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

Instructors:
Roger B. Dannenberg and Greg Ganger

■ 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: Examples ■ Value Representation 101.112 5-3/4 2-7/8 10.1112 63/64 0.111111, 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 • $1/2 + 1/4 + 1/8 + ... + 1/2^{i} + ... \rightarrow 1.0$ \bullet Compare to 0.9999...₁₀ \rightarrow 1.0 • Use notation 1.0 – ϵ

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 0.0101010101[011...2
1/5 0.001100110011[0011]...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 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)* M 2*

Sign bits 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)

Precisions

Single precision: 32 bits

Sexp frac

1 8 23

Double precision: 64 bits

Sexp frac

1 11 52

Extended precision: 80 bits (Intel only)

Sexp frac

1 15 63 or 64

Normalized Values

Condition: exp ≠ 000...0 and exp ≠ 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 . xxx...x

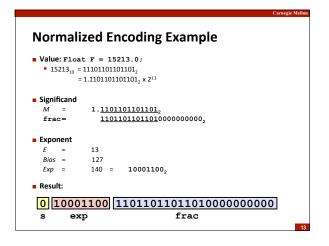
xxx...x: bits of £rac

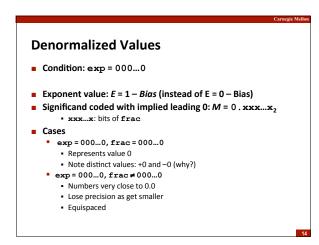
Minimum when 000...0 (M = 1.0)

Maximum when 111...1 (M = 2.0 - ε)

Why does M range from 1 to 2-? Why not 0 to 1-?

Get extra leading bit for "free"





Special Values

Condition: exp = 111...1

Case: exp = 111...1, frac = 000...0

Represents value ∞ (infinity)

Operation that overflows

Both positive and negative

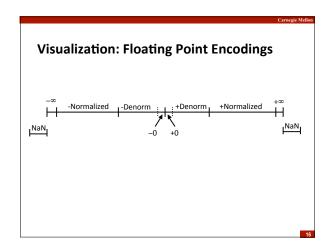
E.g., 1.0/0.0 = -1.0/-0.0 = +∞, 1.0/-0.0 = -∞

Case: exp = 111...1, frac ≠ 000...0

Not-a-Number (NaN)

Represents case when no numeric value can be determined

E.g., sqrt(-1), ∞ - ∞, ∞ * 0



Today: Floating Point

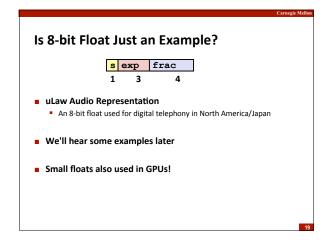
Background: Fractional binary numbers
IEEE floating point standard: Definition
Example and properties
Rounding, addition, multiplication
Floating point in C
Summary

Tiny Floating Point Example

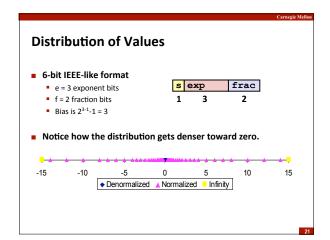
Sexp frac
1 4 3

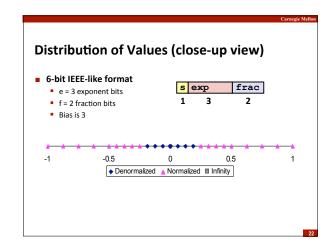
8-bit Floating Point Representation
the sign bit is in the most significant bit.
the next four bits are the exponent, with a bias of 7.
the last three bits are the frac

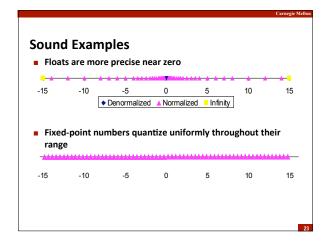
Same general form as IEEE Format
normalized, denormalized
representation of 0, NaN, infinity

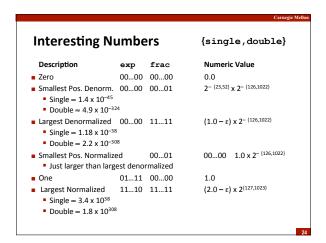


		. /D ''	0 . 1 .)	Carnegie Mello			
Dynamic Range (Positive Only)							
	s exp fr	ac E	Value				
Denormalized numbers	0 0000 00	00 -6	0				
	0 0000 00	1 -6	1/8*1/64 = 1/512	closest to zero			
	0 0000 01	.0 -6	2/8*1/64 = 2/512				
	0 0000 11	.0 -6	6/8*1/64 = 6/512				
	0 0000 11	1 -6	7/8*1/64 = 7/512	largest denorm			
Normalized numbers	0 0001 00	0 -6	8/8*1/64 = 8/512	smallest norm			
	0 0001 00	1 -6	9/8*1/64 = 9/512				
	0 0110 11	.0 -1	,,,				
	0 0110 11		15/8*1/2 = 15/16	closest to 1 below			
	0 0111 00						
	0 0111 00		9/8*1 = 9/8	closest to 1 above			
	0 0111 01	.0 0	10/8*1 = 10/8				
	0 1110 11	0 7	14/8*128 = 224				
	0 1110 11	1 7	15/8*128 = 240	largest norm			
	0 1111 00	00 n/a	inf				
				2			









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

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

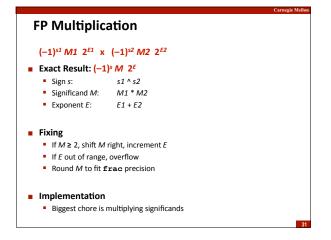
- Default Rounding Mode
 - Hard to get any other kind without dropping into assembly
 - All others are statistically biased
 - Sum of set of positive numbers will consistently be over- or underestimated
- Applying to Other Decimal Places / Bit Positions
 - When exactly halfway between two possible values
 - Round so that least significant digit is even
 - E.g., round to nearest hundredth

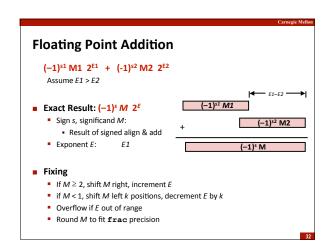
1.2349999 1.23 (Less than half way) 1.2350001 1.24 (Greater than half way) 1.2350000 1.24 (Half way—round up) 1.2450000 1.24 (Half way—round down)

Rounding Binary Numbers

- Binary Fractional Numbers
 - "Even" when least significant bit is 0
 - \blacksquare "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 23/32 10.00011₂ 10.00₂ (<1/2-down) 23/16 10.00110₂ 10.01₂ (>1/2-up) 27/8 10.11100₂ 11.00₂ (1/2-up) 2 1/4 3 25/8 10.10100₂ 10.10₂ (1/2-down) 2 1/2





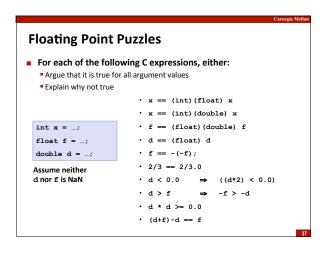
Mathematical Properties of FP Add ■ Compare to those of Abelian Group Closed under addition? Yes · But may generate infinity or NaN Yes Commutative? Associative? No • Overflow and inexactness of rounding 0 is additive identity? Yes Every element has additive inverse **Almost** Except for infinities & NaNs Monotonicity **Almost** • $a \ge b \Rightarrow a+c \ge b+c$? • Except for infinities & NaNs

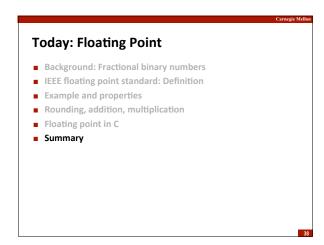
Mathematical Properties of FP Mult ■ Compare to Commutative Ring Closed under multiplication? Yes • But may generate infinity or NaN Multiplication Commutative? Yes • Multiplication is Associative? • Possibility of overflow, inexactness of rounding 1 is multiplicative identity? • Multiplication distributes over addition? No • Possibility of overflow, inexactness of rounding Monotonicity • $a \ge b \& c \ge 0 \Rightarrow a * c \ge b * c$? **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 ≤ 53 bit word size int → float Will round according to rounding mode



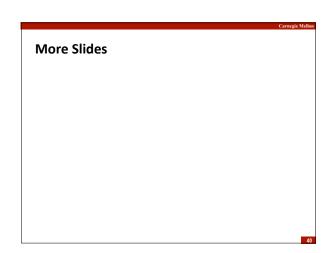


Summary

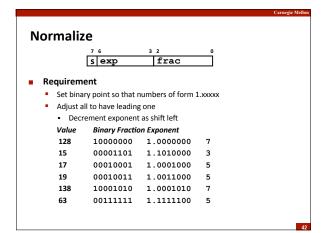
IEEE Floating Point has clear mathematical properties
Represents numbers of form $M \times 2^E$ One can reason about operations independent of implementation

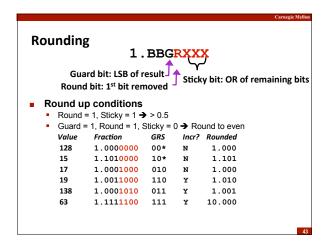
As if computed with perfect precision and then rounded

Not the same as real arithmetic
Violates associativity/distributivity
Makes life difficult for compilers & serious numerical applications programmers



Creating Floating Point Number Steps s exp frac Normalize to have leading 1 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 10000000 128 15 00001101 33 00010001 35 00010011 138 10001010 00111111





Postnormalize Issue Rounding may have caused overflow ■ Handle by shifting right once & incrementing exponent Value Rounded Exp Adjusted 128 1.000 128 7 1.000 7 128 1.101 3 15 1.000 4 16 1.010 4 20 1.001 7 134 10.000 5 1.000/6 64 15 17 19 138 134 63