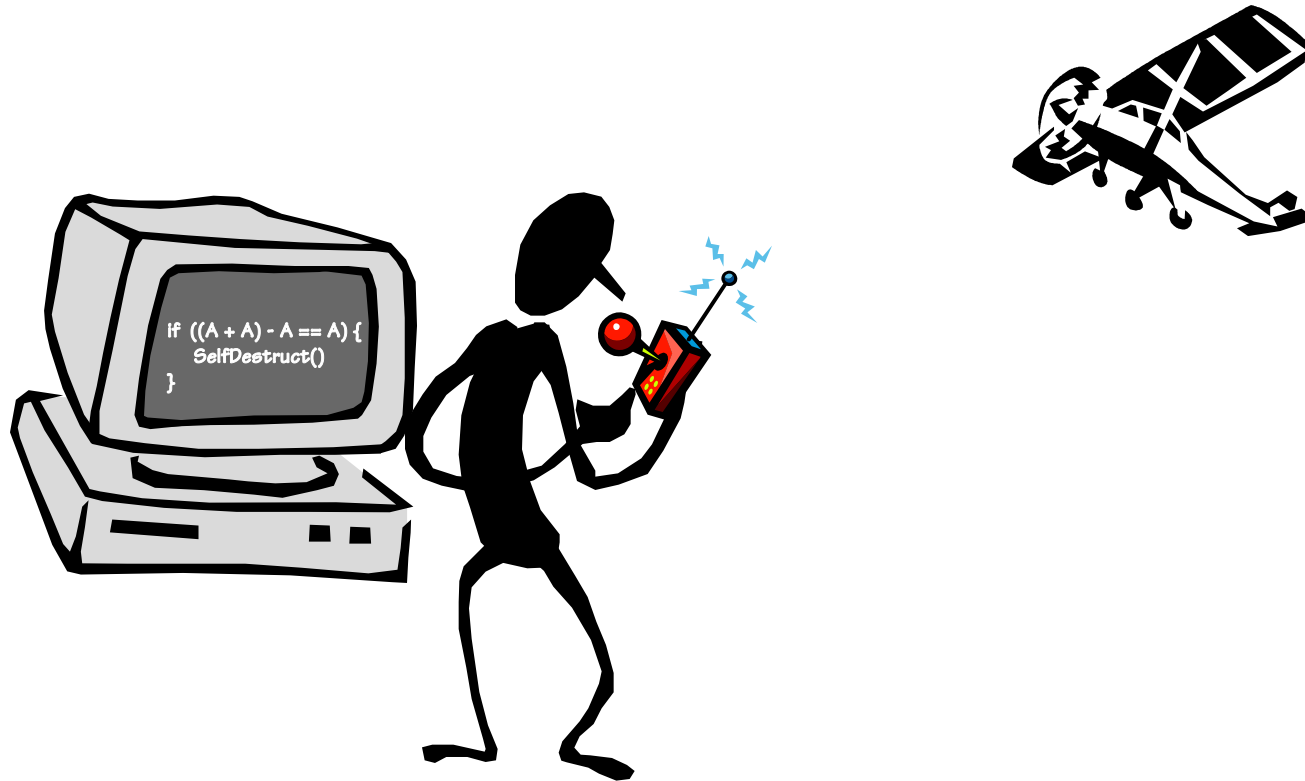


Floating-Point Arithmetic



What is the problem?

Many numeric applications require numbers over a VERY large range. (e.g. nanoseconds to centuries)

Most scientific applications require real numbers (e.g. π)

But so far we only have integers.

We **COULD** implement the fractions explicitly (e.g. $\frac{1}{2}$, $\frac{1023}{102934}$)

We **COULD** use bigger integers

Floating point is a better answer for most applications.

Recall Scientific Notation

- Let's start our discussion of floating point by recalling scientific notation from high school

- Numbers represented in parts:

$$42 = 4.200 \times 10^1$$

$$1024 = 1.024 \times 10^3$$

$$-0.0625 = -6.250 \times 10^{-2}$$

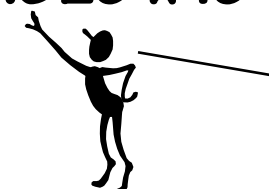
Significant Digits
Exponent

- Arithmetic is done in pieces

$$\begin{array}{r} 1024 \\ - 42 \\ \hline 982 \end{array}$$

$$\begin{array}{r} 1.024 \times 10^3 \\ - 0.042 \times 10^3 \\ \hline 0.982 \times 10^3 \\ \hline 9.820 \times 10^2 \end{array}$$

Before adding, we must match the exponents, effectively "denormalizing" the smaller magnitude number



We then "normalize" the final result so there is one digit to the left of the decimal point and adjust the exponent accordingly.

Multiplication in Scientific Notation

- **Is straightforward:**
 - Multiply together the significant parts
 - Add the exponents
 - Normalize if required

- **Examples:**

$$\begin{array}{r} 1024 \\ \times 0.0625 \\ \hline 64 \end{array}$$

$$\begin{array}{r} 42 \\ \times 0.0625 \\ \hline 2.625 \end{array}$$

$$\begin{array}{r} 1.024 \times 10^3 \\ \times 6.250 \times 10^{-2} \\ \hline 6.400 \times 10^1 \end{array}$$

$$\begin{array}{r} 4.200 \times 10^1 \\ \times 6.250 \times 10^{-2} \\ \hline 26.250 \times 10^{-1} \\ 2.625 \times 10^0 \text{ (Normalized)} \end{array}$$

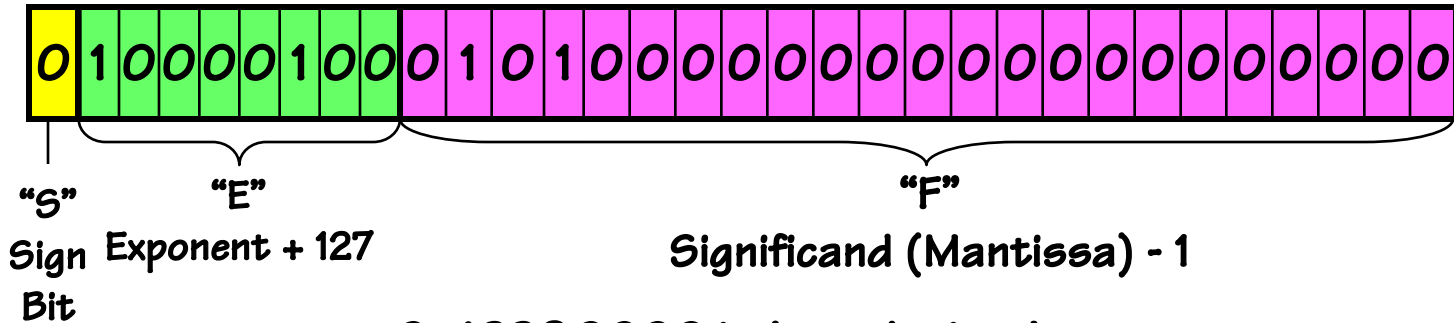
In multiplication, how far is the most you will ever normalize?



In addition?

FP == "Binary" Scientific Notation

- IEEE single precision floating-point format



0x42280000 in hexadecimal

- Exponent: Unsigned "Bias 127" 8-bit integer

$$E = \text{Exponent} + 127$$

$$\text{Exponent} = 10000100 (132) - 127 = 5$$

- Significand: Unsigned fixed binary point with "hidden-one"

$$\text{Significand} = "1" + 0.01010000000000000000000 = 1.3125$$

- Putting it all together

$$N = -1^S (1 + F) \times 2^{E-127} = -1^0 (1.3125) \times 2^5 = 42$$

Example Numbers

- One Sign = +, Exponent = 0, Significand = 1.0

$$-1^0 (1 .0) \times 2^0 = 1$$

$$S = 0, E = 0 + 127, F = 1.0 - '1'$$

0 0111111 000000000000000000000000

0x3f800000

- One-half Sign = +, Exponent = -1, Significand = 1.0

$$-1^0 (1 .0) \times 2^{-1} = \frac{1}{2}$$

$$S = 0, E = -1 + 127, F = 1.0 - '1'$$

0 0111110 000000000000000000000000

0x3f000000

- Minus Two Sign = -, Exponent = 1, Significand = 1.0

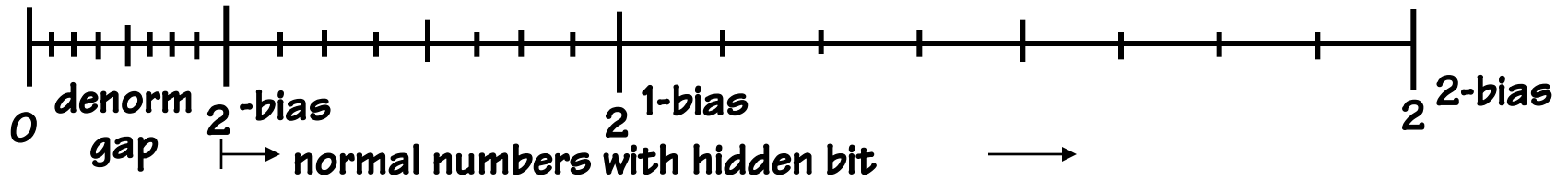
1 10000000 000000000000000000000000

0xc0000000

Infinities

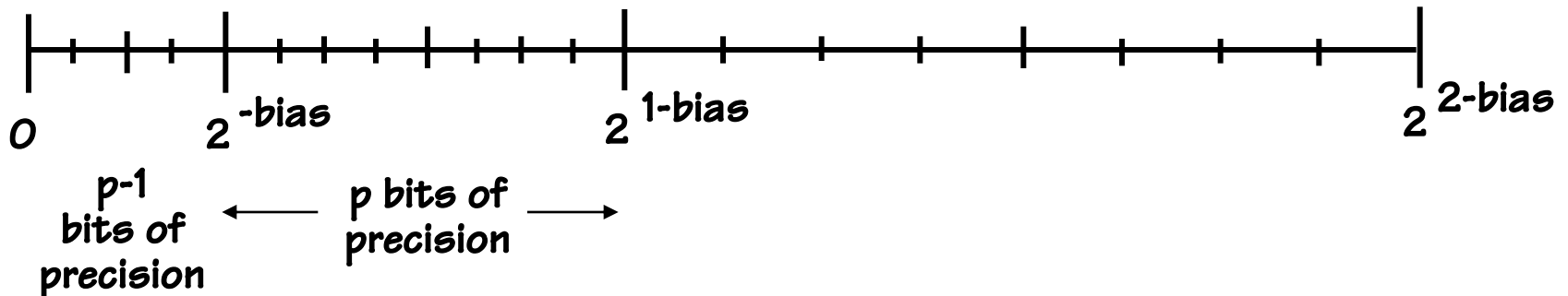
- IEEE floating point also reserves the largest possible exponent to represent “unrepresentable” large numbers
- Positive Infinity: $S = 0, E = 255, F = 0$
 $0\ 11111111\ 0000000000000000000000000000 = +\infty$
 $0x7f800000$
- Negative Infinity: $S = 1, E = 255, F = 0$
 $1\ 11111111\ 0000000000000000000000000000 = -\infty$
 $0xff800000$
- Other numbers with $E = 255$ ($F \neq 0$) are used to represent exceptions or Not-A-Number (NAN)
 $\sqrt{-1}, -\infty \times 42, 0/0, \infty/\infty, \log(-5)$
- It does, however, attempt to handle a few special cases:
 $1/0 = +\infty, -1/0 = -\infty, \log(0) = -\infty$

Low-End of the IEEE Spectrum



The gap between 0 and the next representable normalized number is much larger than the gaps between nearby representable numbers.

IEEE standard uses denormalized numbers to fill in the gap, making the distances between numbers near 0 more alike.



Denormalized numbers have a hidden "0" and a fixed exponent of -126

$$X = (-1)^S 2^{-126} (O.F)$$

NOTE: Zero is represented using 0 for the exponent and 0 for the mantissa. Either, +0 or -0 can be represented, based on the sign bit.

Floating point AIN'T NATURAL

It is CRUCIAL for computer scientists to know that Floating Point arithmetic is NOT the arithmetic you learned since childhood

1.0 is NOT EQUAL to $10 * 0.1$ (Why?)

$$1.0 * 10.0 == 10.0$$

$$0.1 * 10.0 != 1.0$$

$$0.1 \text{ decimal} == 1/16 + 1/32 + 1/256 + 1/512 + 1/4096 + \dots == \\ 0.000110011001100110011\dots$$

In decimal $1/3$ is a repeating fraction $0.333333\dots$

If you quit at some fixed number of digits, then $3 * 1/3 != 1$

Floating Point arithmetic IS NOT associative

$$x + (y + z) \text{ is not necessarily equal to } (x + y) + z$$

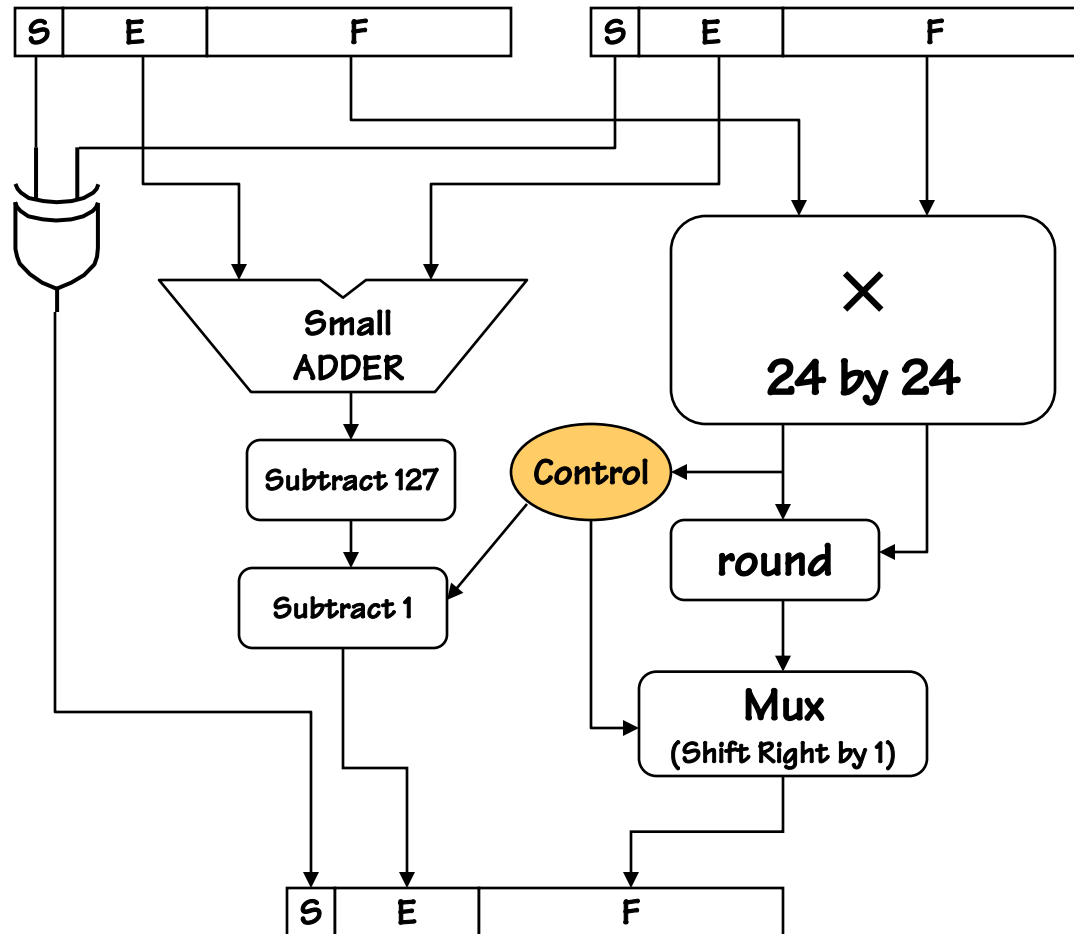
Addition may not even result in a change

$$(x + 1) \text{ MAY } == x$$

Floating Point Disasters

- **Scud Missiles get through, 28 die**
 - In 1991, during the 1st Gulf War, a Patriot missile defense system let a Scud get through, hit a barracks, and kill 28 people. The problem was due to a floating-point error when taking the difference of a converted & scaled integer. (Source: Robert Skeel, "Round-off error cripples Patriot Missile", SIAM News, July 1992.)
- **\$7B Rocket crashes (Ariane 5)**
 - When the first ESA Ariane 5 was launched on June 4, 1996, it lasted only 39 seconds, then the rocket veered off course and self-destructed. An inertial system, produced an floating-point exception while trying to convert a 64-bit floating-point number to an integer. Ironically, the same code was used in the Ariane 4, but the larger values were never generated (<http://www.around.com/ariane.html>).
- **Intel Ships and Denies Bugs**
 - In 1994, Intel shipped its first Pentium processors with a floating-point divide bug. The bug was due to bad look-up tables used in to speed up quotient calculations. After months of denials, Intel adopted a no-questions replacement policy, costing \$300M. (<http://www.intel.com/support/processors/pentium/fdiv/>)

Floating-Point Multiplication



Step 1:

Multiply significands
Add exponents

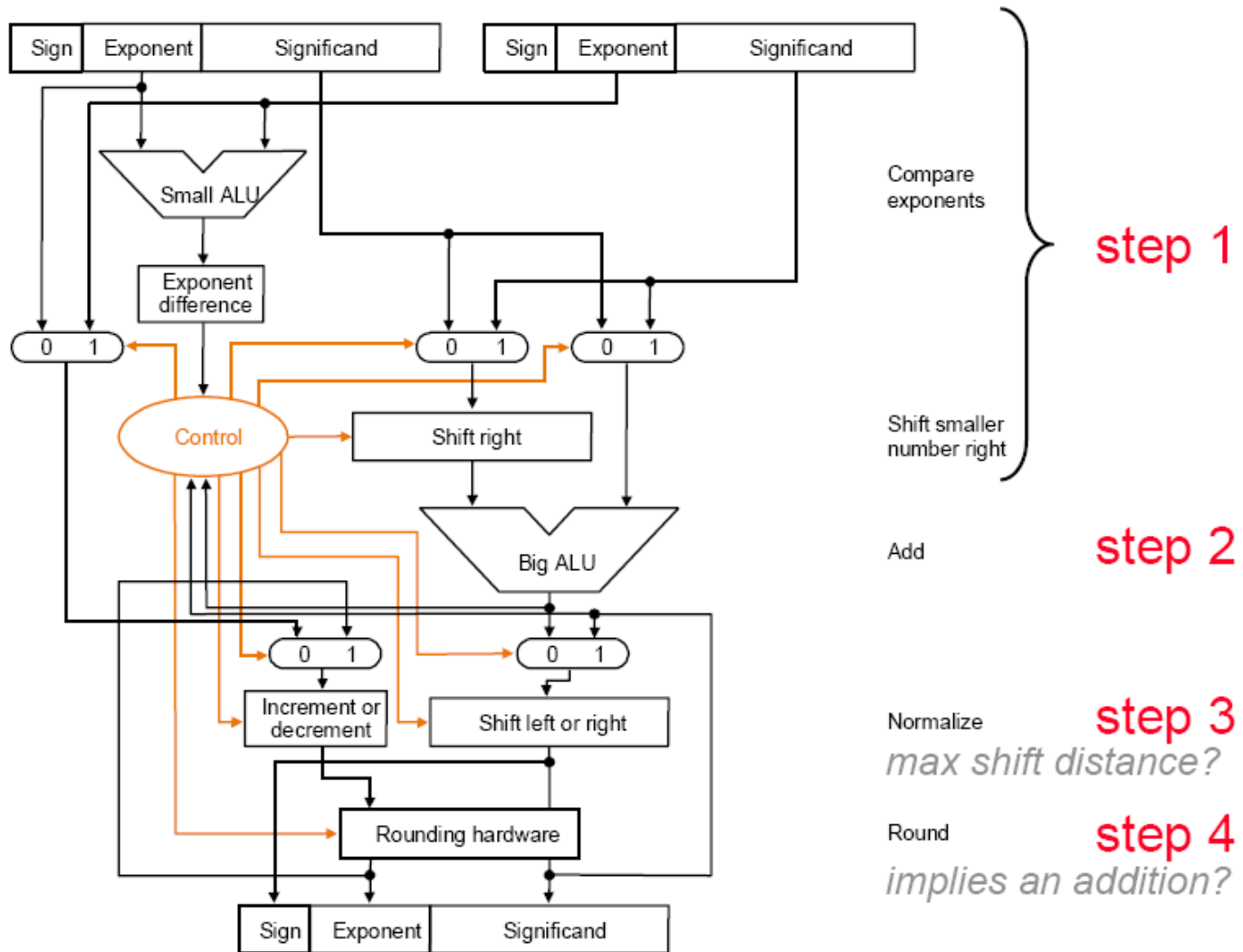
$$E_R = E_1 + E_2 - 127$$

$$\text{Exponent}_R + 127 = \text{Exponent}_1 + 127 + \text{Exponent}_2 + 127 - 127$$

Step 2:

Normalize result
(Result of $[1,2) * [1,2) = [1,4)$
at most we shift right one bit, and fix exponent

Floating-Point Addition



[Figure 3.17 from P&H, COPYRIGHT 2004 Elsevier. ALL RIGHTS RESERVED.]

MIPS Floating Point

Floating point “Co-processor”

32 Floating point registers

separate from 32 general purpose registers

32 bits wide each.

use an even-odd pair for double precision

add.d fd, fs, ft # fd = fs + ft in double precision

add.s fd, fs, ft # fd = fs + ft in single precision

sub.d, sub.s, mul.d, mul.s, div.d, div.s, abs.d, abs.s

l.d fd, address # load a double from address

l.s, s.d, s.s

Conversion instructions

Compare instructions

Branch (bc1t, bc1f)

Chapter Three Summary

Computer arithmetic is constrained by limited precision

Bit patterns have no inherent meaning but standards do exist

two's complement

IEEE 754 floating point

Computer instructions determine “meaning” of the bit patterns

Performance and accuracy are important so there are many complexities in real machines (i.e., algorithms and implementation).

Accurate numerical computing requires methods quite different from those of the math you learned in grade school.