

## Limiti delle rappresentazioni in virgola fissa

• I numeri frazionari irrazionali possono essere usati per rappresentare numeri reali (un reale può sempre essere espresso come valore a cui converge una successione di razionali).

• Ci sono però dei limiti di cui occorre tener conto

- LEGAME range, numero di bit, errore, risoluzione



> l'errore di rappresentazione è legato alla più piccola grandezza rappresentabile LSB -

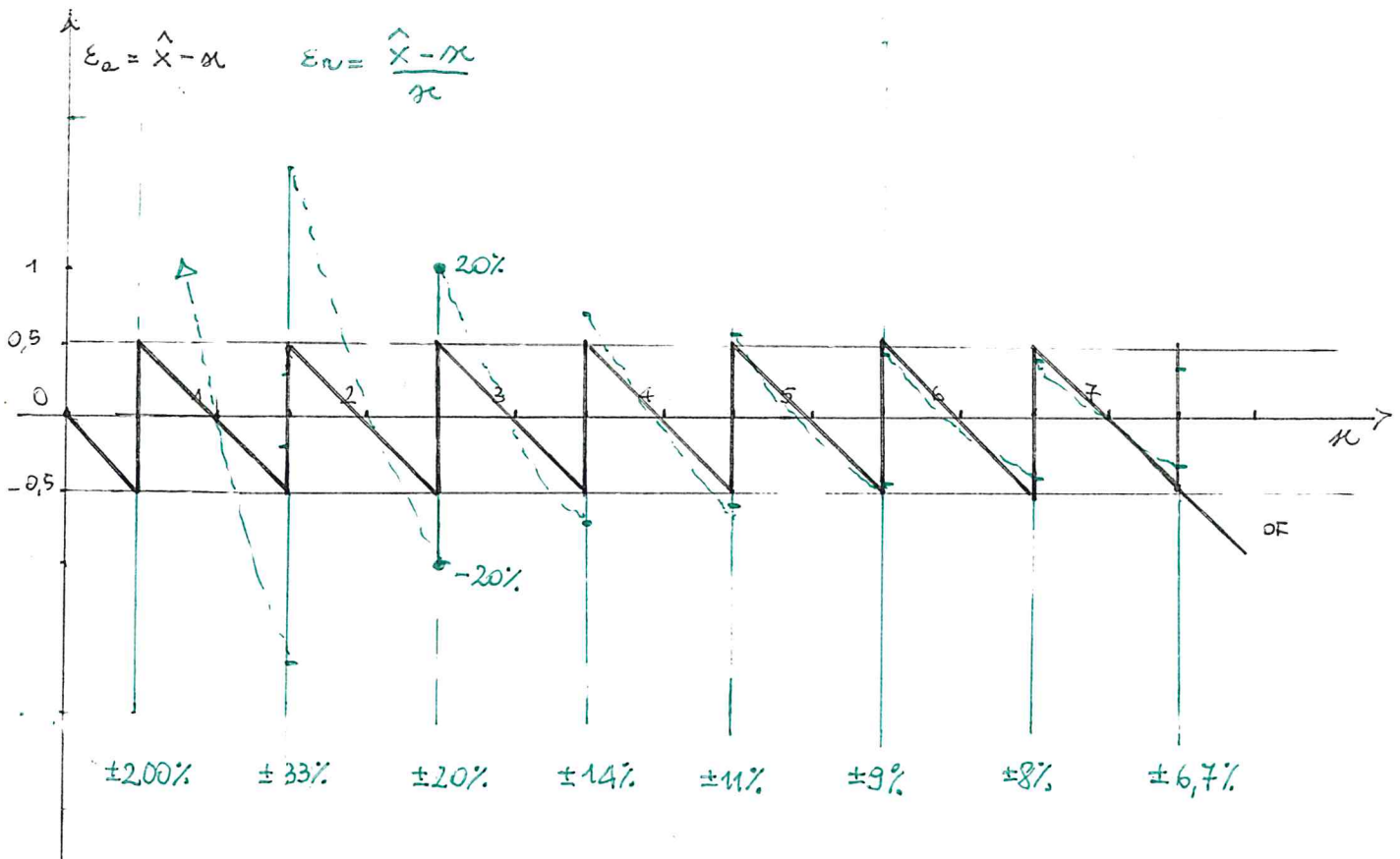
$$\frac{LSB}{2} < \epsilon \leq \frac{LSB}{2} \quad \text{arrotondamento}$$

> Il range coperto dai codici è  $2^m \text{ LSB}$

> L'errore ASSOLUTO è COSTANTE; nel mondo reale però ci interessa di più l'errore RELATIVO

- Un millimetro di errore è trascurabile nella misura di qualche metro  
 è intollerabile se misuriamo pochi millimetri

- Andamento di  $\epsilon_A$  e  $\epsilon_R$  nel caso proposto



• Quindi se intendiamo rappresentare grandezza da 1mm a 10m, con un errore relativo inferiore al 10%, quanti bit servono? (1)  
 Se l'unità di misura desiderata è il metro, che rappresentazione useremo? (2)

(1) Usando l'arrotondamento, dovremo avere una risoluzione di 0,2mm (l'errore è METÀ della risoluzione con arrotondamento)  
 Serviremo  $\frac{10m}{0,2mm} = 50 \times 10^3$  codici come minimo, quindi 16b  
 In questo caso l'unità di misura è 0,2mm e la rappresentazione non avrà parte frazionata.

(2) Se vogliamo le misure in metri serviremo 4b per la parte intera (10 < 16) e 13b per la parte frazionata

$$\frac{1m}{8192} < 0,2mm \quad \text{quindi } (4.13)$$

Osservazione: le rappresentazioni tra 0 e 1mm avranno errore relativo maggiore, divergenti  $\pi \rightarrow \phi$

Un modo più efficace di usare i bit per ridurre  $E_r$

• Dedicare un po' di bit a specificare l'ORDINE DI GRANDEZZA (come nella NOTAZIONE SCIENTIFICA)

Facciamo riferimento al problema precedente - Per avere il 10% (arrotondamento) tra 1 e 2mm ci servono 3 bit.

$$1. \boxed{b_1 b_2 b_3}$$

Per arrivare a 10m ci sono 4:000  $10^4 < 2^{14}$   
 Possiamo specificare l'OrdG<sub>2</sub> con 4 bit. Con 7b otteniamo la specificazione richiesta abbondantemente

$$\boxed{f_1 | f_2 | f_3} \quad \boxed{e_3 e_2 e_1 e_0} \quad x = \left(1 + \frac{f}{8}\right) \cdot 2^e$$

1.000 ÷ 1.875 range con e = 0 (in mm)  
 2.000 ÷ 3.750 con e = 1

⋮

32,8 ÷ 61,4 con e = 15 (in metri)

- La tecnica di separare esponente, cifre significative e con l'eventuale aggiunta del segno prende il nome di rappresentazione in virgola mobile (FLOATING POINT)
  - > Molto diffusa per il calcolo numerico
  - > Per rappresentare numeri con grande range dinamico
- Da diversi decenni, per rendere INTEROPERABILI sistemi diversi e permettere lo scambio di informazione, sono stati proposti STANDARD per la rappresentazione F.P.
  - Uno STANDARD è costituito da un insieme di prescrizioni o definizioni a cui il sistema (o il LINGUAGGIO SOFTWARE) si deve conformare
  - Nel caso degli standard di rappresentazione numerica sono importanti sia l'HW sia il SW che insieme concorrono a definire il comportamento del sistema

Lo standard IEEE 754-2008 - Precisione

- Per approfondire i concetti delle rappresentazioni F.P., usiamo come riferimento lo standard corrente usato dalla maggior parte dei sistemi.

> 58 pagine, 11 capitoli, 2 appendici informative

**Oggetto:** formati e metodi per la gestione dell'aritmetica F.P. in sistemi elettronici di elaborazione

- precisione singola, doppia, estesa, estendibile
- formati di scambio
- condizioni di "eccezione" e modi di gestione

**Scopo:** Sistemi conformi che elaborano gli STESSI DATI dovrebbero portare allo STESSO RISULTATO (indipendentemente dall'implementazione)

**Cose e' incluso:** > Formati (base 2 e base 10) per ELABORAZIONE e per lo SCAMBIO dei DATI

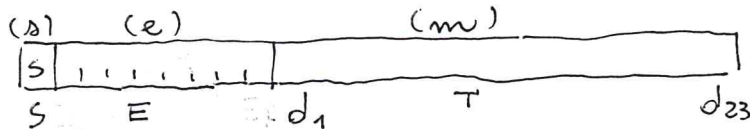
- > Comportamento di addizione, sottrazione, moltiplicazione, divisione,  $\text{mod}$ ,  $\sqrt{\quad}$ , confronti e altre operazioni
- > Conversione INTERI  $\leftrightarrow$  F.P.
- > Conversione tra F.P. diversi
- > Conversione tra F.P. e rappresentazioni esterne (es: stringhe di caratteri)
- > Eccezioni e relative gestione

Formati proposti

- Base 2 o 10
- Diverse cifre per la parte significativa
- Diverso range per l'esponente
- Rappresentazione di  $\pm\infty$
- NaN (di vario tipo) quiet e signaling

BINARY32

- 32 bit divisi in 3 campi



regole di rappresentazione

$$x = (-1)^s \cdot 2^e \cdot m$$

"s"  $\bar{=}$  0 (positivi)  
1 (negativi)

"e" esponente  $e_{min} \leq e \leq e_{max}$   
-126 127

( $e_{min} + e_{max} = 1$   
per tutti i formati)

"m" significando frazionale  $d_0 \cdot d_1 d_2 \dots d_{23}$

- Per rendere univoca la rappresentazione, "m" è massimizzato (riducendo allo stesso tempo "e") fino a che  $m \geq 1$  oppure  $e = e_{min}$

- la rappresentazione di  $e$  è in treslezione

$$e = E - 127$$

quindi i valori "normali" di E sono  $E = e + 127$

$$1 \leq E \leq 254$$

Ricapitolando

S	E	T	valore
S	1..254	T	$x = (-1)^s \cdot 2^{E-127} \cdot (1 + T \cdot 2^{-23})$
S	255	$\emptyset$	$x = \pm\infty$ secondo S
S	255	$\neq \emptyset$	NaN (SNaN: $d_1 = \emptyset$ ; qNaN: $d_1 = 1$ )
S	0	$\neq 0$	$x = (-1)^s \cdot 2^{-126} \cdot T \cdot 2^{-23}$ (non normal.)
S	0	0	$x = \pm\emptyset$ secondo S



## Arrotondamento

#6.5

- lo standard prevede "attributi sulle direzioni dell'arrotondamento".  
Quindi i sistemi devono essere in grado di implementare diverse strategie.

> L'esigenza di arrotondare nasce dal fatto che i risultati di molte operazioni non sono NUMERI DI MACCHINA e quindi occorre rappresentarli con approssimazione

> L'idea è quella di avere il RISULTATO ESATTO, con infinite cifre e senza limiti di RANGE e poi applicare l'arrotondamento (con eventuale OVERFLOW)

## Operazioni

- lo standard prevede diverse operazioni ed esamina il comportamento da intraprendere in seguito e possibili eccezioni. Qui ci limitiamo a riflettere sugli algoritmi da implementare e sui possibili eventi di cui tenere conto nel caso di NUMERI NORMALIZZATI.

## PRODOTTO A · B

$$S_p = \begin{cases} 0 & \text{se uguali } S_A = S_B \\ 1 & \text{se diversi } S_A \neq S_B \end{cases}$$

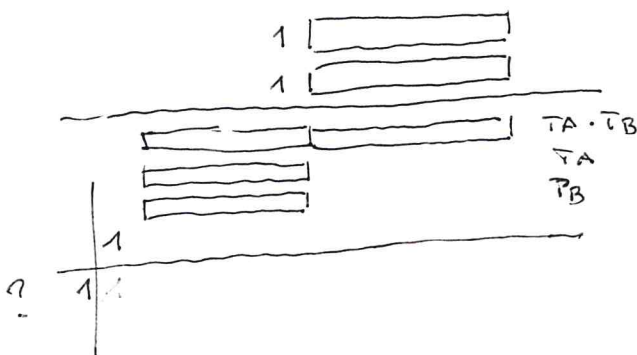
$$E_p = E_A + E_B - 127 \quad \begin{array}{l} \text{ci può essere overflow o underflow se} \\ E_p > 254 \text{ oppure } E_p < 1 \end{array}$$

$$(1 + T_A \cdot 2^{-23}) (1 + T_B \cdot 2^{-23}) = 1 + (T_A + T_B) 2^{-23} + T_A T_B 2^{-46}$$

questo risultato può venire da 1 a poco meno di 4  
se viene  $1 \leq R < 2$  si arrotonda

$$T_p = T_A + T_B + T_A T_B 2^{-23}$$

se viene  $2 \leq R < 4$  si rinnormalizza incrementando  $E_p$



## DIVISIONE A/B

$$A = (-1)^{S_A} 2^{E_A - 127} (1 + T_A 2^{-23})$$

$$B = (-1)^{S_B} 2^{E_B - 127} (1 + T_B 2^{-23})$$

$$\frac{A}{B} = (-1)^{S_R} 2^{E_A - E_B + 127 - 127} \frac{1 + T_A 2^{-23}}{1 + T_B 2^{-23}}$$

con

$$S_R = 1 \text{ se } S_A \neq S_B$$

$$0 \text{ se } S_A = S_B$$

$E_R = E_A - E_B + 127$  ci può essere overflow o underflow

il significando può andare

$$\frac{1}{2} < m_R \leq 2$$

quindi potrebbe essere necessario ( $T_B > T_A$ ) rinormalizzare  
DECREMENTANDO E

## SOMMA e DIFFERENZA

A±B

#6.7

Le due operazioni sono simili (per la differenza basta cambiare il segno di B). Il comportamento cambia se i segni sono CONCORDI o DISCORDI

> Segni concordi  $S_R = S_A = S_B$

Sia  $E_A \geq E_B$

$$A+B = (-1)^{S_R} 2^{E_A-127} \left[ 1 + T_A 2^{-23} + (1 + T_B 2^{-23}) 2^{E_B-E_A} \right]$$

se la permuta viene  $\geq 2$  (poterò sempre  $< 4$ ) si incrementa  $E_R$  e si rinnormalizza - Potrebbe esserci OVERFLOW

> Segni discordi con  $E_A > E_B$ ;  $S_R = S_A$

$$A-B = (-1)^{S_R} 2^{E_A-127} \left[ 1 + T_A 2^{-23} - \underbrace{(1 + T_B 2^{-23}) 2^{E_B-E_A}} \right]$$

→ valore minimo

$$E_B - E_A = -1; T_A = \phi; T_B = 2^{23} - 1 \quad \text{si ha}$$

$$\left[ 1 - \frac{1}{2} - \frac{1}{2} + 2^{-23} \right]$$

Occorre rinnormalizzare ottenendo  $E_R = E_A - 23$  (rischio underflow)

> Segni discordi con  $E_A = E_B$ ;  $T_A > T_B$  (altrimenti si potrebbe avere  $\phi$ )

$$S_R = S_A$$

$$A-B = (-1)^{S_R} 2^{E_A-127} \cdot (T_A - T_B) \cdot 2^{-23}$$

→ valore minimo  $T_A - T_B = 1$  ed è simile al caso limite precedente con  $E_R = E_A - 23$

