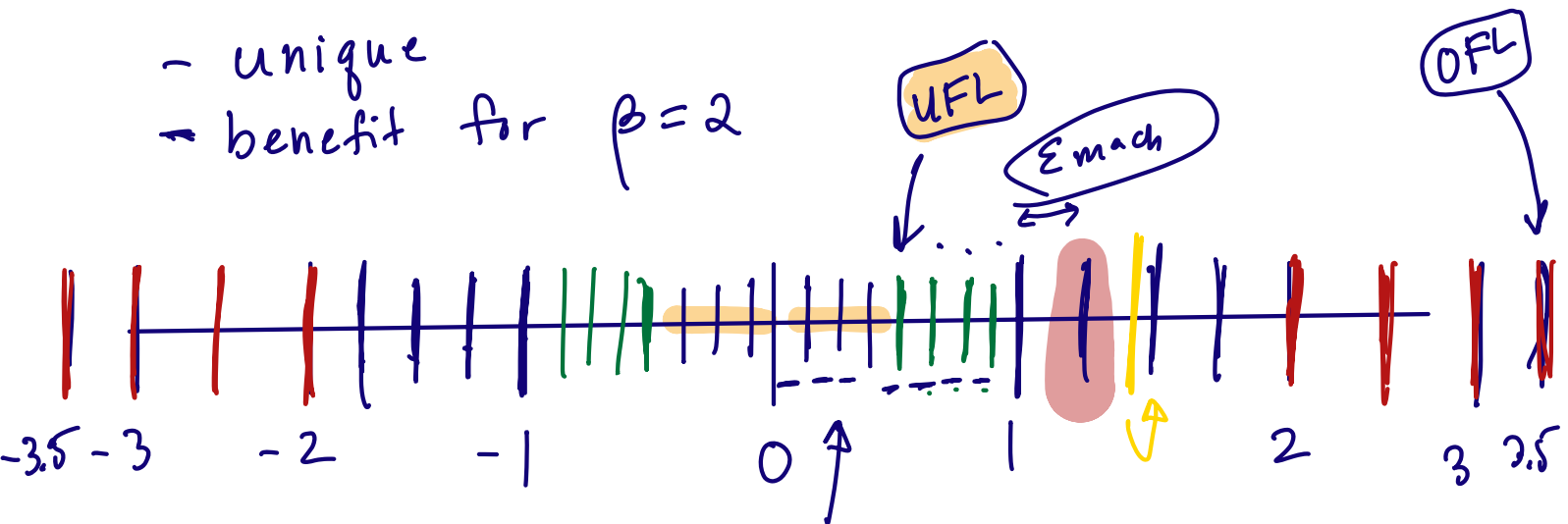


$$\beta = 2 \quad p = 3 \quad L = -1 \quad u = 1$$

Normalization

leading digit in mantissa not 0

- unique
- benefit for $\beta = 2$



subnormals
denormalized #'s

gradual underflow

$$\begin{array}{c} \underline{0} . \underline{0} \quad \dots \quad \underline{1} \\ 0 . 010 \dots 01 \end{array}$$

$$\frac{E}{0}$$

$$0$$

$$\begin{array}{c} \underline{\quad} . \underline{\quad} \underline{\quad} \\ 0 . \overset{\downarrow}{0} \overset{\downarrow}{1} \\ 0 . 10 \end{array} \rightarrow \begin{array}{l} E = -1, 0, 1 \\ \frac{1}{4} \times 2^{-1} = \frac{1}{8} \\ \frac{1}{2} \times 2^{-1} = \frac{1}{4} \end{array}$$

0.11

$$= \frac{3}{8}$$

1.0 →

0.0

01 ↗

denormalized #'s precision loss

Exceptional values

Inf

- #/0

- exceed OFL

NaN

0
0

Inf
Inf

0. Inf

$$\lim_{x \rightarrow 0} \frac{x}{x^2} = \lim_{x \rightarrow 0} \frac{1}{x} \rightarrow \infty$$

$$\lim_{x \rightarrow 0} \frac{x^2}{x} = \lim_{x \rightarrow 0} x \rightarrow 0$$

Rounding

$$x \rightarrow fl(x)$$

chop
round to nearest

	chop	
1.649	1.6	1.6
1.650	1.6	1.6
1.651	1.6	1.7
⋮	⋮	
1.701	1.7	1.7

machine epsilon

$$\epsilon_{mach} = \beta^{-(p-1)}$$

while $((1+\delta) == 1)$

$\delta *= 2$

float (32) $p=24$
double (64) $p=53$

$$2^{-23} \approx 1.19 \times 10^{-7}$$

$$2^{-52} \approx 2.22 \times 10^{-16}$$



unit roundoff

characterizes rounding error

$$\frac{\epsilon_{mach}}{2}$$

$$fl(1.000\underbrace{5}) = \underline{1.001}$$

\times

absolute error

\rightarrow relative error

$$|fl(x) - x|$$

absolute error

$$\frac{|fl(x) - x|}{|x|}$$

relative error

$$fl(x) = s.m \times \beta^E$$

$$\frac{|fl(x) - x|}{|x|} \leq \frac{\frac{1}{2} \epsilon_{mach} \times \cancel{\beta^E}}{\cancel{m} \times \cancel{\beta^E}}$$

$$\leq \frac{1}{2} \epsilon_{mach}$$

if $fl(x)$ is normalized,
 $m \geq 1$

relative error in
storing x as
normalized $f(x)$ is

$$\frac{|f(x) - x|}{|x|} \leq \frac{\epsilon_{\text{mach}}}{2}$$

$$f(0.0055) \\ = 0.006$$

abs err
0.0005

relative err $\frac{1}{6} \frac{\boxed{\frac{0.0005}{0.001}}}{\sim 8\%} = \frac{1}{6} \cdot \frac{1}{2}$

floats 32 $p=24$ $\frac{\epsilon}{2} \quad 1.19 \times 10^{-7}$

double 64 $p=53$ 2.22×10^{-16}

$10^{-5} \%$

0.001
0.0005 ✓

rel. err ~~50%~~
100%

Floating Point Math

$$+ - \quad \times \quad y$$

$$x \quad 1.235 \times 10^5$$

$$y \rightarrow 1.001 \times 10^4$$

$$(1) 2.35 \times 10^4$$

$$1.235 \times 10^5$$

$$0.1001$$

$$1.001$$

$$1.235$$

$$10.000 \dots 1001$$

$$f(x+y) = x$$

$$\underbrace{x + y}_{x} - x = y$$

$$\underbrace{x - x} + y = y$$

cancellation error

$$\begin{array}{r} 1.235217 \\ - 1.23516 \\ \hline 0.000383 \end{array}$$

$$\begin{array}{r} 1.235 \\ - 1.236 \\ \hline - 0.001 \end{array}$$

multiplication

- overflow
- underflow

$$\frac{\cancel{x}^x}{\cancel{x}} = x$$

Responses to student questions after class

$$\underline{1} . \underline{0} \underline{0} \dots \underline{0} \underline{1} \quad p=24$$

$\nearrow 2^{-23}$

general FP system

$$\beta = \text{base } \underset{10}{2} \quad \beta = 2$$

p = digits of precision

L = smallest exponent

$$\underline{-126} \quad (1)$$

U = largest exponent

$$\underline{(127)} \quad (254)$$

biased exponent

$\textcircled{0}$ \downarrow -127	\dots	$\textcircled{255}$ \downarrow 128	\dots	$\underline{\text{bias}}$ 127
$ $				

$$f(x) \quad s \quad m \quad E$$

$$f(x) = s \times \underbrace{m_0 \cdot m_1 \dots m_{p-1}}_m \times \beta^E$$

1	8	23	(+1)	p=24	Ⓢ
	11	52	+1	p=53	