Introduction to Computer Systems

15-213/18-243, fall 2009 3rd Lecture, Sep. 1st

Instructors:

Roger B. Dannenberg and Greg Ganger

Last Time: Integers

- Representation: unsigned and signed
- Conversion, casting
 - Bit representation maintained but reinterpreted
- Expanding, truncating
 - Truncating = mod
- Addition, negation, multiplication, shifting
 - Operations are mod 2^w
- Ordering properties do not hold
 - u > 0 does not mean u + v > v
 - u, v > 0 does not mean u · v > 0

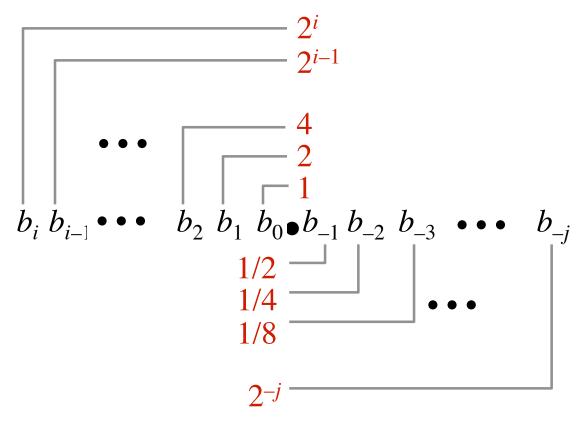
Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

Fractional binary numbers

- What is 1023.405₁₀?
- What is 1011.101₂?

Fractional Binary Numbers



Representation

- Bits to right of "binary point" represent fractional powers of 2
- Represents rational number: $\sum_{k=-j}^{i} b_k \cdot 2^k$

Fractional Binary Numbers: Examples

Value

Representation

5-3/4 101.11₂
2-7/8 10.111₂
63/64 0.111111₂

Observations

- Divide by 2 by shifting right
- Multiply by 2 by shifting left
- Compare to shifting decimal numbers right or left
- Numbers of form **0.111111**...₂ are just below 1.0
 - $1/2 + 1/4 + 1/8 + ... + 1/2^{i} + ... \rightarrow 1.0$
 - Compare to $0.9999..._{10} \rightarrow 1.0$
 - Use notation 1.0ε

Representable Numbers

Limitation

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

Value	Representation
1/3	$0.01010101[01]{2}$
1/5	$0.001100110011[0011]_{\cdots 2}$
1/10	$0.0001100110011[0011]{2}$

Observation

0.1₁₀ has no finite exact binary representation!

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

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 fast in hardware
 - Numerical analysts predominated over hardware designers in defining standard

Floating Point Representation

Numerical Form:

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

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

Encoding

- MSB s is sign bit s
- exp field encodes E (but is not equal to E)
- frac field encodes M (but is not equal to M)

s	exp	frac
	•	

Precisions

■ Single precision: 32 bits

s	exp	frac
1	8	23

Double precision: 64 bits

s	ехр	frac
1	11	52

Extended precision: 80 bits (Intel only)

s	ехр	frac
1	15	63 or 64

Normalized Values

- Condition: $exp \neq 000...0$ and $exp \neq 111...1$
- **Exponent coded as** biased **value:** E = Exp Bias
 - Exp: unsigned value exp
 - $Bias = 2^{e-1} 1$, where e 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: $M = 1 \cdot xxx...x_2$
 - xxx...x: bits of frac
 - Minimum when 000...0 (M = 1.0)
 - Maximum when **111...1** ($M = 2.0 \varepsilon$)
 - Why does M range from 1 to 2-? Why not 0 to 1-?
 - Get extra leading bit for "free"

Normalized Encoding Example

```
■ Value: Float F = 15213.0;

■ 15213<sub>10</sub> = 11101101101101<sub>2</sub>

= 1.1101101101101<sub>2</sub> x 2<sup>13</sup>
```

Significand

```
M = 1.1101101101_2
frac= 1101101101101_000000000_2
```

Exponent

```
E = 13
Bias = 127
Exp = 140 = 10001100_{2}
```

Result:

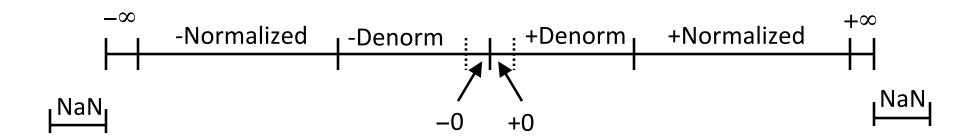
Denormalized Values

- **■** Condition: exp = 000...0
- **Exponent value:** E = 1 Bias (instead of E = 0 Bias)
- Significand coded with implied leading 0: $M = 0 . xxx...x_2$
 - xxx...x: bits of frac
- Cases
 - exp = 000...0, frac = 000...0
 - Represents value 0
 - Note distinct values: +0 and -0 (why?)
 - exp = 000...0, $frac \neq 000...0$
 - Numbers very close to 0.0
 - Lose precision as get smaller
 - Equispaced

Special Values

- Condition: exp = 111...1
- Case: exp = 111...1, frac = 000...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: exp = 111...1, frac ≠ 000...0
 - Not-a-Number (NaN)
 - Represents case when no numeric value can be determined
 - E.g., sqrt(-1), $\infty \infty$, $\infty * 0$

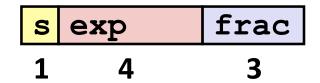
Visualization: Floating Point Encodings



Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

Tiny Floating Point Example



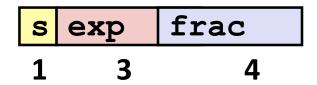
8-bit Floating Point Representation

- 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

Is 8-bit Float Just an Example?



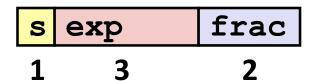
- uLaw Audio Representation
 - An 8-bit float used for digital telephony in North America/Japan
- We'll hear some examples later
- Small floats also used in GPUs!

Dynamic Range (Positive Only)

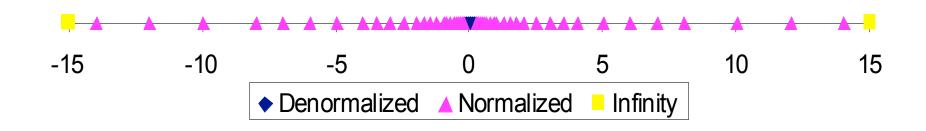
	s exp frac	E	Value
Denormalized numbers	0 0000 000 0 0000 001 0 0000 010	-6 -6 -6	0 1/8*1/64 = 1/512 closest to zero 2/8*1/64 = 2/512
Humbers	0 0000 110 0 0000 111	-6 -6	6/8*1/64 = 6/512 7/8*1/64 = 7/512 largest denorm
	0 0001 000 0 0001 001		8/8*1/64 = 8/512 smallest norm 9/8*1/64 = 9/512
Normalized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

Distribution of Values

- 6-bit IEEE-like format
 - e = 3 exponent bits
 - f = 2 fraction bits
 - Bias is $2^{3-1}-1=3$



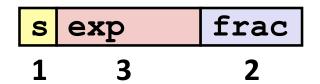
Notice how the distribution gets denser toward zero.

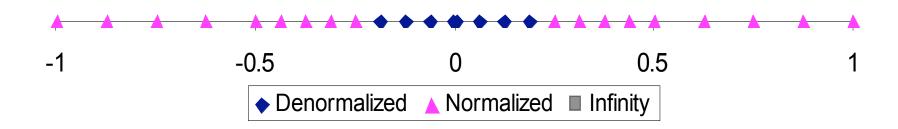


Distribution of Values (close-up view)

6-bit IEEE-like format

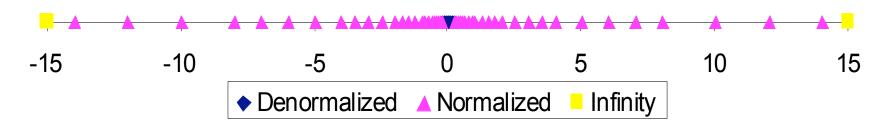
- e = 3 exponent bits
- f = 2 fraction bits
- Bias is 3





Sound Examples

Floats are more precise near zero



Fixed-point numbers quantize uniformly throughout their range



Interesting Numbers

{single,double}

Description	exp	frac	Numeric Value
Zero	0000	0000	0.0
■ Smallest Pos. Denorm.	0000	0001	$2^{-\{23,52\}} \times 2^{-\{126,1022\}}$
■ Single $\approx 1.4 \times 10^{-45}$			
■ Double $\approx 4.9 \times 10^{-32}$	4		
Largest Denormalized	0000	1111	$(1.0 - \varepsilon) \times 2^{-\{126,1022\}}$
■ Single $\approx 1.18 \times 10^{-38}$			
■ Double $\approx 2.2 \times 10^{-30}$	8		
■ Smallest Pos. Normalize	ed	0001	0000 1.0 x $2^{-\{126,1022\}}$
Just larger than larg	est denoi	rmalized	
One	0111	0000	1.0
Largest Normalized	1110	1111	$(2.0 - \varepsilon) \times 2^{\{127,1023\}}$
■ Single $\approx 3.4 \times 10^{38}$			
■ Double $\approx 1.8 \times 10^{308}$			

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
 - Denorm vs. normalized
 - Normalized vs. infinity

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

Floating Point Operations: Basic Idea

$$\mathbf{x} +_{\mathbf{f}} \mathbf{y} = \text{Round}(\mathbf{x} + \mathbf{y})$$

$$\mathbf{x} \times_{\mathbf{f}} \mathbf{y} = \text{Round}(\mathbf{x} \times \mathbf{y})$$

Basic idea

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

Rounding

Rounding Modes (illustrate with \$ rounding)

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

■ What are the advantages of the modes?

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 underestimated

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 wav—round down)

Rounding Binary Numbers

Binary Fractional Numbers

- "Even" when least significant bit is 0
- "Half way" 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	10.000112	10.002	(<1/2—down)	2
2 3/16	10.001102	10.012	(>1/2—up)	2 1/4
2 7/8	10.11 <mark>100</mark> 2	11.002	(1/2—up)	3
2 5/8	10.10 <mark>100</mark> 2	10.102	(1/2—down)	2 1/2

FP Multiplication

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

■ Exact Result: $(-1)^s M 2^E$

• Sign s: s1 ^ s2

■ Significand *M*: *M1* * *M2*

• Exponent *E*: *E*1 + *E*2

Fixing

- If $M \ge 2$, shift M right, increment E
- If E out of range, overflow
- Round M to fit frac precision

Implementation

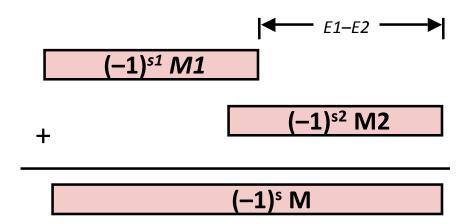
Biggest chore is multiplying significands

Floating Point Addition

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

Assume E1 > E2

- Exact Result: (-1)^s M 2^E
 - Sign s, significand M:
 - Result of signed align & add
 - Exponent E:



Fixing

- If $M \ge 2$, shift M right, increment E
- if *M* < 1, shift *M* left *k* positions, decrement *E* by *k*
- Overflow if E out of range
- Round M to fit frac precision

Mathematical Properties of FP Add

Compare to those of Abelian Group

Closed under addition?
Yes

But may generate infinity or NaN

Commutative?

Associative?

Overflow and inexactness of rounding

• 0 is additive identity?

Every element has additive inverse
Almost

Except for infinities & NaNs

Monotonicity

• $a \ge b \Rightarrow a+c \ge b+c$?

Except for infinities & NaNs

Mathematical Properties of FP Mult

Compare to Commutative Ring

Closed under multiplication?
Yes

But may generate infinity or NaN

• Multiplication Commutative?

• Multiplication is Associative?

Possibility of overflow, inexactness of rounding

■ 1 is multiplicative identity? Yes

• Multiplication distributes over addition?

Possibility of overflow, inexactness of rounding

Monotonicity

$$\bullet$$
 $a \ge b \& c \ge 0 \Rightarrow a *c \ge b *c?$

Almost

Except for infinities & NaNs

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

Floating Point in C

C Guarantees Two Levels

float single precision
double double precision

Conversions/Casting

- Casting between int, float, and double changes bit representation
- Double/float → int
 - Truncates fractional part
 - Like rounding toward zero
 - Not defined when out of range or NaN: Generally sets to TMin
- int → double
 - Exact conversion, as long as int has \leq 53 bit word size
- int → float
 - Will round according to rounding mode

Floating Point Puzzles

- For each of the following C expressions, either:
 - Argue that it is true for all argument values
 - Explain why not true

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

Assume neither d nor f is NaN

```
• x == (int)(float) x
• x == (int) (double) x
• f == (float)(double) f
• d == (float) d
• f == -(-f);
\cdot 2/3 == 2/3.0
• d < 0.0 \Rightarrow ((d*2) < 0.0)
• d > f \Rightarrow -f > -d
• d * d >= 0.0
• (d+f)-d == f
```

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

Summary

- IEEE Floating Point has clear mathematical properties
- **Represents numbers of form** $M \times 2^{E}$
- One 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

More Slides

Creating Floating Point Number

Steps

Normalize to have leading 1

S	ехр	frac	
7	6	3 2	0

- Round to fit within fraction
- Postnormalize to deal with effects of rounding

Case Study

- Convert 8-bit unsigned numbers to tiny floating point format
- Example Numbers

128	10000000
15	00001101
33	00010001
35	00010011
138	10001010
63	00111111

Normalize

7 6	3 2	0
s exp	frac	

Requirement

- Set binary point so that numbers of form 1.xxxxx
- Adjust all to have leading one
 - Decrement exponent as shift left

Value	Binary Fraction Exponent			
128	10000000	1.000000	7	
15	00001101	1.1010000	3	
17	00010001	1.0001000	5	
19	00010011	1.0011000	5	
138	10001010	1.0001010	7	
63	00111111	1.1111100	5	

Rounding

1.BBGRXXX

Guard bit: LSB of result

Round bit: 1st bit removed

Sticky bit: OR of remaining bits

Round up conditions

- Round = 1, Sticky = $1 \rightarrow > 0.5$
- Guard = 1, Round = 1, Sticky = 0 → Round to even

Value	Fraction	GRS	Incr?	Rounded
128	1.000000	00*	N	1.000
15	1.1010000	10*	N	1.101
17	1.0001000	010	N	1.000
19	1.0011000	110	Y	1.010
138	1.0001010	011	Y	1.001
63	1.1111100	111	Y	10.000

Postnormalize

Issue

- Rounding may have caused overflow
- Handle by shifting right once & incrementing exponent

Value	Rounded	Exp	Adjusted	Result
128	1.000	7		128
15	1.101	3		15
17	1.000	4		16
19	1.010	4		20
138	1.001	7		134
63	10.000	5	1.000/6	64