# Floating Point

Math 426

University of Alaska Fairbanks

September 25, 2020

# Rounding modes

9-bit format
mantissa: 4 bits

Compute $(1.0000)_2 \times 2^0 + (1.11)_2 \times 2^{-3}$

1.00111

$$
\begin{array}{r}
1 \\
1.0011 \\
0.0001 \\
\hline
1.0100
\end{array}
$$

Modes:

1. up ($\to \infty$):
2. down ($\to -\infty$):
3. zero ($\to 0$):
4. nearest (default)
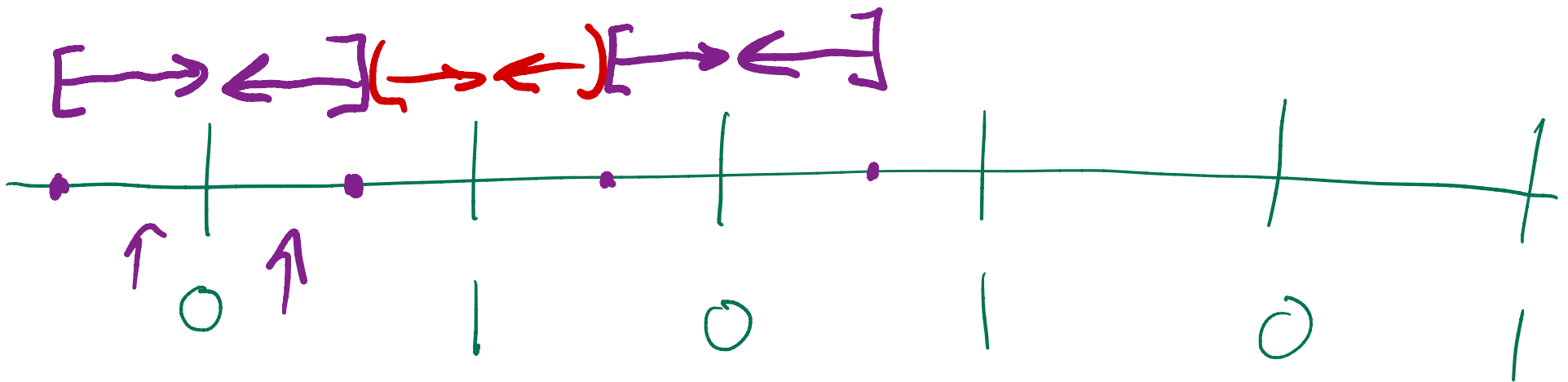
1.0100   a
1.00115

$$
\begin{array}{r}
11 \\
1.0011 \\
0.0001 \\
\hline
1.0100
\end{array}
$$

b         x
          a

# Round to nearest

Special rule if the number is exactly halfway between the two nearest representable numbers: result is the unique nearby representable number with a 0 in its least significant digit.

# Rounding Error

Suppose $2^E \leq x < 2^{E+1}$.

Number line:

$\varepsilon \rightarrow$ machine $\varepsilon$
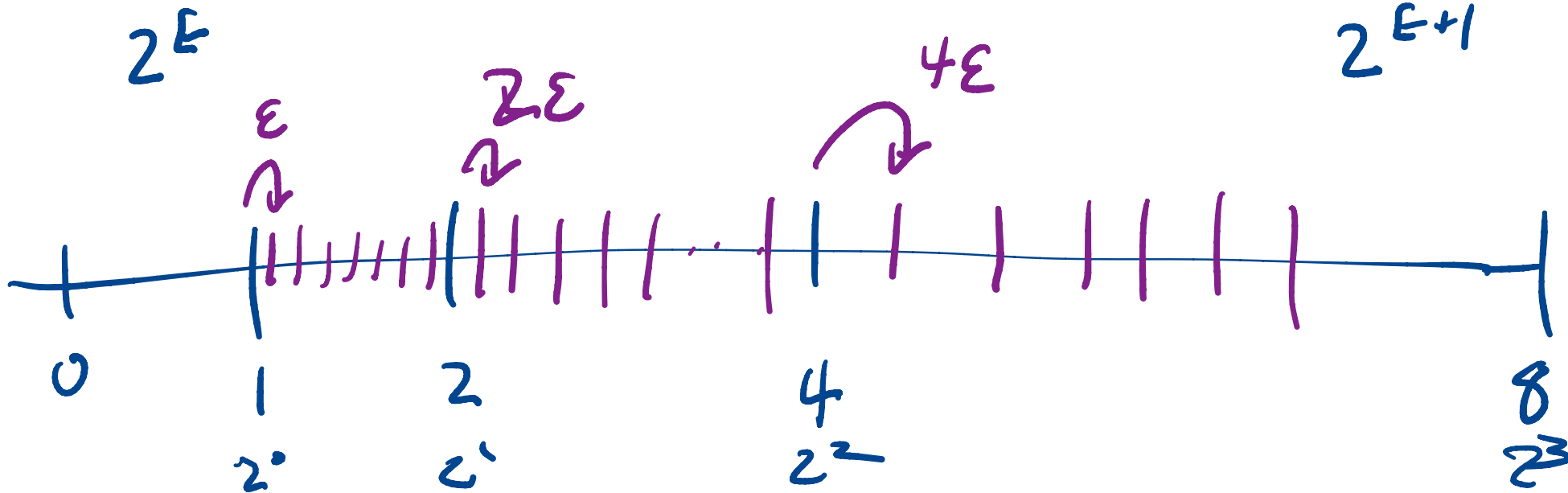
$(-1)^s \; m \cdot \boxed{2^E}$

$\quad \rightarrow$ # of choices $2^k$

$\qquad k$ number of bits

$2^E \varepsilon$



$2^E \qquad\qquad\qquad\qquad\qquad\qquad\qquad 2^{E+1}$

$\varepsilon \qquad\qquad 2\varepsilon \qquad\qquad\qquad 4\varepsilon$



$0 \qquad 1 \qquad\quad 2 \qquad\qquad\qquad 4 \qquad\qquad\qquad\qquad 8$

$\quad\; 2^0 \qquad\; 2^1 \qquad\qquad\; 2^2 \qquad\qquad\qquad\; 2^3$

# Rounding Error

Suppose $2^E \leq x < 2^{E+1}$.

Number line:



$2^E \varepsilon$

$2^E 2\varepsilon$

$2^E$

$2^{E+1}$

all other modes $\longrightarrow$ $\left| x - round(x) \right| \leq 2^E \varepsilon$

round to nearest $\longrightarrow$ $\left| x - round(x) \right| \leq 2^E \varepsilon / 2$

# Rounding Error

Suppose $2^E \leq x < 2^{E+1}$.

Number line:

So: $|x - \mathrm{round}(x)| \leq \epsilon 2^E$.

(or $\epsilon/2 2^E$ for round to nearest)

# Rounding Error

$$\frac{|x - \text{round}(x)|}{|x|}$$

Relative error

Suppose $2^E \leq x < 2^{E+1}$.

$$|x - \text{round}(x)| \leq 2^E \varepsilon$$

$$2^E \leq |x|$$

$$\frac{1}{|x|} \leq \frac{1}{2^E}$$

$$\frac{|x - \text{round}(x)|}{|x|} \leq \frac{2^E \varepsilon}{|x|} \leq \frac{2^E \varepsilon}{2^E}$$

$$= \varepsilon$$

# Rounding Error

$$\frac{|x - \text{round}(x)|}{|x|}$$

Suppose $2^E \leq x < 2^{E+1}$.

$$\frac{|x - \text{round}(x)|}{|x|} \leq \varepsilon$$

$$\left( \leq \frac{\varepsilon}{2} \right) \text{ for round to}$$
$$\text{neorest}$$

# Rounding Error

$$\frac{|x - \text{round}(x)|}{|x|}$$

Suppose $2^E \leq x < 2^{E+1}$.

$1/|x| \leq 2^{-E}$

# Rounding Error

$$\frac{|x - \text{round}(x)|}{|x|}$$

Suppose $2^E \leq x < 2^{E+1}$.

$1/|x| \leq 2^{-E}$

$$\frac{|x - \text{round}(x)|}{|x|} \leq \epsilon 2^E 2^{-E} = \epsilon$$

(or $\epsilon/2$ for round to nearest)

The result of a floating point operation $(+, -, \cdot, /)$ is the correctly rounded value of the exact result.

$$x \oplus y := \text{round}(x + y) = x + y + \text{error}$$

$$\varepsilon \geqslant \frac{|\text{error}|}{|x+y|} = \frac{|\text{round}(x+y) - x+y|}{|x+y|}$$

$$|\text{round}(x+y) - (x+y)| \leqslant \varepsilon |x+y|$$

# IEEE 754 arithmetic

$\ominus \quad \otimes \quad \oslash \qquad\qquad x \otimes y = (x \cdot y)(1 + \delta)$

The result of a floating point operation $(+, -, \cdot, /)$ is the correctly rounded value of the exact result.

$$x \oplus y := \underbrace{\text{round}(x + y)}_{} = x + y + \text{error} \overbrace{\phantom{xxxxxxxx}}^{(x+y)(1+\delta)}$$

$x \oplus y = (x + y)(1 + \delta)$

$\qquad = (x+y) + \underbrace{(x+y)\delta}_{} \qquad \longrightarrow \text{error}$

for some number $\delta$ with $|\delta| < \varepsilon / 2$

$\dfrac{\text{error}}{x+y} = \delta$

# IEEE 754 arithmetic

The result of a floating point operation $(+, -, \cdot, /)$ is the correctly rounded value of the exact result.

$$x \oplus y := \text{round}(x + y) = x + y + \text{error}$$

Similar operations: $\ominus, \otimes, \oslash$.

# Rounding Isn't Easy

Compute $1.0000_2 - 0.00001010_2$ versus $1.0000_2 - 0.00001000_2$ under round to nearest.

$1.0000$
$- 0.0000$

$0.1111\ 11$

$0.000001$

$1.000000$

$0.1111$

$0.0001$

$1.0000$

$0.111111\overset{0}{\cancel{1}}00$
$- 0.00000010$

$0.11111010$

$0.11111$

$1.1111 \times 2^{-1}$

$vs \quad 1.0000 \times 2^{0}$

# Rounding Isn't Easy

Compute $1.0000_2 - 0.00001010_2$ versus $1.0000_2 - 0.00001000_2$ under round to nearest.

Requires extra bits (2 suffice for most: guard bits). Special cases requrie more (sticky bit for flag).