# CSSE132
# Introduction to Computer Systems

6 : Floating point

March 12, 2013

# Today: Floating Point

- IEEE Rounding modes
- Conversion process
- Addition, multiplication
- Floating point in C

# 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)

■

|  | $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 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

|           |      |                         |
|-----------|------|-------------------------|
| 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 = $100\ldots_2$

- Examples
  - Round to nearest 1/4 (2 bits right of binary point)

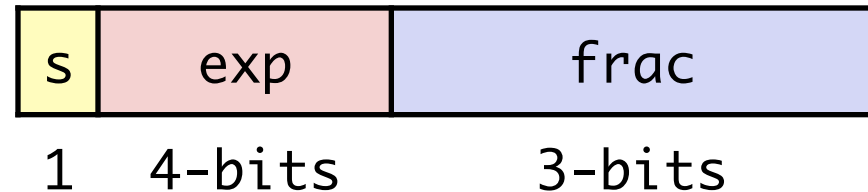| Value | Binary | Rounded | Action | Rounded Value |
|---|---|---|---|---|
| 2 3/32 | $10.00011_2$ | $10.00_2$ | (<1/2—down) | 2 |
| 2 3/16 | $10.00110_2$ | $10.01_2$ | (>1/2—up) | 2 1/4 |
| 2 7/8 | $10.11100_2$ | $11.00_2$ | ( 1/2—up) | 3 |
| 2 5/8 | $10.10100_2$ | $10.10_2$ | ( 1/2—down) | 2 1/2 |

# Today: Floating Point

- IEEE Rounding modes
- Conversion process
- Addition, multiplication
- Floating point in C

# Creating Floating Point Number

- Steps
  - Normalize to have leading 1
  - Round to fit within fraction
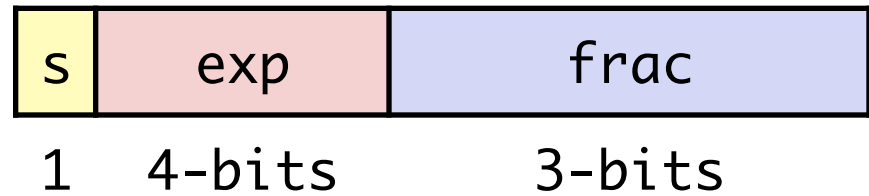  - Postnormalize to deal with effects of rounding

| s | exp | frac |
|---|-----|------|
| 1 | 4-bits | 3-bits |

- 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

| s | exp | frac |
|---|-----|------|
| 1 | 4-bits | 3-bits |

- **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 | 00001101 | 1.1010000 | 3 |
| 17 | 00010001 | 1.0001000 | 4 |
| 19 | 00010011 | 1.0011000 | 4 |
| 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 ➡ > 0.5
  - Guard = 1, Round = 1, Sticky = 0 ➡ Round to even

| Value | Fraction | GRS | Incr? | Rounded |
|-------|----------|-----|-------|---------|
| 128 | 1.0000000 | 000 | N | 1.000 |
| 15 | 1.1010000 | 100 | 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     |

# FP Multiplication

- $(-1)^{s1}$ M1  $2^{E1}$   x   $(-1)^{s2}$ M2  $2^{E2}$
- Exact Result: $(-1)^{s}$ M  $2^{E}$
  - Sign s:                    s1 ^ s2
  - Significand M:        M1 x  M2
  - Exponent E:           E1 + E2

- Fixing
  - If M ≥ 2, shift M right, increment E
  - If E out of range, overflow
  - Round M to fit `frac` precision

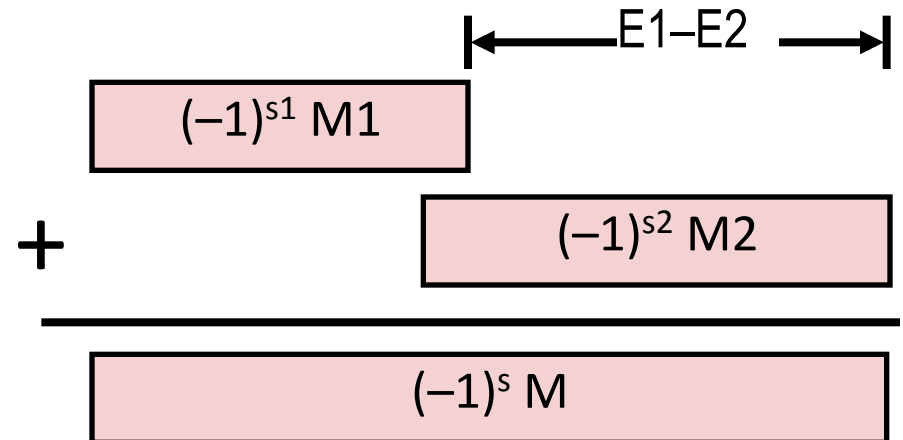- Implementation
  - Biggest chore is multiplying significands

# 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:     E1



- Fixing
  - If M ≥ 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

13

# Today: Floating Point

- IEEE Rounding modes
- Conversion process
- Addition, multiplication
- Floating point in C

# 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

# 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    ⇒   ((d*2) < 0.0)
- d > f      ⇒   -f > -d
- d * d >= 0.0
- (d+f)-d == f