

Recall the format

(Sign). mantissa $\times 2^{\text{Exponent}}$

$M \times 2^E$

$$a = m \times 2^E$$

$$b = p \times 2^F$$

If $E = F$

$$a + b = (m + p) \times 2^E$$

$$a = 1.006 \times 2^1 = 3$$

$$b = 1.006 \times 2^1 = 2$$

Renorm

$$1.01006 \times 2^2 = 5$$

If $E \neq F$

$$a = 3 = 1.1006 \times 2^1$$

$$b = 3/4 = 1.1006 \times 2^{-1}$$

$$1.10000 \times 2^1$$

$$.01100 \times 2^1$$

$$1.11100 \times 2^1$$

$$= \left(1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8}\right) \times 2$$

3 3/4

What about Rounding?

$$a = 1$$

$$b = (1.11)6 \times 2^{-23}$$

(Assume IEEE Single)

add $a + b$

$$1.000 \ 000 \ 000 \ 000 \ 000 \ 000 \ 000 \ 0000$$

$$0.000 \ 000 \ 000 \ 000 \ 000 \ 000 \ 000 \ 0111$$

$$1. \quad \quad \quad 22 \ 06 \quad \quad \quad 111$$

How should we round?

- Floor
- Round to nearest.
- ~~Round~~ Truncate
- Ceil
- Round to an even value.

IEEE Rounding Modes (§ 5.5 in book)

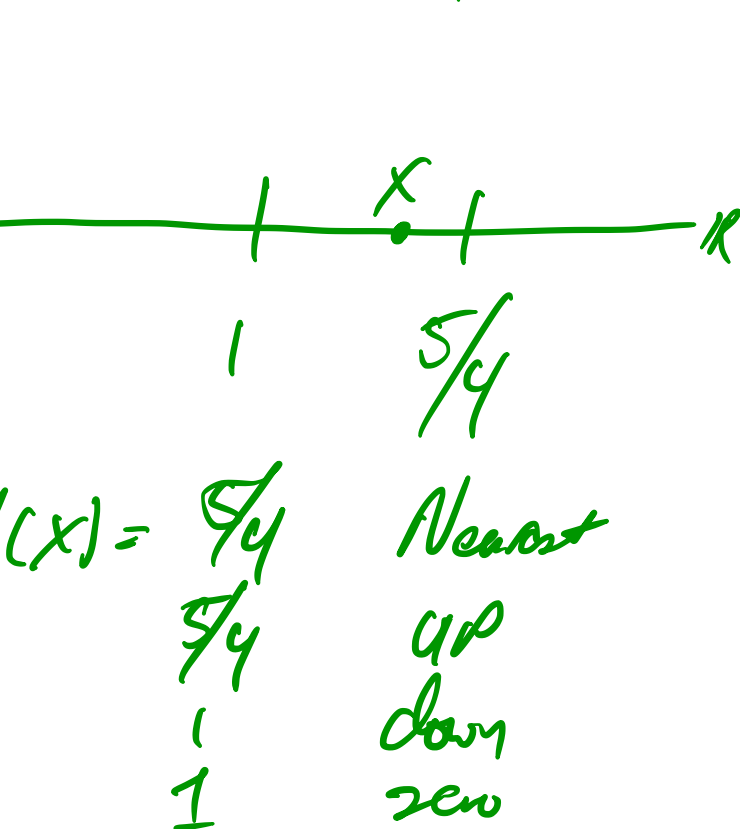
Definition a floating pt # is a value that can be represented exactly in a FP # system.

- Rounding Modes
- Round to Nearest
 - Round up
 - Round down
 - Round to zero

round(x) = FP #

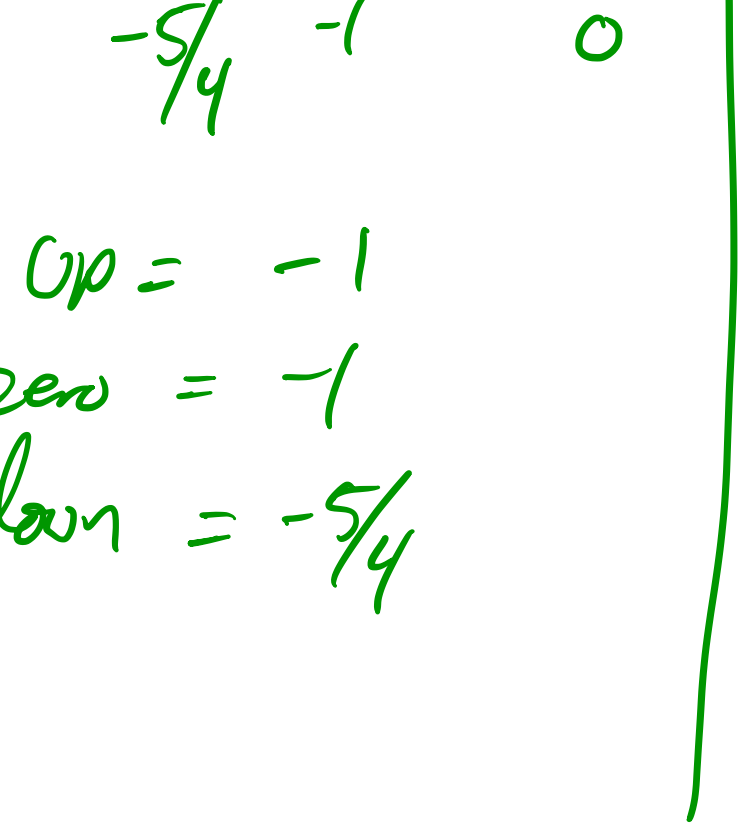
$$x \in \mathbb{R}$$

$$\text{round} : \mathbb{R} \rightarrow \text{FP \#}$$

round(x) = FP #

$$x \in \mathbb{R}$$

$$\text{round} : \mathbb{R} \rightarrow \text{FP \#}$$



Quiz

In our toy FP system

Theorem For all Rounding modes in IEEE

$$\text{round}(x) = \boxed{x(1 + \delta)}$$

FP #

where $|\delta| \leq \epsilon$

machine ϵ .

Pf: Rounding modes only change the last bit.

More notation

$$\overset{\text{FP \#}}{a \oplus b} = \text{round}(a + b)$$

↑ ↑ May not be a FP #

FP # FP #

$$= (a + b)(1 + \delta_1)$$

$\approx | \delta_1 | \leq \epsilon$

$a \oplus b$ is what you get on your computer

IEEE guarantees

$$a \oplus b = (a + b)(1 + \delta_1)$$

$$a \ominus b = (a - b)(1 + \delta_2)$$

$$a \otimes b = (a \cdot b)(1 + \delta_3)$$

$$a \oslash b = (a / b)(1 + \delta_4)$$

Implement this guarantee is HARD