ERRORS of FIXED POINT REPRESENTATION $\left( \text{Hp. } X \left( \beta = 10, t = 4, q = 1 \right) \right)$

### ABSOLUTE ERROR   (ROUNDOFF)

$$E(x) = x - flp(x)$$

error associated to rounding of given real number $x$

"FIXED POINT"

Example.  $x_1 = \frac{10^3}{3} = 333.\overline{3}$

$x_2 = \frac{1}{3} = 0.\overline{3}$

closest member of $X$

CONSTANT ABSOLUTE ERROR ←

$$E(x_1) = 333.\overline{3} - 333.3 = 0.0\overline{3}$$

$$E(x_2) = 0.\overline{3} - 0.3 = 0.0\overline{3}$$

### RELATIVE ERROR   $e(x) = \left| \dfrac{E(x)}{x} \right|$

$$e(x_1) = \frac{0.0\overline{3}}{333.\overline{3}} = \frac{3.\overline{3} \cdot 10^{-2}}{3.\overline{3} \cdot 10^2} = 10^{-4} \Rightarrow \text{FOUR DIGITS of ACCURACY}$$

$$e(x_2) = \frac{0.0\overline{3}}{0.3} = \frac{3.\overline{3} \cdot 10^{-2}}{3.\overline{3} \cdot 10^{-1}} = 10^{-1} \Rightarrow \text{ONE DIGIT of ACCURACY}$$

### Pros.

- Simple ⟹ EASY TO implement
  ↓
- FAST (VIDEOGAMES, MICROCONTROLLERS) ← OLD

### Cons.

- VARIABLE relative error (IDEALLY: CONSTANT relative error)

### FLOATING POINT REPRESENTATION

⟹ GOAL: set of numbers ⟶ WIDER RANGE

                    ⟶ CONSTANT rel. error   } for SAME amount

$\Rightarrow$ GOAL: set of numbers $\nearrow$ WIDER RANGE $\qquad$ } for SAME amount

$\qquad\qquad\qquad\qquad\qquad\qquad$ $\searrow$ CONSTANT. rel. error $\qquad$ of memory

NORMALIZED
REPRESENTATION:
$\downarrow$
for any number $\neq 0$

$$X = \text{sign}(x)\left[\sum_{k=0}^{\infty} d_k \beta^{-k}\right]\beta^{P} \quad \swarrow \in \mathbb{N}$$

$\underbrace{\qquad\qquad}_{m \ "MANTISSA"}$

$d_k$: DIGITS

$0 \leq d_k \leq \beta - 1$

$d_0 \neq 0$

FIRST DIGIT of
mantissa cannot
be zero

$1.12 = 01.12 \ldots$

## FLOATING POINT SET of COMPUTER NUMBERS

UNION
$\downarrow$

$$F(\beta, t, L, U) = \{0\} \cup \left\{ X \in \mathbb{R} = \text{sign}(x)\left[\sum_{k=0}^{t-1} d_k \beta^{-k}\right]\beta^{P} \right\}$$

$\qquad$ $\uparrow$ $\quad$ $\uparrow$
$\qquad$ BASE $\quad$ DIGITS for
$\qquad\qquad$ MANTISSA

$\uparrow m$

LOWER (L)
UPPER (U)
boundaries for
$P$
$P \in [L, U]$

EXAMPLE: $\quad F(\underbrace{\beta = 10, t = 3, L = 0, U = 9}_{4 \text{ DIGITS}})$

$0 \leq d_k \leq \beta - 1$

$\text{MAX}(F): \left[\underset{9\cdot10^0}{d_0\beta^0} + \underset{9\cdot10^{-1}}{d_1\beta^{-1}} + \underset{9\cdot10^{-2}}{d_2\beta^{-2}}\right]\underset{10^9}{\beta^U} = 9.99 \cdot 10^9$

$\updownarrow \sim 10^{10}$

RANGE

$X(\beta=10, t=4, q=7) \sim 10^4$

$d_0 \neq 0$
$\downarrow$

$\text{MIN}(F) \left[1 \cdot 10^0 + 0 \cdot 10^{-1} + 0 \cdot 10^{-2}\right]\underset{10^0}{\beta^L} = 1$

$\dfrac{m \qquad \beta^P}{}$

| $m$ | $\beta^r$ |
|-----|-----------|
| 9.99 | $10^9$ |
| 9.98 | $10^9$ |
| $\vdots$ | $\vdots$ |
| 1.01 | $10^9$ |
| 1.00 | $10^9$ |
| 9.99 | $10^8$ |
| 9.98 | $10^8$ |
| $\vdots$ | $\vdots$ |
| 1.00 | $10^8$ |
| $\vdots$ | $\vdots$ |
| 9.99 | $10^0$ |
| 9.98 | $10^0$ |
| $\vdots$ | |
| 1.00 | $10^0$ |
| 0 | |

$\Delta = 10^7$

$\Delta = 10^6$

$\Delta = 10^{-2}$

ABSOLUTE ERROR
NOT CONSTANT

$$E(x) = x - flp(x)$$

FLOATING POINT
REPRESENT. of $x$

## RELATIVE ROUNDOFF ERROR

$$\ell(x) = \left| \frac{E(x)}{x} \right| = \left| \frac{sign(x)\left[ \sum_{k=0}^{\infty} d_k \beta^{-k} \right]\beta^p - sign(x)\left[ \sum_{k=0}^{t-1} d_k \beta^{-k} \right]\beta^p}{sign(x)\left[ \sum_{k=0}^{\infty} d_k \beta^{-k} \right]\beta^p} \right| =$$

$x - flp(x)$

$+1 // -1$

$\geq 1$

$$= \frac{\sum_{k=t}^{\infty} d_k \beta^{-k}}{\sum_{k=0}^{\infty} d_k \beta^{-k}} < \frac{\beta^{1-t}}{1} = \beta^{1-t}$$

$1 \leq m \leq \beta - 1$

CONSTANT // BOUNDED

$$\Rightarrow \ell(x) < \beta^{1-t}$$

$$1 \leq \overset{\circ}{m} \leq \beta - 1$$

$$\sum_{k=t}^{\infty} d_k \beta^{-k} = d_t \beta^{-t} + d_{t+1} \beta^{-(t+1)} + d_{t+2} \beta^{-(t+2)} + \dots$$

$$= \beta^{-t} \left( d_t + d_{t+1} \beta^{-1} + d_{t+2} \beta^{-2} + \dots \right) < \beta^{-t} \cdot \beta =$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{< \beta}$$

$$= \beta^{1-t}$$

$$e(x) < \beta^{1-t} : \text{machine precision}$$

$$\text{if } t = 3 \implies \beta^{1-t} = 10^{-2} \implies \text{ computer numbers}$$
$$\beta = 10 \qquad\qquad\qquad\uparrow \qquad\qquad \text{accurate to the}$$
$$\text{RELATIVE} \qquad\qquad\qquad \text{SECOND DIGIT}$$
$$\text{ACCURACY}$$

## IEEE 754 STANDARD → FOR FLOATING POINT IMPLEMENTATION

GOAL: reproducibility of results

① BINARY FORMAT + HIDDEN BIT TECHNIQUE
$$\beta = 2$$

Example $F\left( \beta = 2, \ t = 2, \ L = -1, \ U = +2 \right)$

$\downarrow$

mantissa $\quad 0 \leq d_k \leq \underbrace{\beta - 1}_{= 1} \ ; \ d_0 \neq 0 \overset{\text{if } \beta = 2}{\implies} d_0 \text{ ALWAYS } = 1$

without Hidden bit $\qquad$ With Hidden bit
$\quad t = 2 \qquad\qquad\qquad t = 2 \to 2 + "1"$

| $d_0$ | $d_1$ | | $d_0$ | $d_1$ | $d_2$ |
|-------|-------|---|-------|-------|-------|
| 1 | 0 | | 1 | 0 | 0 |

$$\overset{2}{\downarrow}$$
$$1\beta^0 + 0\beta^{-1} + 0\beta^{-2} = (1)_{10}$$

| $a_0$ | $a_1$ |
|-------|-------|
| 1 | 0 |
| 1 | 1 |

| $a_0$ | $a_1$ | $a_2$ |
|-------|-------|-------|
| 1 | 0 | 0 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |
| 1 | 1 | 1 |

$\underbrace{\qquad}_{(m)_2}$

$$1\beta^0 + 0\beta^{-1} + 0\beta^{-2} = (1)_{10}$$
$$1\cdot 2^0 + 0\cdot 2^{-1} + 1\cdot 2^{-2} = (1.25)_{10}$$
$\Rightarrow$
$$1\cdot 2^0 + 1\cdot 2^{-1} + 0\cdot 2^{-2} = (1.5)_{10}$$
$$1\cdot 2^0 + 1\cdot 2^{-1} + 1\cdot 2^{-2} = (1.75)_{10}$$

SPACING MANTISSA $\Delta = (0.25)_{10}$ $\qquad (m)_{10}$

$$L = -1, \quad U = 2$$

| $p=-1$ | $p=0$ | $p=1$ | $p=2$ |
|--------|-------|-------|-------|
| $2^{-1}=\frac{1}{2}$ | $2^0=1$ | $2^1=2$ | $2^2=4$ |
| 0.5 | 1 | 2 | 4 |
| 0.625 | 1.25 | 2.5 | 5 |
| 0.75 | 1.5 | 3 | 6 |
| 0.875 | 1.75 | 3.5 | 7 |
| $\Delta = 0.125$ | 0.25 | 0.5 | 1 |

$m \beta^p$

0.2  0

0   0.5   1   2

UNDERFLOW REGION

$x = 4.3$  ROUNDED TO 4

4

8

7

OVERFLOW REGION

② PRECISION LEVELS

[FP 32]  $\underline{\text{SINGLE PRECISION}}$ | $\underline{\text{DOUBLE PRECISION}}$

32 BITS | 64 BITS

| SIGN | 1 BIT | 1 BIT |
|---|---|---|
| EXPONENT | 8 BIT $(2^8 = 256)$ | 11 BIT $(2^{11} = 2048)$ |
| | $\Rightarrow P \in [-126 ; 127]$ | $\Rightarrow P \in [-1022 ; 1023]$ |
| MANTISSA | 23 BIT | 52 BIT |
| $\varepsilon$ | $\beta^{1-(t+1)} \sim \beta^{-23}$ | $\beta^{1-(t+1)} \sim \beta^{-52}$ |
| | $\sim 10^{-7}$ | $\sim 10^{-16}$ |
| | $\uparrow$ | $\uparrow$ |
| | numbers accurate to the 7th DIGIT | accurate to 16th digits |

FP 16 $\Rightarrow$ Floating Point 16 BIT    "QUANTIZATION"
(HALF)