

## CSC 252: Data Representation

### Floating Point Representation

89

89

## This Week's Action Items

- Read Chapter 2 and start reading Chapter 3
- Finish Quiz 3 on Blackboard
- Finish Assignment 1
  - Due Date: ~~Thursday~~ **Friday** September 10 **11** at 11:59 pm

90

90

## Recap: Bits, Bytes, and Integers

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary
- Representations in memory, pointers, strings

91

## 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 <ul style="list-style-type: none"><li>▪ <math>1/2 + 1/4 + 1/8 + \dots + 1/2^i + \dots \rightarrow 1.0</math></li><li>▪ Use notation <math>1.0 - \epsilon</math></li></ul>

92

## Representable Numbers

- Limitation #1
  - Can only exactly represent numbers of the form  $x/2^k$ 
    - Other rational 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>
- Limitation #2
  - Just one setting of binary point within the  $w$  bits
    - Limited range of numbers (very small values? very large?)

93

## Today: Floating Point

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

94

## 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

95

## Floating Point Representation

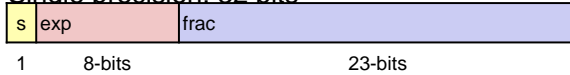
- Numerical Form:  
 $(-1)^s M 2^E$ 
  - Sign bit  $s$  determines whether number is negative or positive
  - Significand  $M$  normally a fractional value in range [1.0,2.0).
  - Exponent  $E$  weights value by power of two
- Encoding
  - MSB  $s$  is sign bit  $s$
  - exp field encodes  $E$  (but is not equal to  $E$ )
  - frac field encodes  $M$  (but is not equal to  $M$ )



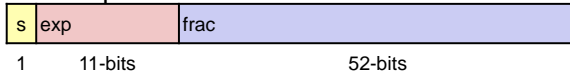
96

## Precision options

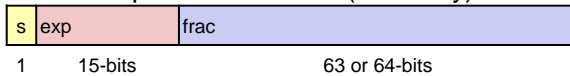
- Single precision: 32 bits



- Double precision: 64 bits



- Extended precision: 80 bits (Intel only)



97

## “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
  - *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:  $M = 1.\text{xxx}\dots\text{x}_2$ 
  - xxx...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”

98

## Normalized Encoding Example $v = (-1)^s M 2^E$ $E = \text{Exp} - \text{Bias}$

- Value: float  $F = 15213.0$ ;
  - $15213_{10} = 11101101101101_2$
  - $= 1.1101101101101_2 \times 2^{13}$

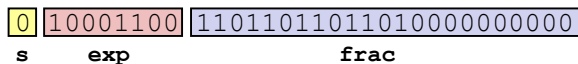
- Significand

$M = 1.\underline{1101101101101}_2$   
 $\text{frac} = \underline{1101101101101000000000}_2$

- Exponent

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

- Result:



99

## 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:  $M = 0.\text{xxx}\dots\text{x}_2$ 
  - xxx...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?)

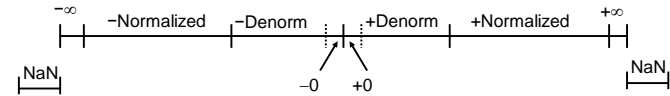
100

## Special Values

- Condition:  $\text{exp} = 111\dots 1$
- Case:  $\text{exp} = 111\dots 1, \text{frac} = 000\dots 0$ 
  - Represents value  $\infty$  (infinity)
  - Operation that overflows
  - Both positive and negative
  - E.g.,  $1.0/0.0 = -1.0/-0.0 = +\infty$ ,  $1.0/-0.0 = -\infty$
- 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.,  $\text{sqrt}(-1)$ ,  $\infty - \infty$ ,  $\infty \times 0$

101

## Visualization: Floating Point Encodings



102

## Special Properties of the 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

103

## Today: Floating Point

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

104

## Floating Point Operations: Basic Idea

- $x \pm_{\epsilon} y = \text{Round}(x + y)$
- $x \times_{\epsilon} y = \text{Round}(x \times y)$
- Basic idea
  - First **compute exact result**
  - Make it fit into desired precision
    - Possibly overflow if exponent too large
    - Possibly **round to fit into frac**

105

## 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 ( $-\infty$ )	\$1	\$1	\$1	\$2	
			–\$2		
– Round up ( $+\infty$ )	\$2	\$2	\$2	\$3	–\$1
– Nearest Even (default)	\$1	\$2	\$2	\$2	
			–\$2		

106

## 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
- 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

107

## 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

Value	Binary	Rounded Value	Action	Rounded Value
2 3/32	10.00 <b>11</b> <sub>2</sub>	10.00 <sub>2</sub>	(<1/2—down)	2
2 3/16	10.00 <b>110</b> <sub>2</sub>	10.01 <sub>2</sub>	(>1/2—up)	2 1/4
2 7/8	10.11 <b>100</b> <sub>2</sub>	11.00 <sub>2</sub>	( 1/2—up)	3
2 5/8	10.10 <b>100</b> <sub>2</sub>	10.10 <sub>2</sub>	( 1/2—down)	2 1/2

108

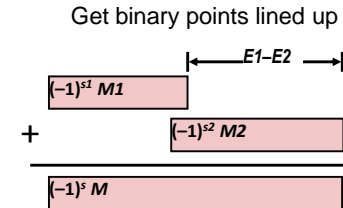
## FP Multiplication

- $(-1)^{s1} M1 2^{E1} \times (-1)^{s2} M2 2^{E2}$
- Exact Result:  $(-1)^s M 2^E$ 
  - Sign  $s$ :  $s1 \wedge s2$
  - Significand  $M$ :  $M1 \times M2$
  - Exponent  $E$ :  $E1 + E2$
- Fixing
  - If  $M \geq 2$ , shift  $M$  right, increment  $E$
  - If  $E$  out of range, overflow
  - Round  $M$  to fit **frac** precision

109

## 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  $E$ :  $E1$
- 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
  - Round  $M$  to fit **frac** precision



110

## 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?
  - Every element has additive inverse? **Yes**
    - Yes, except for infinities & NaNs **Almost**
- Monotonicity
  - $a \geq b \Rightarrow a+c \geq b+c$ ? **Almost**
    - Except for infinities & NaNs

111

## 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

112

## Today: Floating Point

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

113

## 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

114

## 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 == (double)(float) d`
- `f == -(-f)`;
- `2/3 == 2/3.0`
- `d < 0.0` ⇒ `((d*2) < 0.0)`
- `d > f` ⇒ `-f > -d`
- `d * d >= 0.0`
- `(d+f)-d == f`

115

## 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

116