

CS 105

"Tour of the Black Holes of Computing!"

Floating Point

Topics

- IEEE Floating-Point Standard
- Rounding
- Floating-Point Operations
- Mathematical Properties

Floating-Point Puzzles

- For each of the following C expressions, either:

- Argue that it is true for all argument values
- Explain why it is not true, ideally with an example

```
int x = foo();
float f = bar();
double d = baz();
```

Assume neither
d nor f is NaN

Assume a 32-bit
machine

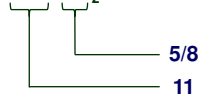
- $x == (\text{int})(\text{float})\ x$
- $x == (\text{int})(\text{double})\ x$
- $f == (\text{float})(\text{double})\ f$
- $d == (\text{float})\ d$
- $f == -(-f)$
- $2/3 == 2/3.0$
- $d < 0.0 \Rightarrow ((d*2) < 0.0)$
- $d > f \Rightarrow -f > -d$
- $d * d \geq 0.0$
- $(d+f)-d == f$

- 2 -

CS 105

Fractional binary numbers

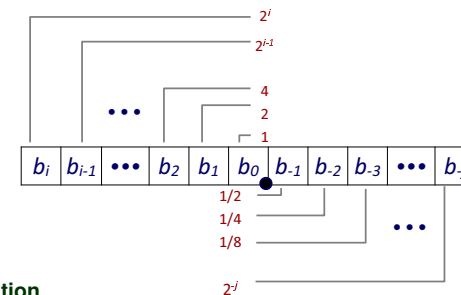
What is 1011.101_2 ?



- 3 -

CS 105

Fractional Binary Numbers



Representation

- Bits to right of "binary point" represent fractional powers of 2

- Represents rational number: $\sum_{k=-j}^i b_k \times 2^k$

- 4 -

CS 105

Fractional Binary Numbers: Examples



Value	Representation
$5 \frac{3}{4}$	101.11_2
$2 \frac{7}{8}$	10.111_2
$1 \frac{7}{16}$	1.0111_2

Observations

- Divide by 2 by shifting right (unsigned)
- Multiply by 2 by shifting left
- Numbers of form $0.11111\dots_2$ are just below 1.0
 - $1/2 + 1/4 + 1/8 + \dots + 1/2^i + \dots \rightarrow 1.0$
 - Use notation $1.0 - \epsilon$

- 5 -

CS 105

Representable Numbers



Limitation #1

- Can only exactly represent numbers of the form $x/2^k$
 - Other rational numbers have repeating bit representations

Value	Representation	Decimal Representation
$1/3$	$0.0101010101[01]\dots_2$	$0.333333333\dots$
$1/5$	$0.001100110011[0011]\dots_2$	$0.200000000\dots$
$1/10$	$0.0001100110011[0011]\dots_2$	$0.100000000\dots$

Limitation #2

- Just one setting of binary point within the w bits
 - Limited range of numbers (very small values? very large?)

- 6 -

CS 105

IEEE Floating Point



IEEE Standard 754

- Established in 1985 as uniform standard for floating-point arithmetic
 - Before that, many idiosyncratic formats
- Supported by all major CPUs

Driven by numerical concerns

- Nice standards for rounding, overflow, underflow
- Hard to make go fast
 - Numerical analysts predominated over hardware types in defining standard

- 7 -

CS 105

Floating-Point Representation



Numerical Form

- $-1^s M 2^E$
 - Sign bit s determines whether number is negative or positive (negative zero representable)
 - Significand M normally a fractional value in range $[1.0, 2.0)$.
 - Exponent E weights value by a power of two

Encoding



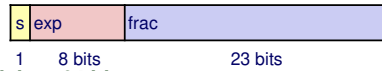
- MSB is sign bit
- exp field encodes E (emphasis on "encodes")
- frac field encodes M (likewise)

- 8 -

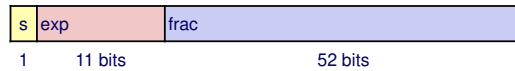
CS 105

Precision Options (Not to Scale)

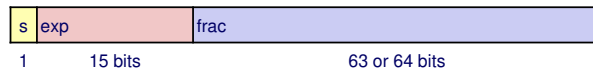
Single precision: 32 bits



Double precision: 64 bits



Extended precision: 80 bits (Intel only)



- 9 -

CS 105

“Normalized” Values

$$v = (-1)^s M 2^E$$

When: $\text{exp} \neq 000\dots 0$ and $\text{exp} \neq 111\dots 1$

Exponent coded as a *biased* value: $E = \text{Exp} - \text{Bias}$

- Exp : unsigned value of exp field
- $\text{Bias} = 2^{k-1} - 1$, where k is number of exponent bits
 - Single precision: 127 (Exp: 1...254, E: -126...127)
 - Double precision: 1023 (Exp: 1...2046, E: -1022...1023)

Significand coded with implied leading 1: $M = 1.\text{xxx}\dots\text{x}_2$

- $\text{xxx}\dots\text{x}$: bits of frac field
- Minimum when $\text{frac} = 000\dots 0$ ($M = 1.0$)
- Maximum when $\text{frac} = 111\dots 1$ ($M = 2.0 - \epsilon$)
- Get extra leading bit for “free”

- 10 -

CS 105

Normalized Encoding Example

Value

`float f = 15213.0;`

■ $15213_{10} = 11101101101101_2 = 1.1101101101101_2 \times 2^{13}$

Significand

$M = 1.1101101101101_2$
 $\text{frac} = 11011011011010000000000_2$

Exponent

$E = 13$
 $\text{Bias} = 127$
 $\text{Exp} = 140 = 10001100_2$

Floating-Point Representation (Class 02):

Hex: 4 6 6 D B 4 0 0
 Binary: 0100 0110 0110 1101 1011 0100 0000 0000
 140: 100 0110 0
 15213: 1110 1101 1011 01

- 11 -

CS 105

Denormalized Values

$$v = (-1)^s M 2^E$$

$$E = 1 - \text{Bias}$$

Condition: $\text{exp} = 000\dots 0$

Exponent value: $E = 1 - \text{Bias}$ (instead of $E = 0 - \text{Bias}$)

Significand coded with implied leading 0: $M = 0.\text{xxx}\dots\text{x}_2$

- $\text{xxx}\dots\text{x}$: bits of frac
- Cases
- $\text{exp} = 000\dots 0$, $\text{frac} = 000\dots 0$
 - Represents zero value
 - Note distinct values: +0 and -0 (why?)
 - $\text{exp} = 000\dots 0$, $\text{frac} \neq 000\dots 0$
 - Numbers closest to 0.0
 - Equispaced

- 12 -

CS 105

Special Values



Condition: $\text{exp} = 111\dots 1$

Case: $\text{exp} = 111\dots 1$, $\text{frac} = 000\dots 0$

- Represents value ∞ (infinity)
- Operation that overflows
- Both positive and negative
- E.g., $1.0/0.0 = -1.0/-0.0 = +\infty$, $1.0/-0.0 = -\infty$

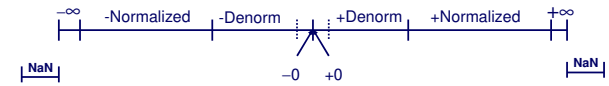
Case: $\text{exp} = 111\dots 1$, $\text{frac} \neq 000\dots 0$

- Not-a-Number (NaN)
- Represents case when no numeric value can be determined
- E.g., $\text{sqrt}(-1)$, $\infty - \infty$, $\infty \times 0$

- 13 -

CS 105

Visualization: Floating-Point Encodings



- 14 -

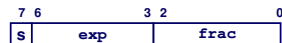
CS 105

Tiny Floating-Point Example



8-bit floating-point representation

- The sign bit is in the most significant bit.
- The next four bits are the exponent, with a bias of 7.
- The last three bits are the frac
- Same general form as IEEE format
 - Normalized, denormalized
 - Representation of 0, NaN, infinity



- 15 -

CS 105

Values Related to the Exponent



Exp	exp	E	2^E	
0	0000	-6	1/64	(denorms)
1	0001	-6	1/64	
2	0010	-5	1/32	
3	0011	-4	1/16	
4	0100	-3	1/8	
5	0101	-2	1/4	
6	0110	-1	1/2	
7	0111	0	1	
8	1000	+1	2	
9	1001	+2	4	
10	1010	+3	8	
11	1011	+4	16	
12	1100	+5	32	
13	1101	+6	64	
14	1110	+7	128	
15	1111	n/a		(inf, NaN)

- 16 -

CS 105

Dynamic Range

$$v = (-1)^s M 2^E$$

$$n: E = \text{Exp} - \text{Bias}$$

$$d: E = 1 - \text{Bias}$$

	s	exp	frac	E	Value
Denormalized numbers	0	0000	000	-6	0
	0	0000	001	-6	$1/8 * 1/64 = 1/512$ ← closest to zero
	0	0000	010	-6	$2/8 * 1/64 = 2/512$
	...				
	0	0000	110	-6	$6/8 * 1/64 = 6/512$
	0	0000	111	-6	$7/8 * 1/64 = 7/512$ ← largest denorm
Normalized numbers	0	0001	000	-6	$8/8 * 1/64 = 8/512$ ← smallest norm
	0	0001	001	-6	$9/8 * 1/64 = 9/512$
	...				
	0	0110	110	-1	$14/8 * 1/2 = 14/16$
	0	0110	111	-1	$15/8 * 1/2 = 15/16$ ← closest to 1 below
	0	0111	000	0	$8/8 * 1 = 1$
	0	0111	001	0	$9/8 * 1 = 9/8$ ← closest to 1 above
	0	0111	010	0	$10/8 * 1 = 10/8$
	...				
	0	1110	110	7	$14/8 * 128 = 224$
	0	1110	111	7	$15/8 * 128 = 240$ ← largest norm
	0	1111	000	n/a	inf

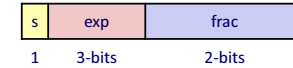
- 17 -

CS 105

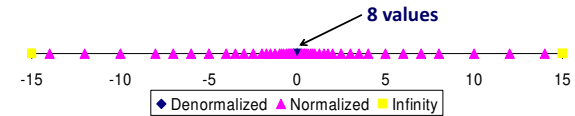
Distribution of Values

6-bit IEEE-like format

- e = 3 exponent bits
- f = 2 fraction bits
- Bias is $2^{3-1} - 1 = 3$



Notice how the distribution gets denser toward zero.



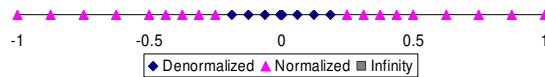
- 18 -

CS 105

Distribution of Values (close-up view)

6-bit IEEE-like format

- e = 3 exponent bits
- f = 2 fraction bits
- Bias is 3



- 19 -

CS 105

Interesting Numbers

Description	exp	frac	Numeric Value
Zero	00...00	00...00	0.0
Smallest Pos. Denorm.	00...00	00...01	$2^{-(23,52)} \times 2^{-(126,1022)}$
			■ Single (float) $\approx 1.4 \times 10^{-45}$
			■ Double $\approx 4.9 \times 10^{-324}$
Largest Denormalized	00...00	11...11	$(1.0 - \epsilon) \times 2^{-(126,1022)}$
			■ Single (float) $\approx 1.18 \times 10^{-38}$
			■ Double $\approx 2.2 \times 10^{-308}$
Smallest Pos. Normalized	00...01	00...00	$1.0 \times 2^{-(126,1022)}$
			■ Just larger than largest denormalized
One	01...11	00...00	1.0
Largest Normalized	11...10	11...11	$(2.0 - \epsilon) \times 2^{(127,1023)}$
			■ Single (float) $\approx 3.4 \times 10^{38}$
			■ Double $\approx 1.8 \times 10^{308}$

- 20 -

CS 105

Special Properties of Encoding



FP zero same as integer zero

- All bits = 0

Can (almost) use unsigned integer comparison

- Must first compare sign bits
- Must consider $-0 = 0$
- NaNs problematic
 - Will be greater than any other values
 - What should comparison yield?
- Otherwise OK
 - Denormalized vs. normalized
 - Normalized vs. infinity

– 21 –

CS 105

Floating Point Operations: Basic Idea



$$x +_f y = \text{Round}(x + y)$$

$$x \times_f y = \text{Round}(x \times y)$$

Basic idea

- First **compute exact result**
- Make it fit into desired precision
 - Possibly overflow if exponent too large
 - Possibly **round to fit into *frac***

– 22 –

CS 105

Rounding



Rounding Modes (illustrated with \$ rounding)

	\$1.40	\$1.60	\$1.50	\$2.50	–\$1.50
■ Towards zero	\$1	\$1	\$1	\$2	–\$1
■ Round down ($-\infty$)	\$1	\$1	\$1	\$2	–\$2
■ Round up ($+\infty$)	\$2	\$2	\$2	\$3	–\$1
■ Nearest Even (default)	\$1	\$2	\$2	\$2	–\$2

– 23 –

CS 105

Closer Look at Round-To-Even



Default rounding mode

- Hard to get any other kind without dropping into assembly
- All others are statistically biased
 - Sum of set of positive numbers will consistently be over- or under-estimated
 - Need randomness

Applying to other decimal places / bit positions

- When exactly halfway between two possible values:
 - Round so that least significant digit is even
- E.g., round to nearest hundredth

1.2349999	1.23	(Less than half way)
1.2350001	1.24	(Greater than half way)
1.2350000	1.24	(Half way—round up)
1.2450000	1.24	(Half way—round down)

– 24 –

CS 105

Rounding Binary Numbers

Binary fractional numbers

- “Even” when least significant bit is 0
- Halfway when bits to right of rounding position = 100...₂

Examples

- Round to nearest 1/4 (2 bits right of binary point)

Value	Binary	Rounded	Action	Rounded Value
2 3/32	10.00011 ₂	10.00 ₂	(<1/2—down)	2
2 3/16	10.00110 ₂	10.01 ₂	(>1/2—up)	2 1/4
2 7/8	10.11100 ₂	11.00 ₂	(1/2—up)	3
2 5/8	10.10100 ₂	10.10 ₂	(1/2—down)	2 1/2

– 25 –

CS 105

FP Multiplication

Operands

$$(-1)^{s1} M1 2^{E1} * (-1)^{s2} M2 2^{E2}$$

Exact Result

$$(-1)^s M 2^E$$

- Sign s : $s1 \wedge s2$
- Significand M : $M1 * M2$
- Exponent E : $E1 + E2$

Fixing

- If $M \geq 2$, shift M right, increment E
- If E out of range, overflow
- Round M to fit ϵ_{rac} precision

Implementation

- 26 – ■ Biggest chore is multiplying significands

CS 105

FP Addition

Operands

$$(-1)^{s1} M1 2^{E1}$$

$$(-1)^{s2} M2 2^{E2}$$

- Assume $E1 > E2$

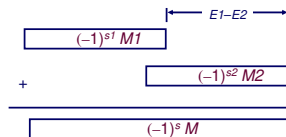
Exact Result

$$(-1)^s M 2^E$$

- Sign s , significand M :
 - Result of signed align & add
- Exponent E : $E1$

Fixing

- If $M \geq 2$, shift M right, increment E
- If $M < 1$, shift M left k positions, decrement E by k
- Overflow if E out of range
- 27 – ■ Round M to fit ϵ_{rac} precision



– 27 –

CS 105

Mathematical Properties of FP Add

Compare to those of Abelian Group

- Closed under addition? **Yes**
 - But may generate infinity or NaN
- Commutative? **Yes**
- Associative? **No**
 - Overflow and inexactness of rounding
 - $(3.14 + 1e10) - 1e10 = 0$, $3.14 + (1e10 - 1e10) = 3.14$
- 0 is additive identity? **Yes**
- Every element has additive inverse? **Almost**
 - Yes, except for infinities & NaNs

Monotonicity

- $a \geq b \Rightarrow a+c \geq b+c$? **Almost**
 - Except for infinities & NaNs

– 28 –

CS 105

Mathematical Properties of FP Mult



Compare to Commutative Ring

- Closed under multiplication? Yes
 - But may generate infinity or NaN
- Multiplication Commutative? Yes
- Multiplication is Associative? No
 - Possibility of overflow, inexactness of rounding
 - Ex: $(1e20 * 1e20) * 1e-20 = \text{inf}$, $1e20 * (1e20 * 1e-20) = 1e20$
- 1 is multiplicative identity? Yes
- Multiplication distributes over addition? No
 - Possibility of overflow, inexactness of rounding
 - $1e20 * (1e20 - 1e20) = 0.0$, $1e20 * 1e20 - 1e20 * 1e20 = \text{NaN}$

Monotonicity

- $a \geq b$ & $c \geq 0 \Rightarrow a * c \geq b * c$? Almost
 - Except for infinities & NaNs

- 29 -

CS 105

Floating Point in C



C Guarantees Two Levels

float single precision
double double precision

Conversions

- Casting between int, float, and double changes numeric values
- Double Or float to int
 - Truncates fractional part
 - Like rounding toward zero
 - Not defined when out of range
 - » Generally saturates to TMin or TMax
- int to double
 - Exact conversion, as long as int has ≤ 53 -bit word size
- int to float
 - Will round according to rounding mode

- 30 -

CS 105

Answers to Floating-Point Puzzles



```
int x = ...;
float f = ...;
double d = ...;
```

Assume neither
d nor f is NaN

- | | |
|---|---------------------------|
| • <code>x == (int)(float) x</code> | No: 24 bit significand |
| • <code>x == (int)(double) x</code> | Yes: 53 bit significand |
| • <code>f == (float)(double) f</code> | Yes: increases precision |
| • <code>d == (float) d</code> | No: loses precision |
| • <code>f == -(-f)</code> | Yes: Just change sign bit |
| • <code>2/3 == 2/3.0</code> | No: <code>2/3 == 0</code> |
| • <code>d < 0.0 \Rightarrow ((d*2) < 0.0)</code> | Yes! |
| • <code>d > f \Rightarrow -f > -d</code> | Yes! |
| • <code>d * d >= 0.0</code> | Yes! |
| • <code>(d+f)-d == f</code> | No: Not associative |

- 31 -

CS 105

Ariane 5



- Exploded 37 seconds after liftoff
- Cargo worth \$500 million

Why

- Computed horizontal velocity as floating-point number
- Converted to 16-bit integer
- Worked OK for Ariane 4
- Overflowed for Ariane 5
 - Used same software



- 32 -

CS 105

Summary



IEEE floating point has clear mathematical properties

- Represents numbers of form $M \times 2^E$
- Can reason about operations independent of implementation
 - As if computed with perfect precision and then rounded
- Not the same as real arithmetic
 - Violates associativity/distributivity
 - Makes life difficult for compilers & serious numerical applications programmers