

# Module II

## Floating Point

**Dr. Arijit Roy**

**Computer Science and Engineering Group  
Indian Institute of Information Technology Sri City**

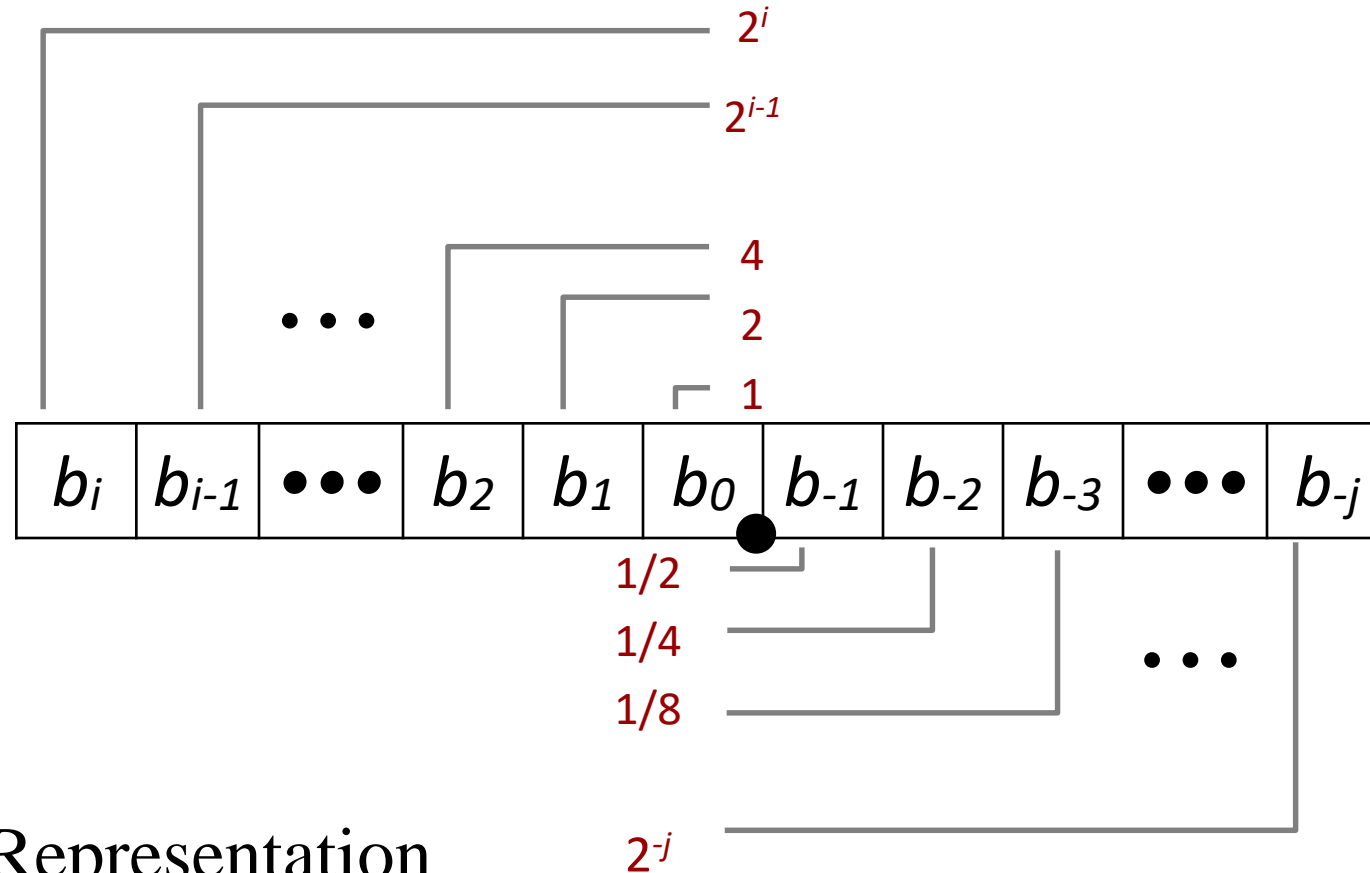
# Floating Point

- 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  $1011.101_2$ ?

# Fractional Binary Numbers



- Representation

- Bits to right of “binary point” represent fractional powers of 2

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

# Fractional Binary Numbers: Examples

■ Value	Representation
5.75	$101.11_2$
2.875	$10.111_2$
0.984375	$1.0111_2$

## ■ Observations

- Divide by 2 by shifting right
- Multiply by 2 by shifting left
- Numbers of form  $0.111111\dots_2$  are just below 1.0
  - $1/2 + 1/4 + 1/8 + \dots + 1/2^i + \dots \rightarrow 1.0$

# Representable Numbers

- Limitation

- Can represent  $x \times 2^y$
- Other rational numbers have repeating bit representations

e.g.,

1110111.10010

In order to arrange in a proper form, we shift it

$1.11011110010 \times 2^6$

• 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>

# Motivation: A real story

On February 25, 1991, during the Gulf War, an American Patriot Missile battery in Dhahran, Saudi Arabia, failed to track and intercept an incoming Iraqi Scud missile. The Scud struck an American Army barracks, killing 28 soldiers and injuring around 100 other people. A report of the General Accounting office, GAO/IMTEC-92-26, entitled Patriot Missile Defense: Software Problem Led to System Failure at Dhahran, Saudi Arabia reported on the cause of the failure. It turns out that the cause was an inaccurate calculation of the time since boot due to computer arithmetic errors.



# Motivation: A real story

- Specifically, the time in tenths of second as measured by the system's internal clock was multiplied by  $1/10$  to produce the time in seconds. This calculation was performed using a 24 bit fixed point register. In particular, the value  $1/10$ , which has a non-terminating binary expansion, was chopped at 24 bits after the radix point. The Patriot battery had been up around 100 hours, and an easy calculation shows that the resulting time error due to the magnified chopping error was about 0.34 seconds. (The number  $1/10$  equals  $1/24 + 1/25 + 1/28 + 1/29 + 1/212 + 1/213 + \dots$ . In other words, the binary expansion of  $1/10$  is  $0.000110011001100110011001100\dots$ . Now the 24 bit register in the Patriot stored instead  $0.00011001100110011001100$  introducing an error of  $0.00000000000000000000000011001100\dots$  binary, or about  $0.0000000095$  decimal. Multiplying by the number of tenths of a second in 100 hours gives  $0.0000000095 \times 100 \times 60 \times 60 \times 10 = 0.34$ .)
- A Scud travels at about 1,676 meters per second, and so travels more than half a kilometer in this time. This was far enough that the incoming Scud was outside the "range gate" that the Patriot tracked. Ironically, the fact that the bad time calculation had been improved in some parts of the code, but not all, contributed to the problem, since it meant that the inaccuracies did not cancel, as discussed here.



# 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

# IEEE Floating-Point Format

single: 8 bits

single: 23 bits

double: 11 bits

double: 52 bits



$$x = (-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$$

Significand

- S: sign bit (0  $\Rightarrow$  non-negative, 1  $\Rightarrow$  negative)
- Normalize significand:  $1.0 \leq |\text{significand}| < 2.0$ 
  - Always has a leading pre-binary-point 1 bit, so no need to represent it explicitly (hidden bit)
  - Significand is Fraction with the “1.” restored
- Exponent: excess representation: actual exponent + Bias
  - Ensures exponent is unsigned
  - Single: Bias = 127; Double: Bias 1023

# Floating-Point Example

- Represent  $-0.75$ 
  - $-0.75 = -0.11_2 \times 2^0 = (-1)^1 \times 1.1_2 \times 2^{-1}$
  - $S = 1$
  - Fraction =  $1000\dots00_2$
  - Exponent =  $-1 + \text{Bias}$ 
    - Single:  $-1 + 127 = 126 = 01111110_2$
    - Double:  $-1 + 1023 = 1022 = 01111111110_2$
- Single:  $1011111101000\dots00$
- Double:  $1011111111101000\dots00$

General representation:

$$(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - 127)}$$

# Reason for bias

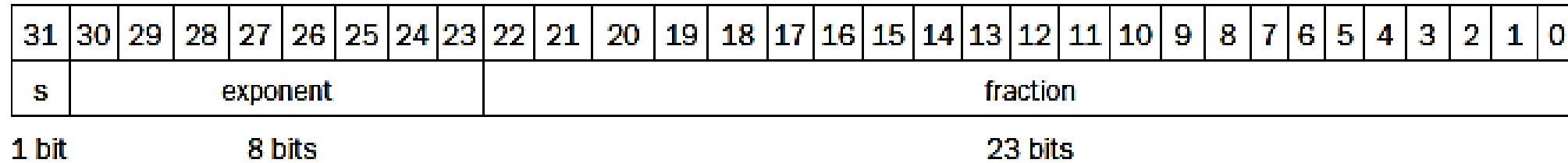
Negative exponents pose a challenge to simplified sorting. If we use two's complement or any other notation in which negative exponents have a 1 in the most significant bit of the exponent field, a negative exponent will look like a big number. For example,  $1.0_{\text{two}} \times 2^{-1}$  would be represented as

[illegible]

(Remember that the leading 1 is implicit in the significand.) The value  $1.0_{\text{two}} \times 2^{+1}$  would look like the smaller binary number

[illegible]

# Single Precision

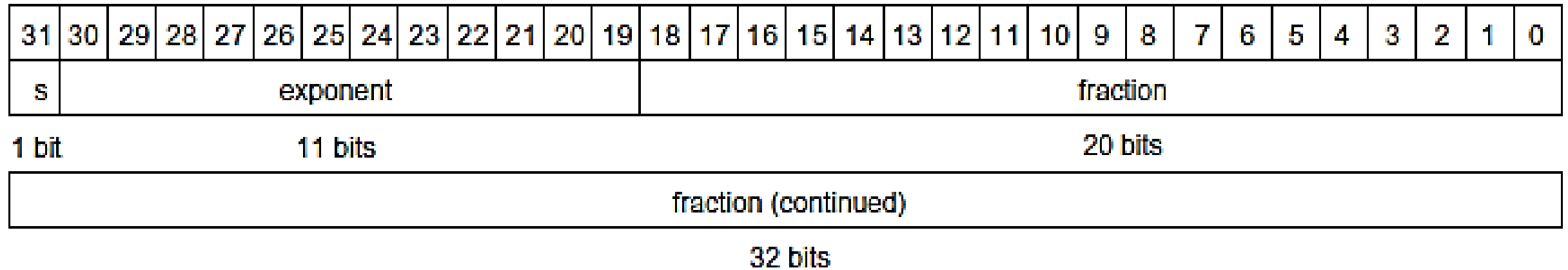


- **Overflow** - the exponent is too large to be represented in the exponent field
- **Underflow** - nonzero fraction has become so small that it cannot be represented

# Single-Precision Range

- Exponents 00000000 and 11111111 reserved
- Smallest value
  - Exponent: 00000001  
 $\Rightarrow$  actual exponent =  $1 - 127 = -126$
  - Fraction: 000...00  $\Rightarrow$  significand = 1.0
  - $\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$
- Largest value
  - exponent: 11111110  
 $\Rightarrow$  actual exponent =  $254 - 127 = +127$
  - Fraction: 111...11  $\Rightarrow$  significand  $\approx 2.0$
  - $\pm 2.0 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$

# Double Precision





# Double-Precision Range

- Exponents 0000...00 and 1111...11 reserved
- Smallest value
  - Exponent: 000000000001  
 $\Rightarrow$  actual exponent =  $1 - 1023 = -1022$
  - Fraction: 000...00  $\Rightarrow$  significand = 1.0
  - $\pm 1.0 \times 2^{-1022} \approx \pm 2.2 \times 10^{-308}$
- Largest value
  - Exponent: 11111111110  
 $\Rightarrow$  actual exponent =  $2046 - 1023 = +1023$
  - Fraction: 111...11  $\Rightarrow$  significand  $\approx 2.0$
  - $\pm 2.0 \times 2^{+1023} \approx \pm 1.8 \times 10^{+308}$

# Floating-Point Example

- What number is represented by the single-precision float

11000000101000...00

- $S = 1$
  - Fraction =  $01000...00_2$
  - Exponent =  $10000001_2 = 129$
- $x = (-1)^1 \times (1 + 01_2) \times 2^{(129 - 127)}$   
 $= (-1) \times 1.25 \times 2^2$   
 $= -5.0$

Single precision		Double precision		Object represented
Exponent	Fraction	Exponent	Fraction	
0	0	0	0	0
0	Nonzero	0	Nonzero	$\pm$ denormalized number
1–254	Anything	1–2046	Anything	$\pm$ floating-point number
255	0	2047	0	$\pm$ infinity
255	Nonzero	2047	Nonzero	NaN (Not a Number)

# Single Precision Examples

- Denormalized Numbers

Type	Sign	Actual Exponent	Exp (biased)	Exponent field	Fraction field	Value
Zero	0	−126	0	0000 0000	000 0000 0000 0000 0000 0000	0.0
Negative zero	1	−126	0	0000 0000	000 0000 0000 0000 0000 0000	−0.0
One	0	0	127	0111 1111	000 0000 0000 0000 0000 0000	1.0
Minus One	1	0	127	0111 1111	000 0000 0000 0000 0000 0000	−1.0
Smallest denormalized number	*	−126	0	0000 0000	000 0000 0000 0000 0000 0001	$\pm 2^{-23} \times 2^{-126} = \pm 2^{-149} \approx \pm 1.4 \times 10^{-45}$
"Middle" denormalized number	*	−126	0	0000 0000	100 0000 0000 0000 0000 0000	$\pm 2^{-1} \times 2^{-126} = \pm 2^{-127} \approx \pm 5.88 \times 10^{-39}$
Largest denormalized number	*	−126	0	0000 0000	111 1111 1111 1111 1111 1111	$\pm (1 - 2^{-23}) \times 2^{-126} \approx \pm 1.18 \times 10^{-38}$
Smallest normalized number	*	−126	1	0000 0001	000 0000 0000 0000 0000 0000	$\pm 2^{-126} \approx \pm 1.18 \times 10^{-38}$
Largest normalized number	*	127	254	1111 1110	111 1111 1111 1111 1111 1111	$\pm (2 - 2^{-23}) \times 2^{127} \approx \pm 3.4 \times 10^{38}$
Positive infinity	0	128	255	1111 1111	000 0000 0000 0000 0000 0000	$+\infty$
Negative infinity	1	128	255	1111 1111	000 0000 0000 0000 0000 0000	$-\infty$
Not a number	*	128	255	1111 1111	non zero	NaN

\* Sign bit can be either 0 or 1 .

# Precisions in Intel

- Extended precision: 80 bits (Intel only)



# Special Properties of Encoding

- Floating Point 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

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



# Rounding

- Rounding Modes (illustrate with \$ rounding)

Positive number downward  
and negative number upward

Both positive negative  
numbers downward

- Towards zero
- Round down ( $-\infty$ )
- Round up ( $+\infty$ )
- Nearest Even (default)

\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
\$1	\$1	\$1	\$2	-\$1
\$1	\$1	\$1	\$2	-\$2
\$2	\$2	\$2	\$3	-\$1
\$1	\$2	\$2	\$2	-\$2

Both positive negative  
numbers upward

Attempt to find the  
closest match – least  
significant digit is even

- What are the advantages of the modes?

# Closer Look at Round-To-Even

- Default Rounding Mode: **Why Even? Why not consistently round values halfway between two representable value upward or downward?**
  - All others are statistically biased
    - Sum of set of positive numbers will consistently be over- or under- estimated
    - The average of a set of numbers that we rounded by this means would always rounded number would wither be slight higher or lower than the average of the number itself.
- 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
  - “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
$2 \frac{3}{32}$	$10.00\textcolor{red}{011}_2$	$10.00_2$	(< $1/2$ —down)	2
$2 \frac{3}{16}$	$10.00\textcolor{red}{110}_2$	$10.01_2$	(> $1/2$ —up)	$2 \frac{1}{4}$
$2 \frac{7}{8}$	$10.11\textcolor{red}{100}_2$	$11.00_2$	( $1/2$ —up)	3
$2 \frac{5}{8}$	$10.10\textcolor{red}{100}_2$	$10.10_2$	( $1/2$ —down)	$2 \frac{1}{2}$

# Floating-Point Addition

Consider a 4-digit decimal example

$$9.999 \times 10^1 + 1.610 \times 10^{-1}$$

1. Align decimal points

Shift number with smaller exponent

$$9.999 \times 10^1 + 0.016 \times 10^1$$

2. Add significands

$$9.999 \times 10^1 + 0.016 \times 10^1 = 10.015 \times 10^1$$

3. Normalize result & check for over/underflow

$$1.0015 \times 10^2$$

4. Round and renormalize if necessary. Assume only four digits are allowed for significand and two digits for exponent

$$1.002 \times 10^2$$

# Floating-Point Addition

Now consider a 4-digit binary example

$$1.000_2 \times 2^{-1} + -1.110_2 \times 2^{-2} \quad (0.5 + -0.4375)$$

## 1. Align binary points

Shift number with smaller exponent

$$1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1}$$

## 2. Add significands

$$1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1} = 0.001_2 \times 2^{-1}$$

## 3. Normalize result & check for over/underflow

$$1.000_2 \times 2^{-4}, \text{ (no over/underflow } 127 \geq -4 \geq -126)$$

## 4. Round and renormalize if necessary

$$1.000_2 \times 2^{-4} \text{ (no change) } = 0.0625$$

# Floating-Point Multiplication

- Consider a 4-digit decimal example
  - $1.110 \times 10^{10} \times 9.200 \times 10^{-5}$
- 1. Add exponents
  - For biased exponents, subtract bias from sum
  - New exponent =  $10 + -5 = 5$
- 2. Multiply significands
$$1.110 \times 9.200 = 10.212 \Rightarrow 10.212 \times 10^5$$
- 3. Normalize result & check for over/underflow
$$1.0212 \times 10^6$$
- 4. Round and renormalize if necessary
$$1.021 \times 10^6$$
- 5. Determine sign of result from signs of operands
$$+1.021 \times 10^6$$

Biased  $10 + 127 = 137$ , and  $-5 + 127 = 122$ ,

New exponent

$137 + 122 = 259$  Wrong !!!

$(10 + 127) + (-5 + 127) = (5 + 2 * 127) = 259$

*we must subtract the bias from the sum:*

New exponent  $137 + 122 - 127 =$

$259 - 127 = 132 = (5 + 127)$

# Floating-Point Multiplication

- Now consider a 4-digit binary example
  - $1.000_2 \times 2^{-1} \times -1.110_2 \times 2^{-2}$  ( $0.5 \times -0.4375$ )
- 1. Add exponents
  - Unbiased:  $-1 + -2 = -3$
  - Biased:  $(-1 + 127) + (-2 + 127) = -3 + 254 - 127 = -3 + 127$
- 2. Multiply significands
  - $1.000_2 \times 1.110_2 = 1.110000_2 \Rightarrow 1.110_2 \times 2^{-3}$
- 3. Normalize result & check for over/underflow
  - $1.110_2 \times 2^{-3}$  (no change) &  $127 \geq -3 \geq -126$  no over/underflow
- 4. Round and renormalize if necessary
  - $1.110_2 \times 2^{-3}$  (no change)
- 5. Determine sign:  $+ve \times -ve \Rightarrow -ve$ 
  - $-1.110_2 \times 2^{-3} = -0.21875$

# 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  $\rightarrow$  int
    - Truncates fractional part
    - Like rounding toward zero
    - Not defined when out of range or NaN: Generally sets to TMin
  - int  $\rightarrow$  double
    - Exact conversion, as long as int has  $\leq 53$  bit word size
  - int  $\rightarrow$  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);`
- `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