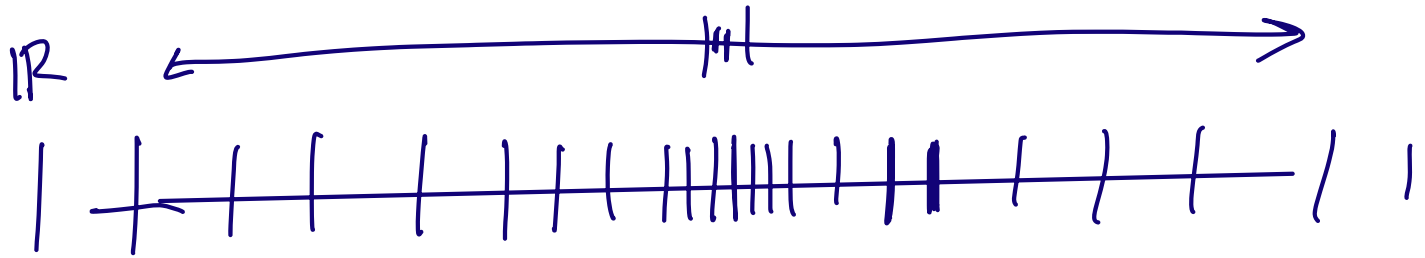


Floating Point Computation



finite precision

$$\underline{1.340} \times 10^{(4)} \quad 13400$$

floating point system

- β base
- P # of digits of precision
- $[L, U]$ exponent range
 $E \in [L, U]$

IEEE 754

	β	P	L	U	
SP floats	2	24	-126	127	$(1+8+23)$
DP doubles	2	53	-1022	1023	$(1+11+52)$

-126
2
2¹²⁷
0...255

floating point number x

$$x = \pm \left(d_0 + \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \dots + \frac{d_{P-1}}{\beta^{P-1}} \right) \times \beta^E$$

$d_0 d_1 \dots d_{P-1}$ mantissa
 $L \leq E \leq U$ exponent

$$0 \leq d_i < \beta$$

$$\beta = 2 \quad p = 3 \quad \underline{L = -1} \quad U = 1$$

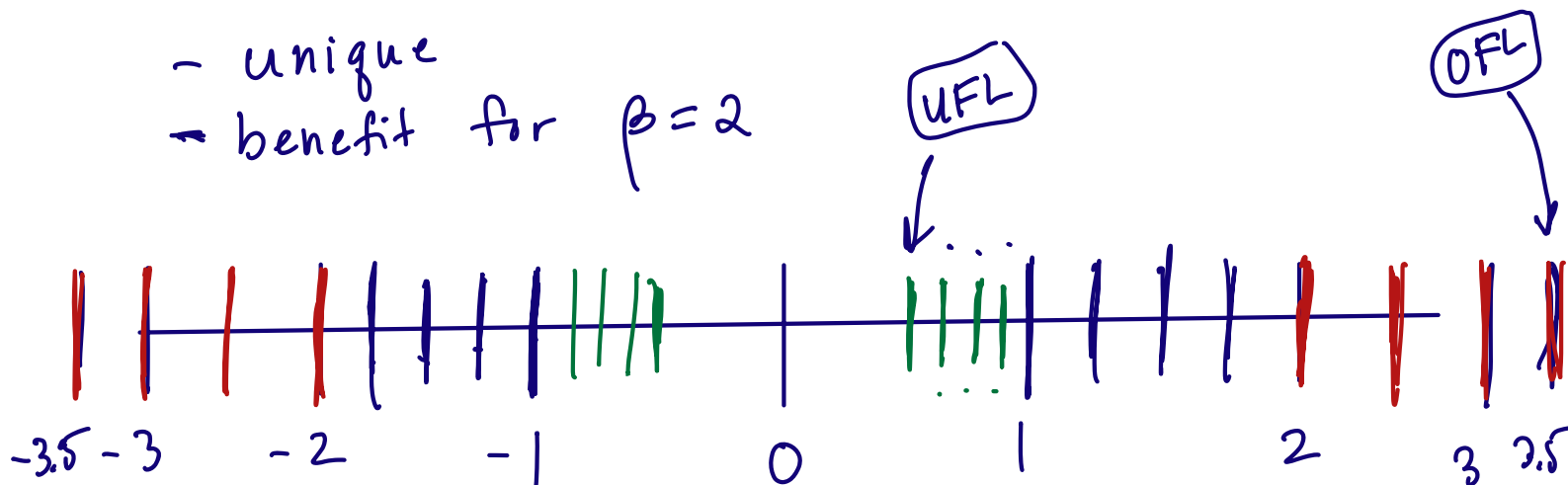
	<u>mantissa</u>	<u>E</u>	<u>number?</u>	
+ -	0.00	-1	± 0	
	0.00	0	± 0	
	0.00	+1	± 0	
	0.01	-1	0.001	
	0.01	0	0.01	
	0.01	+1	0.1	
	0.10	-1	0.010	
	0.10	0	0.10	
	0.10	+1	0.10	
	⋮			

not unique - lots of duplicates!

Normalization

leading digit in mantissa not 0

- unique
- benefit for $\beta = 2$



-1

<u>m antissa</u>	<u>E</u>	<u>number</u>
0.00	.	± 0
1.00	-1	± 0.1
1.00	0	± 1
1.00	+1	$\pm 10.$
1.01	-1	± 0.101
1.01	0	1.01
1.01	+1	10.1
1.10	-1	\vdots
	0	\vdots
	+1	
<hr/>		
1.11	+1	11.1

0.00
0.01
0.10
0.11

$1 \times \boxed{.01 \times 2^{-1}}$

001
1/2 1/4 1/8

$1.00 \times 2^{-1} = .1 \quad \frac{1}{2}$
 \vdots
 $1.11 \times 2^{-1} = .111 = 7/8$
~~10.11~~
 $1.00 \times 2^0 = 1$
 \uparrow
 $.01 \times 2^0 = \frac{1}{4}$

overflow level $\underbrace{(\beta-1) \cdot (\beta-1) \dots (\beta-1)}_P \times \beta^u$

underflow level $1.0 \dots 0 \times \beta^L$

UFL

$$1.00 \times 2^{-1}$$

$$.01 = \frac{1}{4}$$

OFL

$$1.11 \times 2^1$$

$$11.1 = 3.5$$

$$\left(\frac{4}{3}\right) - 1 = \left(\frac{1}{3}\right) \quad \times$$

$$\frac{4}{3} = 1 + \frac{1}{3} \quad \checkmark$$

4 digit decimal

$$1.333(3 \dots 3) - 1$$

$$1.333(-1) = 0.333 \rightarrow$$

$$\frac{1}{3} =$$

$$3.333 \times 10^{-1}$$

$$3.330 \times 10^{-1}$$

$$3.3 \times 10^{-1} \stackrel{?}{=} 3.3 \times 10^{-1}$$

$$1.333 = 1.333 \quad \checkmark$$

$$\begin{array}{r} 3.333 \times 10^{(1)} \\ 1. \times 10^{(0)} \\ \hline \end{array}$$

$$\begin{array}{r} 0.333 \times 10^0 \\ + 1.0 \times 10^0 \\ \hline 1.333 \end{array}$$

$$\varepsilon = 1e-10$$

$$x == y$$

$$\text{abs}(x - y) < \varepsilon$$

$x \neq y$
"equal"