeri floating-point

PRECISIONE FINITA DEL SISTEMA ARITMETICO FLOATING-POINT



Errore di roundoff

statico

dinamico

raccomandato
rappresentazione
in memoria

risultato di opera-
zioni aritmetiche

Misure di errore

Se x indica il valore esatto ed \tilde{x} una sua approssimazione ($\tilde{x}=fl(x)$), per misurare l'accuratezza di \tilde{x} rispetto ad x si usano:

Errore assoluto $E_A(\tilde{x}) = |x - \tilde{x}|$

Errore relativo $E_R(\tilde{x}) = \left| \frac{x - \tilde{x}}{x} \right|$ per $x \neq 0$

Quali informazioni danno gli errori E_A ed E_R sul grado di approssimazione?

Errore assoluto E_A



***cifre
decimali
corrette***

Errore relativo E_R



***cifre
significative
corrette***

... cosa significa ???

$$x = 12.34567 = 1.234567 e^{+02}$$

cifre intere

cifre decimali
o frazionarie

cifre significative

ordine di grandezza

Proprietà 1

Se \tilde{x} è un'approssimazione di x corretta a p cifre decimali, allora si ha $|x - \tilde{x}| < 10^{-p}$

Proprietà 2

Se \tilde{x} è un'approssimazione di x corretta a p cifre significative, allora si ha $\frac{|x - \tilde{x}|}{|x|} < 10^{-p+1}$

Esempio: proprietà 1

$$x = 12.34999$$

$$y = \tilde{x} = 12.34$$

$p=2$ cifre decimali
corrette

Errore assoluto in $y = E_A(y) = 9.9900e-003 < 10^{-2} = 10^{-p}$

$$x = 12.34999$$

$$y = \tilde{x} = 12.35$$

$p=1$ cifra decimale
corretta

?

$p=4$



Errore assoluto in $y = E_A(y) = 1.0000e-005 < 10^{-4} = 10^{-p}$???

Esempio: proprietà 2

$$x = 12.34999$$

$$y = \tilde{x} = 12.34 \quad p=4 \text{ cifre significative corrette}$$

$$\text{Errore relativo} = E_R(y) = 8.0891e-004 < 10^{-3} = 10^{-p+1}$$

$$x = 12.34999$$

$$y = \tilde{x} = 12.35 \quad p=3 \text{ cifre significative corrette}$$

? ← $p=5$ ↑

$$\text{Errore relativo} = E_R(y) = 8.0972e-007 < 10^{-6} = 10^{-p+1} \quad ???$$

Anche se le due proprietà precedenti valgono per una sola implicazione (\Rightarrow), nella pratica si suppone l'equivalenza (\Leftrightarrow), nel senso che dall'ordine di grandezza degli errori si hanno informazioni sulle cifre (decimali o significative) corrette di un'approssimazione rispetto al corrispondente valore esatto.

Esempio: il seguente programma

```
#include <stdio.h>
#include <math.h> // per pow() e atan() [arcotangente]

void main()
{double double_x, rel_err, ass_err; float float_x; char i;
  for (i=-4; i<5; i=i+2)
    {double_x=4*atan(1.0)*pow(10,i);
     float_x=(float) double_x;

     ass_err=fabs(double_x-float_x);
     rel_err=ass_err/fabs(double_x);

     printf("double = %22.16e\n",double_x);
     printf("single = %22.16e\n",float_x);
     printf("errore assoluto = %8.3e\n",ass_err);
     printf("errore relativo = %8.3e\n",rel_err);
     puts("\n\n");
    }
}
```

calcola gli *errori di rappresentazione* del tipo float ...

double (valore esatto!)	float	double (valore esatto!)
0. <u>00031415926</u> 535...	<u>3.1415926</u> 059...e-004	<u>3.1415926</u> 535...e-004
0. <u>03141592</u> 6535...	<u>3.141592</u> 8155...e-002	<u>3.141592</u> 6535...e-002
3. <u>141592</u> 6535...	<u>3.141592</u> 7410...e-000	<u>3.141592</u> 6535...e-000
314. <u>1592</u> 6535...	<u>3.141592</u> 7124...e+002	<u>3.141592</u> 6535...e+002
31415. <u>92</u> 6535...	<u>3.141592</u> 5781...e+004	<u>3.141592</u> 6535...e+004

Errore ass.	Errore rel.
4.76... <u>e-012</u>	1.51... <u>e-008</u>
1.61... <u>e-009</u>	5.15... <u>e-008</u>
8.74... <u>e-008</u>	2.78... <u>e-008</u>
5.88... <u>e-006</u>	1.87... <u>e-008</u>
7.54... <u>e-004</u>	2.40... <u>e-008</u>

cifre decimali corrette

cifre significative corrette

Errore di rappresentazione (tipo `float`)

Errore ass.	Errore rel.
4.76...e-012	1.51...e-008
1.61...e-009	5.15...e-008
8.74...e-008	2.78...e-008
5.88...e-006	1.87...e-008
7.54...e-004	2.40...e-008

Perché l'errore relativo ha sempre lo stesso ordine di grandezza?

Perché dipende dalle cifre significative della mantissa (nel tipo `float` del C la precisione è di 24 bit)!

Perché l'ordine di grandezza è 10^{-8} ?

precisione = numero di cifre_p significative rappresentate

rappresentazione (binaria) fl.p.
a bit implicito su **t** bit per la
mantissa

precisione (binaria) = **t+1**

Nell'Aritmetica Standard l'epsilon macchina $\epsilon_{\text{mach}} = 2^{-t}$

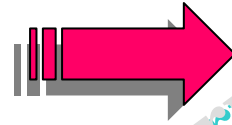
è tale che $\frac{|x - fl(x)|}{|x|} \leq \frac{\epsilon_{\text{mach}}}{2}$

Qual'è la **precisione binaria**?

$$-\log_2 \left(\frac{|x - fl(x)|}{|x|} \right) \geq -\log_2 \left(\frac{\epsilon_{\text{mach}}}{2} \right) = t + 1$$

precisione = numero di cifre _{β} significative rappresentate

rappresentazione (binaria) fl.p.
a bit implicito su t bit per la
mantissa



precisione (binaria) = $t+1$

Qual'è la **precisione decimale equivalente**?

$$-\log_{10} \left(\frac{|x - fl(x)|}{|x|} \right) \geq -\log_{10} \left(\frac{\epsilon_{\text{mach}}}{2} \right) \quad (\epsilon_{\text{mach}} = 2^{-t})$$

$\beta=2$, bit implicito		precisione decimale equivalente	
float	$t=23$ bit	s.p.	$-\log_{10}(\dots) \approx 6.9 \approx 7 \div 8$ cifre decimali
double	$t=52$ bit	d.p.	$-\log_{10}(\dots) \approx 15.9 \approx 16 \div 17$ cifre decimali

perciò l'errore relativo nel tipo float è sempre dell'ordine di 10^{-8}

[illegible]

formato	s.p.	d.p.
%f	0.666667	0.666667
%e	6.666667e-001	6.666667e-001
%8.3f	0.667	0.667
%8.3e	6.667e-001	6.667e-001
%21.15f	0.666666686534882	0.6666666666666667
%21.15e	6.666666865348816e-001	6.6666666666666666e-001

$$\frac{|x - fl(x)|}{|x|}$$

accuratezza statica

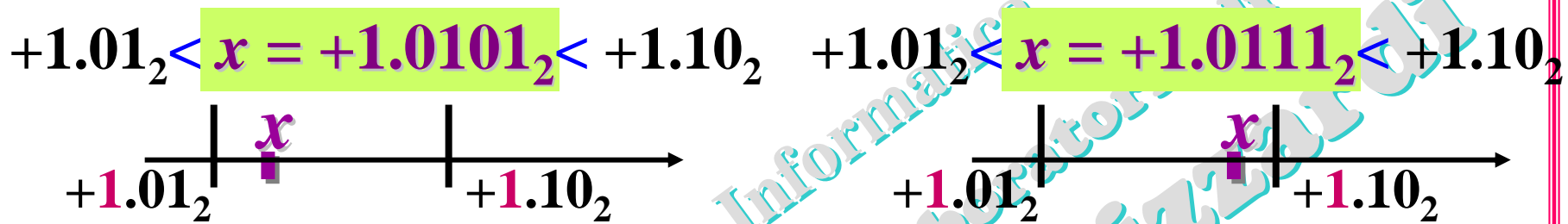
Errore di roundoff di rappresentazione

È l'errore introdotto nel passare da un numero reale x a precisione infinita al suo rappresentante floating-point in memoria $fl(x)$

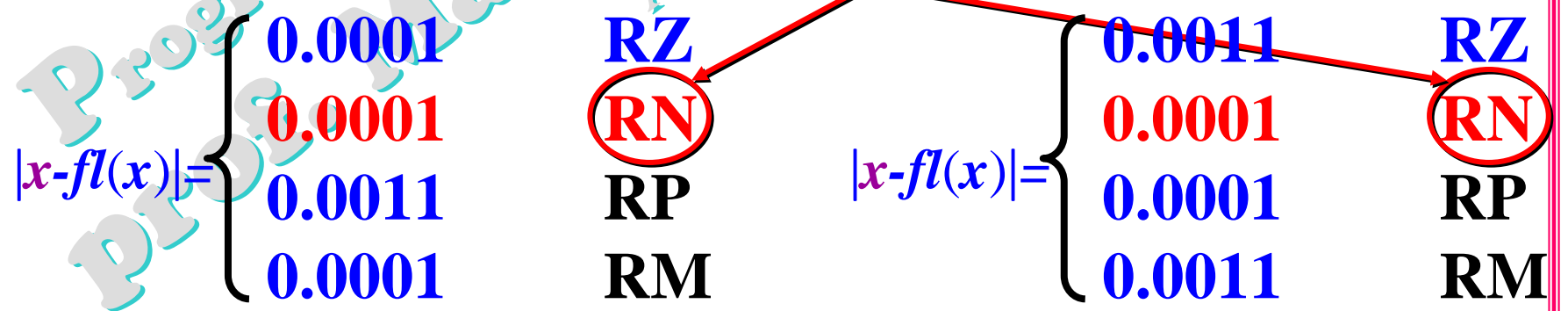
Dipende dalla precisione del Sistema Aritmetico Floating-point e dallo schema di *rounding* utilizzato

Esempio: schemi IEEE St. 754

$x \notin F(\beta=2, t=2 \text{ (bit implicit)}, E_{\min}=0, E_{\max}=3)$



L'errore è sempre il minimo per **RN**



Proprietà
matematica: $\forall a \in \mathbb{R}, \forall \varepsilon > 0 \Rightarrow a + \varepsilon > a$

Nel Sistema Aritmetico Floating-point non vale!

$$\forall a \in \mathbf{F}(\beta, t, E_{\min}, E_{\max}),$$

$$\exists \varepsilon > 0 : a \oplus \varepsilon = a$$

↑ addizione floating-point

Fissato un numero floating-point a : qual è il più piccolo numero floating-point che dà contributo nell'addizione floating-point con a , cioè

$$\text{ulp}(a) = \min\{0 < \varepsilon \in \mathbf{F}(\beta, t, E_{\min}, E_{\max}) : a \oplus \varepsilon > a\}$$

u.l.p. = Unit in the Last Place

Approssimazione di $ulp(1)$

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

```
void main()
```

```
{ float eps, epsp1; int n;
```

```
  eps=1; n=0;
```

```
  epsp1=eps+1;
```

```
  while (epsp1>1)
```

```
  { eps=eps/2; n++;
```

```
    epsp1=eps+1;
```

```
  }
```

```
  eps=2.0f*eps; n--;
```

```
  printf("num. divisioni per 2=%d\nepsilon=%e \n",n,eps);
```

```
}
```

genera successione
 $\{2^{-n}\}$

calcola $\{1+2^{-n}\}$

num. divisioni per 2=23
epsilon=1.192093e-007

$u.l.p.(1) = \text{Epsilon macchina } \epsilon_{mach}$

$\epsilon_{mach} = \min\{0 < \epsilon \in F(\beta, t, E_{min}, E_{max}) : 1 \oplus \epsilon > 1\}$

dove \oplus sta per *addizione floating-point*

bit implicito + round to nearest $\Rightarrow \epsilon_{mach} = 2^{-t}$

Esempio: *Epsilon macchina* ε_{mach} in singola precisione

mostra_float_double.c

FLT_EPSILON

34000000

1.19...e-7

0

011 0100 0

000 0000 0000 0000 0000 0000

+1

3f800000

0

011 1111 1

000 0000 0000 0000 0000 0000

1+FLT_EPS

3f800001

0

011 1111 1

000 0000 0000 0000 0000 0001

1+.6FLT_EPS

3f800001

0

011 1111 1

000 0000 0000 0000 0000 0001

1+.5FLT_EPS

3f800000

0

011 1111 1

000 0000 0000 0000 0000 0000

ESEMPIO 7b: *Epsilon macchina* ϵ_{mach} in doppia precisione

DBL_EPSILON	3cb00000 00000000		
2.22...e-16	0	011 1100 1011	0000 0000 0000... 0000 0000

+1	3ff00000 00000000		
	0	011 1111 1111	0000 0000 0000... 0000 0000

1+DBL_EPS	3ff00000 00000001		
	0	011 1111 1111	0000 0000 0000... 0000 0001

1+DBL_EPS/2	3ff00000 00000000		
	0	011 1111 1111	0000 0000 0000... 0000 0000

Se si elimina la variabile di appoggio `epsp1` dal programma ...

```
#include <stdio.h>
void main()
{ float eps; int n;
  eps=1; n=0;
  while (eps+1>1)
    eps=eps/2; n++;
  eps=2.0f*eps; n--;
  printf("num. divisioni per 2=%d\nepsilon=%e \n",n,eps);
}
```

epsilon macchina



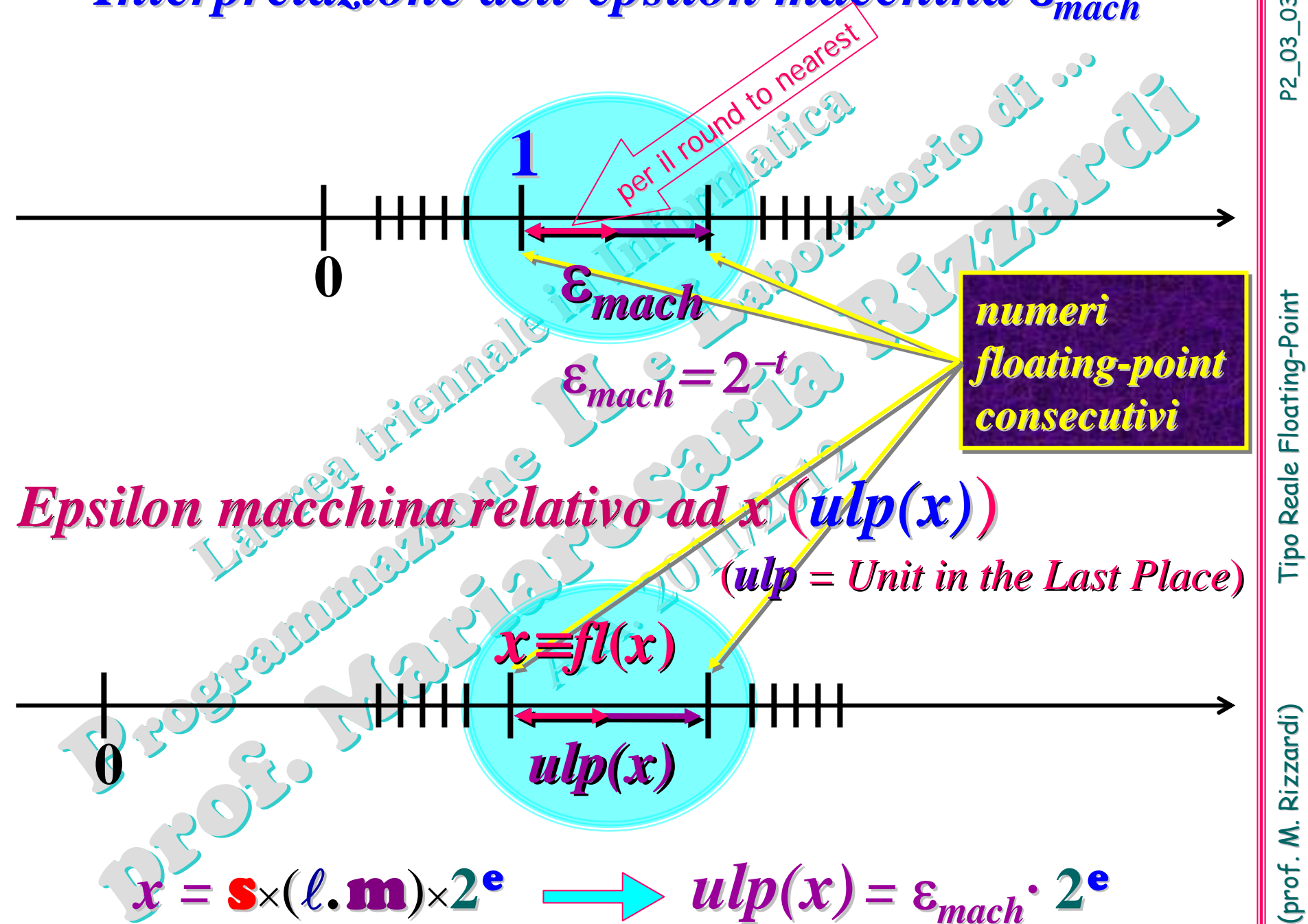
```
num. divisioni per 2=63
epsilon=1.084202e-019
```



Senza la variabile di appoggio `epsp1`, il ciclo `while` viene eseguito nei **registri** dell'unità aritmetica che hanno un campo mantissa maggiore di quello della memoria:

precisione_registro = **precisione_formato IEEE std Extended**
(cioè del tipo `long double` del C).

Interpretazione dell'epsilon macchina ϵ_{mach}



Nel **S.A. Std.**, se x è un numero reale rappresentabile e $fl(x)$ indica il corrispondente numero floating point:

$fl(x)$ risultato di
massima accuratezza



$$fl(x) = x(1 + \delta),$$

$$|\delta| \leq \frac{1}{2} \epsilon_{mach}$$

dove $\delta = \frac{|x - fl(x)|}{|x|}$, δ **errore relativo**

Mediante il bit implicito e lo schema del round to nearest, il **S.A. IEEE Std 754** assicura la massima accuratezza nella rappresentazione in memoria dei numeri reali (**massima accuratezza statica**).

Come assicurare nelle operazioni aritmetiche risultati di massima accuratezza (**massima accuratezza dinamica**) ?

accuratezza dinamica

Si dimostra che nei registri dell'unità aritmetica per assicurare risultati di massima accuratezza (*massima accuratezza dinamica*) sono sufficienti:

- ❑ 1 bit per rappresentare il bit implicito (ℓ)
- ❑ 2 guard-bit (g) oltre la mantissa
- ❑ 1 sticky-bit *

* sticky bit = bit che vale 1 se per esso transita almeno un bit=1



In tal modo con *registri lunghi* ($t+3+ 1$ *sticky*) *bit* si ottengono gli stessi risultati aritmetici che si otterrebbero con registri di lunghezza ...“infinita” (economia di costi).

Esempio 1: uso dei guard bit

cancellazione in s.p. $x-y$

$x=1$ 3f800000		
0	011 1111 1	000 0000 0000 0000 0000 0000
$y=1-2^{-24}$ 3f7fffff		
0	011 1111 0	111 1111 1111 1111 1111 1111

nei registri, dopo l'allineamento degli operandi ...

guard bit
sticky bit

x	0	011 1111 1	1	000 0000 0000 0000 0000 0000	0	0	0
y	0	011 1111 1	0	111 1111 1111 1111 1111 1111	1	0	0
$x-y$	0	011 1111 1	0	000 0000 0000 0000 0000 0000	1	0	0
$x-y$	0	011 0011 1	1	000 0000 0000 0000 0000 0000	0	0	0

normalizzazione

Esempio 1: uso dei guard bit

```
1  // Esempio 1: cancellazione in s.p. x y
2  #include <stdio.h>
3  #include <math.h>
4  #define MAX_LEN 32
5
6  void mostra_32_bit(int num, short bit[32])
7  {
8      char k;
9      for (k=0; k<32; k++) bit[k]=0;
10     for (k=0; k<32; k++)
11     {
12         bit[k] = num & 1;
13         num = num >> 1;
14     }
15     for (k=31; k>=0; k--)
16         (k==31 | k==23) ? printf("%1d ", bit[k]) : printf("%1d", bit[k]);
17     puts("\n");
18 }
19
20 int main()
21 {
22     union sp
23     {
24         float f;
25         int n;
26     } x, y, z;
27     short bit[MAX_LEN];
28
29     x.f=1.0f;
30     y.f=1.0f-(float)pow(2.0, -24);
31     z.f=x.f-y.f;
32
33     printf("\nx = %e\thex = %08x\n\n", x.f, x.n); mostra_32_bit(x.n, bit);
34     printf("\ny = %e\thex = %08x\n\n", y.f, y.n); mostra_32_bit(y.n, bit);
35     printf("\nz = %e\thex = %08x\n\n", z.f, z.n); mostra_32_bit(z.n, bit);
36
37     return 0;
38 }
```

x = 1.000000 hex = 3f800000
0 01111111 000000000000000000000000

y = 9.999999e-001 hex = 3f7fffff
0 01111110 111111111111111111111111

z = 5.960464e-008 hex = 33800000
0 01100111 000000000000000000000000

attenzione: cosa cambia se si usa il formato %f?

Esempio 2: uso dello sticky bit

cancellazione in s.p. $x-y$

$x=1$	3f800000
0	011 1111 1 000 0000 0000 0000 0000 0000
$y=(1+2^{-23})2^{-25}$	33000001
0	011 0011 0 000 0000 0000 0000 0000 0001

nei registri, dopo l'allineamento degli operandi ...

x	0	011 1111 1	1	000 0000 0000 0000 0000 0000	0	0	0
y	0	011 1111 1	0	000 0000 0000 0000 0000 0000	0	1	1
$x-y$	0	011 1111 1	0	111 1111 1111 1111 1111 1111	1	0	1
$x-y$	0	011 1111 0	1	111 1111 1111 1111 1111 1111	0	1	0

← = *normalizzazione*

Senza lo sticky-bit, nemmeno un registro a lunghezza doppia avrebbe fornito un risultato di massima accuratezza!

Esempio 2a: uso dello sticky bit

```
1  /* Esempio 2a: cancellazione in s.p. x-y */
2  #include <stdio.h>
3  #include <math.h>
4  #define MAX_LEN 32
5
6  void mostra_32_bit(int num, short bit[32])
7  {
8      char k;
9      for (k=0; k<32; k++) bit[k]=0;
10     for (k=0; k<32; k++)
11     {
12         bit[k] = num & 1;
13         num = num >> 1;
14     }
15     for (k=31; k>=0; k--)
16         (k==31 | k==23) ? printf("%1d ",bit[k]) : printf("%1d",bit[k]);
17     puts("\n");
18 }
19
20 int main()
21 {
22     union sp
23     {
24         float f;
25         int n;
26     } x, y, z;
27     short bit[MAX_LEN];
28
29     x.f=1.0f;
30     y.f=(1.0f+(float)pow(2.0, -23))*(float)pow(2.0, -25);
31     z.f=x.f-y.f;
32
33     printf("\nx = %e\thex = %08x\n\n",x.f,x.n); mostra_32_bit(x.n, bit);
34     printf("\ny = %e\thex = %08x\n\n",y.f,y.n); mostra_32_bit(y.n, bit);
35     printf("\nz = %e\thex = %08x\n\n",z.f,z.n); mostra_32_bit(z.n, bit);
36
37     return 0;
38 }
```

```
x = 1.000000e+000      hex = 3f800000
0 01111111 000000000000000000000000

y = 2.980233e-008      hex = 33000001
0 01100110 000000000000000000000001

z = 9.999999e-001      hex = 3f7fffff
0 01111110 111111111111111111111111
```

versione 2: passo ... passo (simulazione senza sticky bit)

```

2 #include <stdio.h>
3 #include <math.h>
4 #define MAX_LEN 32
5 void mostra_32_bit(int num, short bit[])
6 {
7     char k;
8     for (k=0; k<32; k++) bit[k] = (num & 1); num = num >> 1;
9     for (k=0; k<32; k++)
10     {
11         bit[k] = num & 1; num = num >> 1;
12     }
13     for (k=31; k>=0; k--)
14         (k==31 | k==23) ? printf("%d ", bit[k]) : printf("%d", bit[k]);
15     puts("\n");
16 }
17 int main()
18 {
19     union sp
20     {
21         float f;
22         int n;
23     } x, y, z;
24     short bit[MAX_LEN];
25
26     puts("\nPasso passo ...\n");
27     x.f=1.0f;
28     printf("\nx = %e\thex = %08x\n\n", x.f, x.n); mostra_32_bit(x.n, bit);
29
30     y.f=(1.0f+(float)pow(2.0, -23))*(float)pow(2.0, -25);
31     printf("\ny = %e\thex = %08x\n\n", y.f, y.n); mostra_32_bit(y.n, bit);
32
33     puts("Allinea y a x");
34     y.f=y.f*(float)pow(2.0, -25); // allinea y a x
35     printf("\ny = %e\thex = %08x\n\n", y.f, y.n); mostra_32_bit(y.n, bit);
36
37     z.f=x.f-y.f;
38     printf("\nz = %e\thex = %08x\n\n", z.f, z.n); mostra_32_bit(z.n, bit);
39     return 0;

```

Passo passo ...

x = 1.000000e+000 hex = 3f800000

0 01111111 000000000000000000000000000000

y = 2.980233e-008 hex = 33000001

0 01100110 000000000000000000000000000001

Allinea y a x

y = 8.881785e-016 hex = 26800001

0 01001101 000000000000000000000000000001

z = 1.000000e+000 hex = 3f800000

0 01111111 000000000000000000000000000000

Esempio 2b: uso dello sticky bit

Tipo

(prof. M. Rizzardi)

Esercizi

1

Scrivere delle function C per calcolare rispettivamente l'epsilon macchina della singola, della doppia precisione e della precisione long double, visualizzando ad ogni passo i singoli bit.

esempio: singola precisione

n= 0	eps+1 =	0	10000000	000000000000000000000000
n= 1	eps+1 =	0	01111111	100000000000000000000000
n= 2	eps+1 =	0	01111111	010000000000000000000000
n= 3	eps+1 =	0	01111111	001000000000000000000000
n= 4	eps+1 =	0	01111111	000100000000000000000000
n= 5	eps+1 =	0	01111111	000010000000000000000000
n= 6	eps+1 =	0	01111111	000001000000000000000000
n= 7	eps+1 =	0	01111111	000000100000000000000000
n= 8	eps+1 =	0	01111111	000000010000000000000000
n= 9	eps+1 =	0	01111111	000000001000000000000000
n=10	eps+1 =	0	01111111	000000000100000000000000
n=11	eps+1 =	0	01111111	000000000010000000000000
n=12	eps+1 =	0	01111111	000000000001000000000000
n=13	eps+1 =	0	01111111	000000000000100000000000
n=14	eps+1 =	0	01111111	000000000000010000000000
n=15	eps+1 =	0	01111111	000000000000001000000000
n=16	eps+1 =	0	01111111	000000000000000100000000
n=17	eps+1 =	0	01111111	000000000000000010000000
n=18	eps+1 =	0	01111111	000000000000000001000000
n=19	eps+1 =	0	01111111	000000000000000000100000
n=20	eps+1 =	0	01111111	000000000000000000010000
n=21	eps+1 =	0	01111111	000000000000000000001000
n=22	eps+1 =	0	01111111	000000000000000000000100
n=23	eps+1 =	0	01111111	000000000000000000000010

2

Scrivere una function C per calcolare dalla definizione l'ULP(x) dove x è il parametro reale float di input. [liv. 3]

Generando in modo random i bit^(*) di un numero reale x (double x), determinare i bit della corrispondente variabile flx (float flx ; $flx=(float) x$). Se il numero x è rappresentabile nel tipo *float* calcolarne l'errore assoluto E_A e relativo E_R (considerando come esatto il double x e come approssimante il float flx) dalle formule

$$E_A(fl x) = |x - fl x| \qquad E_R(fl x) = \frac{|x - fl x|}{|x|} \qquad [\text{liv. 2}]$$

(*) La funzione C **rand()**, in **stdlib.h**, restituisce un intero pseudorandom minore o uguale di **RAND_MAX** (= $32767_{10} = 7fff_{16} = 0111\ 1111\ 1111\ 1111_2$ - massimo intero positivo su 16 bit).

Come generare a caso i 64 bit di una variabile double?

- ☐ Generando i **singoli bit random**:
 - come resto della divisione per 2 [... come?];
 - associando 0 agli interi < 16384 e 1 agli interi > 16383;
- ☐ Generando **4 sequenze random di 16 bit** ($4 \times 16 = 64$) [... come?]

vedere: Uso di **rand()** in Materiale di supporto