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



Tipo Reale Floating-Point

orof. M. Rizzardi)



Errore di roundoff

statico

dinamico

risultato di operazioni aritmetiche

rappresentazione in memoria

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

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

Quali informazioni danno gli errori E_A ed E_R sul grado di approssimazione 3

Errore assoluto E significative corrette

... cosa significa ???

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

Proprietà 2

Se \hat{x} è un'approssimazione di x corretta a p significative, allora si ha

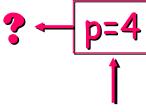
Esempio: proprietà 1

$$x = 12.34999$$
 $y = \tilde{x} = 12.34$ p=2 cifre decimality correcte

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

$$x = 12.34999$$

$$=$$
 \tilde{x} $=$ 12.35 p=1 cifra decimale corretta



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

Esempio: proprietà 2

$$x = 12.34999$$

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

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

$$x = 12.34999$$

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

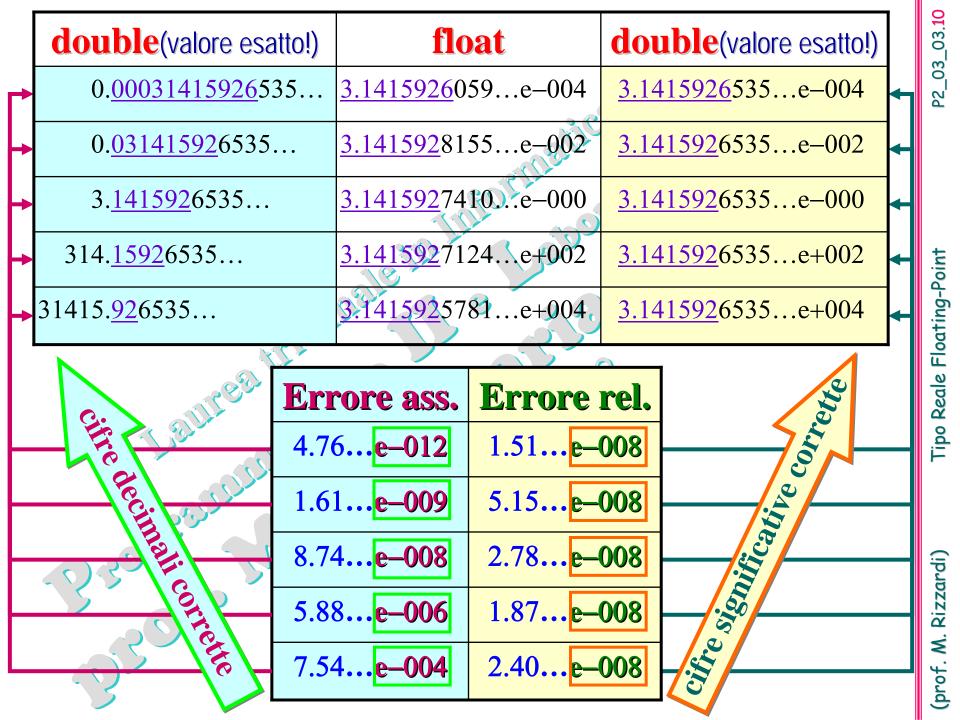
Errore relativo =
$$E_R(y)$$
 = 8.0972e-007 < 10^{-6} = 10^{-p+1} ???

Anche se le due proprietà precedenti valgono per una sola implicazione (=>), nella pratica si suppone l'equivalenza (←⇒), 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)
     {\text{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 ...



Errore di rappresentazione (tipo float)

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

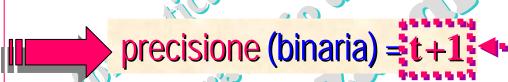
Perché l'errore relativo ha sempre lo stesso ordine di grandezza?

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?

$precisione = numero di cifre_{\beta} significative rappresentate$

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



Nell'Aritmetica Standard l'epsilon macchina

$$\epsilon_{\rm mach} = 2^-$$

è tale che x-fl(x) $\ge \frac{\varepsilon_{\text{mach}}}{2}$

Qual'è la precisione binaria?

$$-\log_2\left(\frac{|x-f(x)|}{|x|}\right) \ge -\log_2\left(\frac{\varepsilon_{\text{mach}}}{2}\right) = t+1$$

$precisione = numero di cifre_{\beta} significative rappresentate$

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



precisione (binaria) = t+1

Qual'è la precisione decimale equivalente?

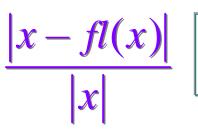
$$-\log_{10}\left(\frac{|x-f(x)|}{|x|}\right) \ge -\log_{10}\left(\frac{\varepsilon_{\text{mach}}}{2}\right) \left(\varepsilon_{\text{mach}} = 2^{-t}\right)$$

β=2, bit implicito precisione decimale equivalente float t=23 bit s.p. $-\log_{10}(...)\approx6.9\approx7\div8$ cifre decimali double t=52 bit d.p. $-\log_{10}(...)\approx15.9\approx16\div17$ cifre decimali

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

... a che serve? È inutile visualizzare più cifre della precisione decimale equivalente

iormato-	s.p.	a.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	0.6666666666667
%21.15e	6.666666 865348816	6.666666666666
%ZI.13E	e-001	e-001



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

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

Nel Sistema Aritmetico Floating-point non vale!

$$\forall a \in \mathbf{F}(\beta, t, E_{\min}, E_{\max}),$$
 $\exists \epsilon > 0 : a \oplus \epsilon = a$
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

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

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

```
#include <stdio.h>
                            Approssimazione di ulp(1)
void main()
  float eps, epsp1; int n;
    eps=1; n=0;
                                        genera successione
    epsp1=eps+1;
                                               {2<sup>-n</sup>}
    while (epsp1>1)
    { eps=eps/2; n++;
                             calcola \{1+2^{-n}\}
      epsp1=eps+1;
    eps=2.0f*eps; n--;
printf("num. divisioni per 2=%d\nepsilon=%e \n",n,eps);
                          num. divisioni per 2=23
                          epsilon=1.192093e-007
```

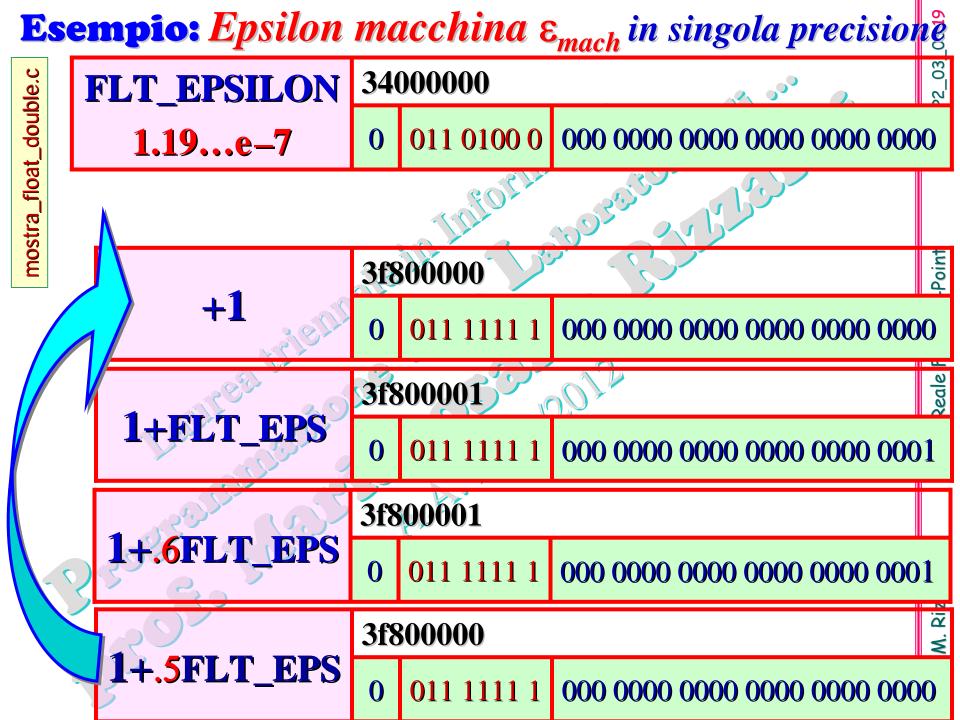
u.l.p.(1) = Epsilon macchina ε_{mach}

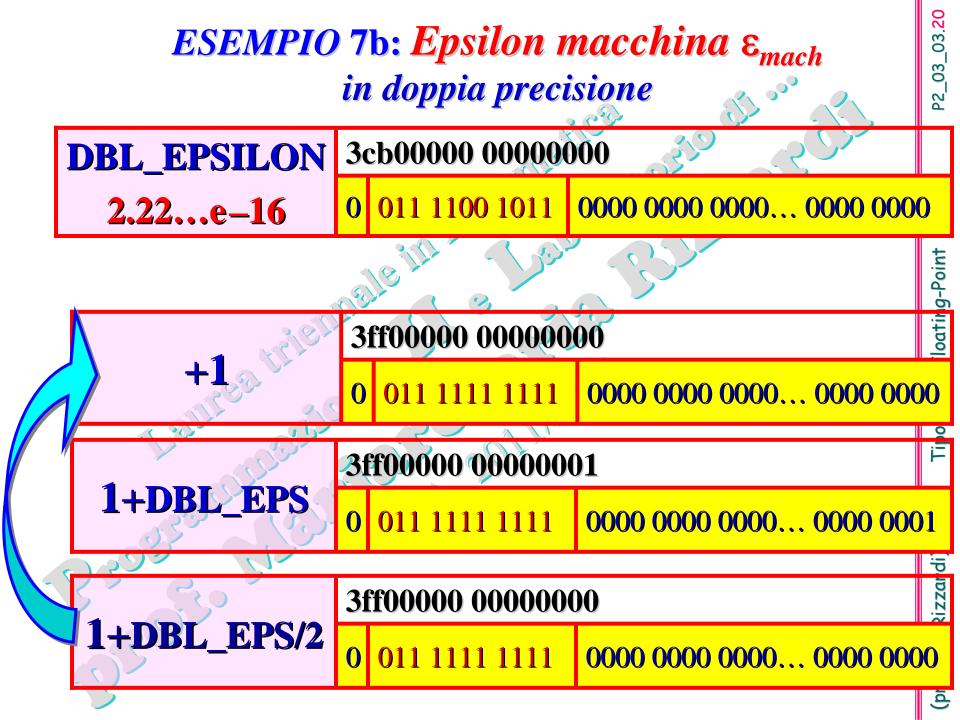
$$\mathbf{\varepsilon}_{mach} = min\{0 < \mathbf{\varepsilon} \in \mathbf{F}(\beta, t, \mathbf{E}_{min}, \mathbf{E}_{max}) : 1 \oplus \mathbf{\varepsilon} > 1\}$$

dove

sta per addizione floating-point







```
#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);
}
```

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



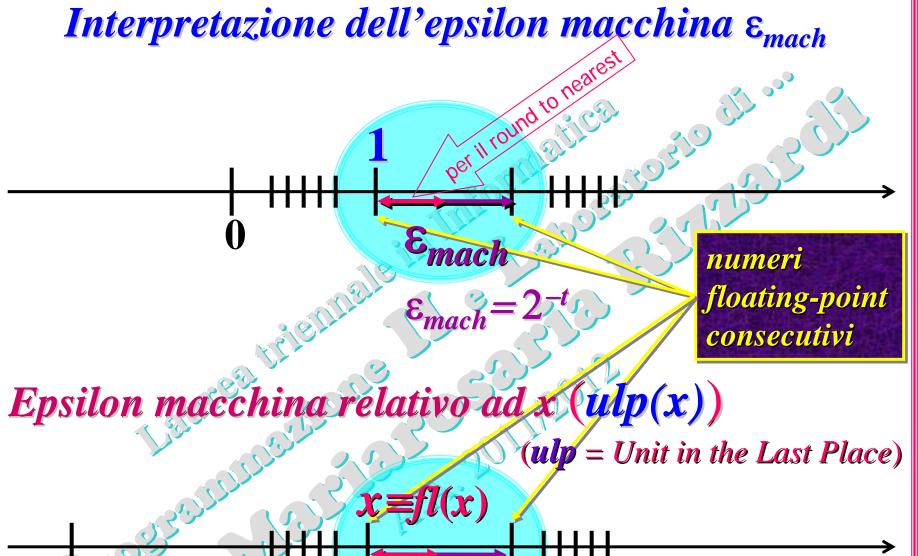
num. divisioni per 2=63 epsilon=1.084202e-019



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

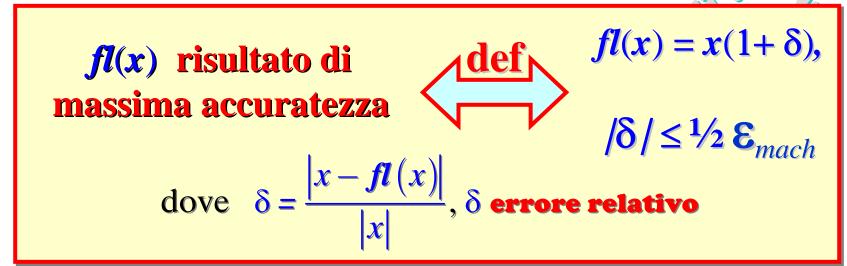
(cioè del tipo long double del C).



 $x = s \times (\ell.m) \times 2^e \longrightarrow ulp(x) = \varepsilon_{mach} \cdot 2^e$

ulp(x)

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



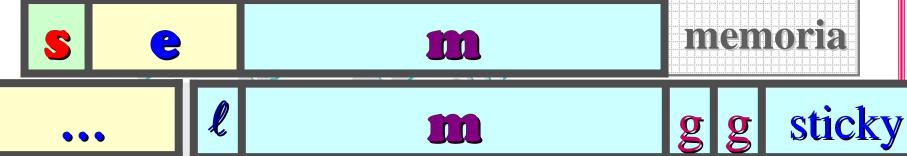
Mediante il bit implicito e lo schema del round to nearest, il *S.A. IEEE Std 754* assicura la massima accuratezza nella rappresentazione in memoria dei numeri reali (*massima accuratezza statica*).

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

accuratezza dinamica

Si dimosta 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



registri ALU

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





3f7fffff

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	111 1111 1111 1111 1111	1	0	0
x-y	0	011, 1111 1	0	000 0000 0000 0000 0000	1	0	0
<i>x</i> - <i>y</i>	0	011 0011 1	1	000 0000 0000 0000 0000 0000	0	0	0

normalizzazione

```
/ psembro r. cancerrazione in s.b. v /
                                                                                     P2_03_03.<mark>26</mark>
      #include <stdio.h>
                                 Esempio 1: uso dei guard bit
      #include <math.h>
      #define MAX LEN 32
 5
      void mostra 32 bit(int num, short bit[32])
          char k;
          for (k=0; k<32; k++) bit[k]=0;
          for (k=0; k<32; k++)
10
              bit[k] = num & 1;
               num = num >> 1;
11
12
                                                                                    g-Point
13
          for (k=31; k>=0; k--)
14
              (k==31 | k==23) ? printf("%1d ",bit[k]) : printf("%1d",bit[k]);
15
          puts("\n");
16
                                         x = 1.000000
                                                                   hex = 3f800000
17
                                           18
      int main()
19
          union sp
              float f;
20
                                           = 9.999999e-001
                                                                   hex = 3f7fffff
21
              int
                  n;
                                            01111110 1111111111111111111111111
22
          } x, y, z;
23
          short bit[MAX LEN];
24
                                         z = 5.960464e-008
                                                                   hex = 33800000
25
          x. f=1.0 f;
                                         26
          y. f=1.0 f-(float) pow(2.0, -24);
27
          z.f=x.f-y.f;
                                                                                     Rizzardi
28
          printf("\nx = %e\thex = %08x\n\n", x. f, x. n); mostra 32_bit(x. n, bit);
29
30
          printf("\ny = \ensuremath{\$e}\thex = \ensuremath{\$08x\n\n",y.f,y.n}); mostra 32 bit(y.n, bit);
                                                                                     (prof. M.
          printf("\nz = \ensuremath{\$\epsilon}\thex = \ensuremath{\$08x\n\n",z.f,z.n}); mostra 32 bit(z.n, bit);
31
32
33
      return 0;
                           attenzione: cosa cambia se si usa il formato £?
34
```

guard

Esempio 2: uso dello sticky bit

cancellazione in s.p. x-y



nei registri, dopo l'allineamento degli operandi .

inci registri, dopo ramineamento degil operanar				+			
x	0	011 1111 1		000 0000 0000 0000 0000	0	0	0
y	0	011 1111 1	0	000 0000 0000 0000 0000 0000	0	1	1
<i>x</i> - <i>y</i>	0	011 1111 1	0	111 1111 1111 1111 1111 1111 1111 1111 1111	1	0	1
x-y	0	011 1111 0	1	111 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: cancellazione in s.p. x-y */
                                                                       P2_03_03.28
#include <stdio.h>
                         Esempio 2a: uso dello sticky bit
#include <math.h>
#define MAX LEN 32
void mostra 32 bit(int num, short bit[32])
   char k;
   for (k=0; k<32; k++) bit[k]=0;</pre>
   for (k=0; k<32; k++)
        bit[k] = num & 1;
        num = num >> 1;
                                                                      ng-Point
   for (k=31; k>=0; k--)
       (k==31 | k==23) ? printf("%1d ",bit[k]) : printf("%1d",bit[k]);
   puts("\n");
                          x = 1.000000e+000
                                                       hex = 3f800000
                          int main()
   union sp
                                                      hex = 33000001
                          y = 2.980233e-008
       float f;
                          int
            n;
   } x, y, z;
                          z = 9.999999e-001
                                                       hex = 3f7fffff
   short bit[MAX LEN];
                            x. f=1.0 f;
   y. f = (1.0 f + (float) pow(2.0, -23)) * (float) pow(2.0, -25);
                                                                      (prof. M. Rizzardi)
   z.f=x.f-y.f;
   printf("\nx = %e\thex = %08x\n\n", x.f, x.n); mostra 32 bit(x.n, bit);
   printf("\ny = %e\thex = %08x\n\n", y, f, y, n); mostra 32 bit(y, n, bit);
   printf("\nz = %e\thex = %08x\n\n", z.f, z.n); mostra 32 bit(z.n, bit);
return 0;
```

1

5 6

8 9

10

11

12

13

14 15

16 17

18

19

20

21

22

23

24

25

26

27

28 29

30

31

32 33

```
#include <stdio.h>
     #include <math.h> versione 2: passo ... passo (simulazione senza sticky bit)
     #define MAX LEN 32
                              Passo passo ...
     void mostra 32 bit(int num, s
                              x = 1.0000000e+000 hex = 3f800000
        char k;
        for (k=0; k<32; k++) bit[]
                              8
        for (k=0; k<32; k++)
            bit[k] = num & 1; num
                              y = 2.980233e-008 hex = 33000001
10
11
        for (k=31; k>=0; k--)
                              (k==31 \mid k==23) ? pri:
12
        puts("\n");
13
                              Allinea y a x
14
15
     int main()
                              y = 8.881785e-016 hex = 26800001
16
        union sp
                              float f;
            int
                n;
        } x, y, z;
                              z = 1.000000e + 000 hex = 3f800000
        short bit[MAX LEN];
                              puts("\nPasso passo ...\n");
        x. f=1.0 f;
        printf("\nx = %e\thex = %08x\n\n", x. f, x. n); mostra 32 bit(x. n, bit);
268
        y. f = (1.0 f + (float) pow(2.0, -23)) * (float) pow(2.0, -25);
        printf("\ny = %e\thex = %08x\n\n", y. f, y. n); mostra 32 bit(y. n, bit);
        puts("Allinea y a x");
        y. f=y.f*(float)pow(2.0, -25); // allinea y a x
        printf("\ny = %e\thex = %08x\n\n", y.f, y.n); mostra 32 bit(y.n, bit);
         z.f=x.f-v.f;
        printf("\nz = %e\thex = %08x\n\n", z.f, z.n); mostra 32 bit(z.n, bit)
```

Rizzardi)

(prof. M.

Gli esempi sull'uso dei guard bit e dello sticky bit continuano a valere anche quando si usano variabili long double

Esempio 2a

Esempio 2b

Passo passo ...

Allinea y a x

Esercizi

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.

eps+1 =10000000 n=0eps+1 =01111111 n=1eps+1 =01111111 n=2n=3eps+1 =01111111 0001000000000000000000 eps+1 =n=401111111 eps+1 =01111111 00001000000000000000000 n=5eps+1 =00000100000000000000000 n=601111111 eps+1 =01111111 n=7n=8eps+1 =01111111 0000001000000000000000 n= 9 n=10 n=11 n=12 n=13 n=14 n=15 0000000100000000000000 n=9eps+1 =01111111 000000001000000000000 eps+1 =01111111 eps+1 =01111111 000000000100000000000 eps+1 =01111111 000000000010000000000 eps+1 =01111111 000000000001000000000 eps+1 =01111111 000000000000100000000 n=15 eps+1 =01111111 0000000000000100000000 n=16 01111111 0000000000000010000000 eps+1 =n=17eps+1 =01111111 0000000000000001000000 n=18eps+1 =01111111 0000000000000000100000 n=19eps+1 =01111111 0000000000000000010000 01111111 n=20eps+1 =0000000000000000001000 eps+1 =01111111 0000000000000000000100 n=21n=22eps+1 =01111111 n=23eps+1 =01111111 00000000000000000000000000000001

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

Tipo Reale

[liv. 2]

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

$$\boldsymbol{E}_{\boldsymbol{A}}(\mathrm{flx}) = |\mathrm{x} - \mathrm{flx}|$$
 $\boldsymbol{E}_{\boldsymbol{R}}(\mathrm{flx}) = \frac{|\mathrm{x} - \mathrm{flx}|}{|\mathrm{x}|}$

(*) La funzione C rand(), in stdlib.h, restituisce un intero pseudorandom minore o uguale di RAND_MAX (= 32767₁₀ = 7fff₁₆ = 0111 1111 1111 1111₂ - 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 × 16 = 64) [... come?]

vedere: Uso di rand () in Materiale di supporto