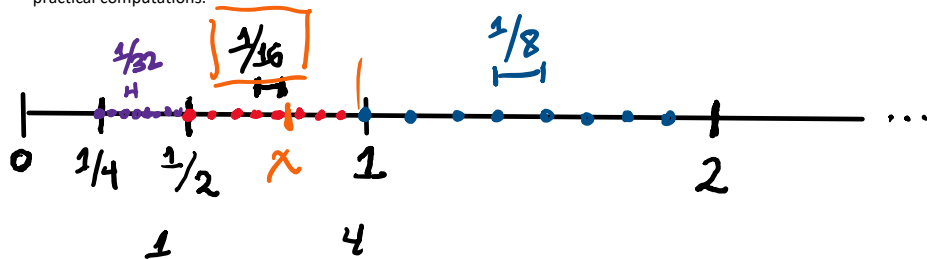**Class 02: Wed, August 28**

**Recall:** After introducing the course, we described the differences between math with real numbers on pen and paper and in the computer. First big thing to consider: we have limited memory budget per real number. Once we have assigned a number of "bits" (zeroes and ones) for our budget, we discussed how to best build a system to do practical computations.

~~64 bits memory.~~

$$x = \pm \underline{\phantom{xxxx}} . \underline{\phantom{xxxx}} |$$

~~not practical!~~

$$x = \pm (0.d_1 d_2 d_3 d_4)_2 \times 2^e \qquad e \in \{-3, -2, -1, 0, 1, 2, 3, 4\}$$

has to be 1

$$e = (e_1 e_2 e_3)_2 - 3$$

bias

$x$ true value

$fl(x) \rightarrow$ closest rounding

- MAX NUMBER?
- MIN NUMBER?
- # of floats for $e = 0$:
- TOTAL # of floats:

## DEFINITIONS / NOTATION:

$fl(x) \rightarrow$ rounding $x$ to nearest fl pt #.

abs error $\rightarrow |x - fl(x)| < \frac{1}{32}$

rel error $\rightarrow \dfrac{|x - fl(x)|}{|x|} < \dfrac{\frac{1}{32}}{\frac{1}{2}} = \boxed{\frac{1}{16}}$

| Uniform bound for rel error |

$\rightarrow$ Spacing between # in $[\frac{1}{2}, 1)$ $(e=0)$

"Machine epsilon"

## DOUBLE PRECISION (IEEE 754)

DOUBLE PRECISION (IEEE 754)

↳ 64 bits: 
- 1 sign.
- 52 mantissa.
- 11 exponent.

$$x = \pm(0.1\, d_2\, d_3 \cdots d_{52})_2 \times 2^e$$

$$e = \{0, 1, \ldots, 2047\} - 1023$$

$$= \{\cancel{-1023}, -1022, \ldots, 1023, \cancel{1024}\}$$

0                                           $\pm\infty$, NaN

$$x_{MAX} \approx 10^{308}$$
$$x_{MIN} \approx 10^{-308}$$

$$\varepsilon_{MACH} = 2^{-52} \approx 2 \times 10^{-16}$$

$$\log_{10}(\varepsilon_{MACH}) \approx -16$$

- DOUBLE ⟷ "16 decimal digits (of rel. accuracy"

- 32 bit — "SINGLE PRECISION"

$$\varepsilon_{MACH} \approx 10^{-8}, \quad 8 \text{ decimal digits.}$$

ARITHMETIC OPERATIONS

$$\begin{array}{l} x + y \\ x - y \end{array} \left\{ \begin{array}{l} \text{INPUTS} \\ fl(x), fl(y) \end{array} \right.$$

$$\left.\begin{array}{c} x - y \\ xy \\ x/y \end{array}\right\}$$

$fl(x), fl(y)$

$\overset{\curvearrowright}{\boxed{+}} - fl\left(fl(x) + fl(y)\right)$

$\underline{\phantom{REL ERROR}}$

REL ERROR
INPUTS

$\left(\leq \varepsilon_{MACH}\right)$

REL ERROR
OUTPUTS?

## LOSS OF PRECISION

x, y have the same sign.

$\hookrightarrow$ x + y  ✓
$\qquad$ xy  ✓*
$\qquad$ x/y  ✓*

$\left.\begin{array}{c} \\ \\ \\ \end{array}\right\}$

$x - y$  ✗

same sign
same magnitude.

---

x = 0.52345
y = 0.52343
5 digits rel' acc.

$x - y = 0.00002$
$\qquad = 0.2 \times 10^{-4}$