

## 15-213

*"The course that gives CMU its Zip!"*

# Floating Point

Sept 5, 2007

**Topics**

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

lecture-03.ppt 15-213, F07

## Fractional Binary Numbers

**Representation**

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

- 3 - 15-213, F07

## Frac. Binary Number Examples

Value	Representation
5-3/4	101.11 <sub>2</sub>
2-7/8	10.111 <sub>2</sub>
63/64	0.111111 <sub>2</sub>

**Observations**

- Divide by 2 by shifting right
- Multiply by 2 by shifting left
- Numbers of form 0.111111...<sub>2</sub> just below 1.0
  - $1/2 + 1/4 + 1/8 + \dots + 1/2^j + \dots \rightarrow 1.0$
  - Use notation  $1.0 - \epsilon$

- 4 - 15-213, F07

## Representable Numbers

**Limitation**

- Can only exactly represent numbers of the form  $x/2^k$ 
  - and limited based on number of bits
- Other numbers have repeating bit representations

Value	Representation
1/3	0.0101010101[01]... <sub>2</sub>
1/5	0.001100110011[0011]... <sub>2</sub>
1/10	0.0001100110011[0011]... <sub>2</sub>

- 5 - 15-213, F07

## 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 won over hardware types in defining standard

- 6 - 15-213, F07

## 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 is sign bit
- exp field encodes  $E$
- frac field encodes  $M$

- 7 - 15-213, F07

## Floating Point Precisions

### Encoding



- MSB is sign bit
- exp field encodes  $E$
- frac field encodes  $M$

### Sizes

- Single precision: 8 exp bits, 23 frac bits
  - 32 bits total
- Double precision: 11 exp bits, 52 frac bits
  - 64 bits total
- Extended precision: 15 exp bits, 63 frac bits
  - Only found in Intel-compatible machines
  - Stored in 80 bits
  - » 1 bit wasted

- 8 -

15-213, F07

## “Normalized” Numeric Values

### Condition

- exp  $\neq$  000...0 and exp  $\neq$  111...1

### Exponent coded as *biased* value

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

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

### Significand coded with implied leading 1

$$M = 1.\text{xxx}\dots\text{x}_2$$

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

- 9 -

15-213, F07

## Normalized Encoding Example

### Value

Float  $F = 15213.0$ ;

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

### Significand

$M = 1.1101101101101_2$

frac = 11011011011010000000000<sub>2</sub>

### Exponent

$E = 13$

Bias = 127

Exp = 140 = 10001100<sub>2</sub>

### Floating Point Representation:

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

- 10 -

15-213, F07

## Denormalized Values

### Condition

- exp = 000...0

### Value

- Exponent value  $E = -\text{Bias} + 1$
- Significand value  $M = 0.\text{xxx}\dots\text{x}_2$ 
  - xxx...x: bits of frac

### Cases

- exp = 000...0, frac = 000...0
  - Represents value 0
  - Note that have distinct values +0 and -0
- exp = 000...0, frac  $\neq$  000...0
  - Numbers closest to 0.0, evenly spaced
  - “Gradual underflow”

- 11 -

15-213, F07

## Special Values

### Condition

- exp = 111...1

### Cases

- exp = 111...1, frac = 000...0
  - Represents value  $\infty$  (infinity)
  - Operation that overflows
  - Both positive and negative
  - E.g.,  $1.0/0.0 = -1.0/-0.0 = +\infty$ ,  $1.0/-0.0 = -\infty$
- exp = 111...1, frac  $\neq$  000...0
  - Not-a-Number (NaN)
  - Represents case when no numeric value can be determined
  - E.g.,  $\text{sqrt}(-1)$ ,  $\infty - \infty$ ,  $\infty * 0$

- 12 -

15-213, F07

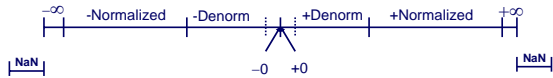
## Interesting Numbers

Description	exp	frac	Numeric Value
Zero	00...00	00...00	0.0
Smallest Pos. Denorm.	00...00	00...01	$2^{-(23,52)} \times 2^{-(126,1022)}$ <ul style="list-style-type: none"> <li>■ Single <math>\approx 1.4 \times 10^{-45}</math></li> <li>■ Double <math>\approx 4.9 \times 10^{-324}</math></li> </ul>
Largest Denormalized	00...00	11...11	$(1.0 - \epsilon) \times 2^{-(126,1022)}$ <ul style="list-style-type: none"> <li>■ Single <math>\approx 1.18 \times 10^{-38}</math></li> <li>■ Double <math>\approx 2.2 \times 10^{-308}</math></li> </ul>
Smallest Pos. Normalized	00...01	00...00	$1.0 \times 2^{-(126,1022)}$ <ul style="list-style-type: none"> <li>■ Just larger than largest denormalized</li> </ul>
One	01...11	00...00	1.0
Largest Normalized	11...10	11...11	$(2.0 - \epsilon) \times 2^{(127,1023)}$ <ul style="list-style-type: none"> <li>■ Single <math>\approx 3.4 \times 10^{38}</math></li> <li>■ Double <math>\approx 1.8 \times 10^{308}</math></li> </ul>

- 13 -

15-213, F07

## Summary of Floating Point Real Number Encodings



- 14 -

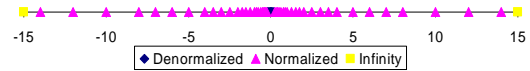
15-213, F07

## Distribution of Values

### 6-bit IEEE-like format

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

Notice how the distribution gets denser toward zero.



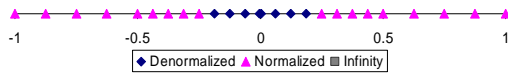
- 15 -

15-213, F07

## Distribution of Values (close-up view)

### 6-bit IEEE-like format

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



- 16 -

15-213, F07

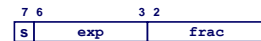
## A Tiny Floating Point Example

### 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



- 17 -

15-213, F07

## Values Related to the Exponent

Exp	exp	E	2 <sup>E</sup>	
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)

- 18 -

15-213, F07

## Dynamic Range

	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	

- 19 -

15-213, F07

## Special Properties of IEEE 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

- 20 -

15-213, F07

## Floating Point Operations

### Conceptual View

- First compute exact result
- Make it fit into desired precision
  - Possibly overflow if exponent too large
  - Possibly round to fit into  $\epsilon_{\text{rac}}$

### Rounding Modes (illustrate with \$ rounding)

	\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
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

Note:

- Round down: rounded result is close to but no greater than true result.
- Round up: rounded result is close to but no less than true result.

- 21 -

15-213, F07

## 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)

- 22 -

15-213, F07

## Rounding Binary Numbers

### Binary Fractional Numbers

- "Even" when least significant bit is 0
- Half way when bits to right of rounding position =  $100\dots_2$

### 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

- 23 -

15-213, F07

## FP Multiplication

### Operands

$$(-1)^{s_1} M_1 2^{E_1} * (-1)^{s_2} M_2 2^{E_2}$$

### Exact Result

- Sign  $s$ :  $s_1 \wedge s_2$
- Significand  $M$ :  $M_1 * M_2$
- Exponent  $E$ :  $E_1 + E_2$

### Fixing

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

### Implementation

- Biggest chore is multiplying significands

- 24 -

15-213, F07

## FP Addition

### Operands

$$(-1)^{s_1} M_1 2^{E_1} + (-1)^{s_2} M_2 2^{E_2}$$

- Assume  $E_1 > E_2$

### Exact Result

- Sign  $s$ , significand  $M$ :
  - Result of signed align & add
- Exponent  $E$ :  $E_1$

### 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
- Round  $M$  to fit  $\epsilon_{\text{rac}}$  precision

- 25 -

15-213, F07

## Curious Excel Behavior

	Number	Subtract 16	Subtract .3	Subtract .01
Default Format	16.31	0.31	0.01	-1.2681E-15
Currency Format	\$16.31	\$0.31	\$0.01	(\$0.00)

- Spreadsheets use floating point for all computations
- Some imprecision for decimal arithmetic
- Can yield nonintuitive results to an accountant!

- 26 -

15-213, F07

## 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 or NaN
    - Generally sets to TMin
- int to double
  - Exact conversion, as long as int has  $\leq 53$  bit word size
- int to float
  - Will round according to rounding mode

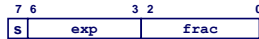
- 27 -

15-213, F07

## Creating Floating Point Number

### Steps

- 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
 

128	10000000
15	00001111
17	00010001
19	00010011
138	10001010
63	00111111

- 28 -

15-213, F07

## Normalize



### 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	1.0000000	7
15	00001111	1.1110000	3
17	00010001	1.0001000	4
19	00010011	1.0011000	4
138	10001010	1.0001010	7
63	00111111	1.1111100	5

- 29 -

15-213, F07

## Rounding

1 . BBGRXXX

Guard bit: LSB of result  
Round bit: 1<sup>st</sup> bit removed  
Sticky bit: OR of remaining bits

### Round up conditions

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

Value	Fraction	GRS	Incr?	Rounded
128	1.0000000	000	N	1.000
15	1.1110000	100	N	1.111
17	1.0001000	010	N	1.000
19	1.0011000	110	Y	1.010
138	1.0001010	111	Y	1.001
63	1.1111100	111	Y	10.000

- 30 -

15-213, F07

## 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	1.111	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

- 31 -

15-213, F07

## 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

- 32 -

15-213, F07

## 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

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

- 33 -

15-213, F07

## 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
- 0 is additive identity? YES
- Every element has additive inverse ALMOST
  - Except for infinities & NaNs

### Monotonicity

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

- 34 -

15-213, F07

## Math. 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
- 1 is multiplicative identity? YES
- Multiplication distributes over addition? NO
  - Possibility of overflow, inexactness of rounding

### Monotonicity

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

- 35 -

15-213, F07