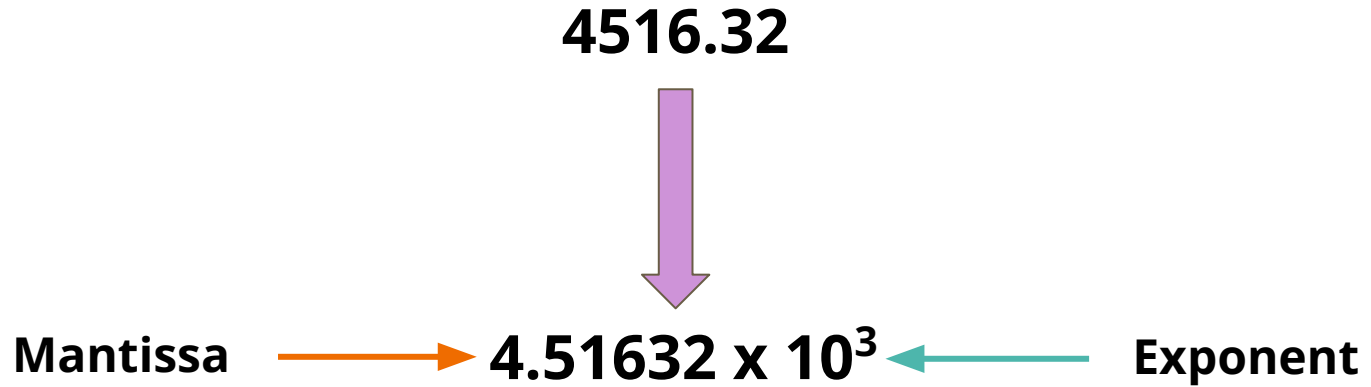

Bits of Architecture

— Floating Point Numbers —

How Do We Represent Represent Non-Integral Numbers?

Recall Scientific Notation



Floating Point Numbers

IEEE-754 Floating Point Standard

- Scientific notation in binary
- Encode non-integral numbers in the bits we have
 - 16b = Half Precision
 - 32b = Single Precision (Floats)
 - 64b = Double Precision (Doubles)
- Similar to instructions (break down bits into field)
 - Sign
 - Exponent
 - Mantissa

Fields for Floats (32b)

Sign	Exponent	Mantissa
1b	8b	23b
-1^{sign}	$2^{(\text{Exponent} - 127)}$	$1 + \text{Mantissa}$

Example

Fields for Floats (32b)

1	0111 1110	000 0000 0000 0000 0000 0000
$-1^1 = -1$	$(126 - 127) = -1$	$(1) + 0 = 1$
-1	2^{-1}	1

-0.5_{10}

But Wait...

Fields for Floats (32b)



But why is this biased?

- Makes comparison easier
- Recall, negative numbers look like large unsigned numbers

Why is there a "1 +"?

- Our number is normalized
- No 0 before decimal point (like scientific notation)

Special Representations

Special Representations

+/- 0		
1 or 0	0000 0000	000 0000 0000 0000 0000 0000
+/- inf		
1 or 0	1111 1111	000 0000 0000 0000 0000 0000
NaN		
1 or 0	1111 1111	Non-Zero
Denormal Numbers		
1 or 0	0000 0000	Non-Zero

Is IEEE-754 Floating Point the Only Way?

No!

- Fixed point
 - Not great, but possible
- Variations on IEEE Floating Point
 - E.g., bfloat
- Posits/Universal Numbers (unums)
 - Interesting research from John Gustafson c. 2015
 - Drop-in replacement for current floating point format