# Floating Point – rounding and operations

These slides adapted from materials provided by the textbook

# 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

- $x +_f y = Round(x + y)$

- $x \times_f y = Round(x \times y)$

- **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 |

# 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

| | | |
|---|---|---|
| 7.8949999 | 7.89 | (Less than half way) |
| 7.8950001 | 7.90 | (Greater than half way) |
| 7.8950000 | 7.90 | (Half way—round up) |
| 7.8850000 | 7.88 | (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 |
|-------|--------|---------|--------|---------|
| 2 3/32 | $10.00\underline{011}_2$ | $10.00_2$ | (<1/2—down) | 2 |
| 2 3/16 | $10.00\underline{110}_2$ | $10.01_2$ | (>1/2—up) | 2 1/4 |
| 2 7/8 | $10.11\underline{100}_2$ | $11.00_2$ | ( 1/2—up) | 3 |
| 2 5/8 | $10.10\underline{100}_2$ | $10.10_2$ | ( 1/2—down) | 2 1/2 |

# Rounding Binary Numbers To Even

| Value | Binary | Rounded | Action |
|-------|--------|---------|--------|
| 2 7/8 | $10.11\underline{100}_2$ | $11.00_2$ | ( 1/2—up) |
| 2 5/8 | $10.10\underline{100}_2$ | $10.10_2$ | ( 1/2—down) |

# Creating Floating Point Number

- **Steps**

  - Normalize to have leading 1

  - Round to fit within fraction

  - Postnormalize to deal with effects of rounding

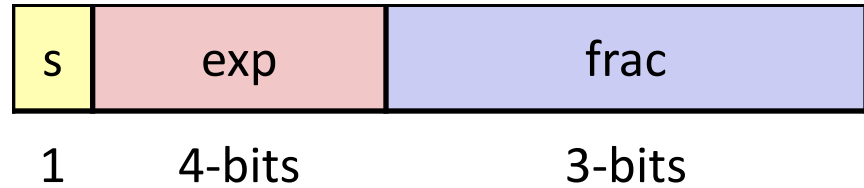| s | exp | frac |
|---|-----|------|
| 1 | 4-bits | 3-bits |

- **Case Study**

  - Convert 8-bit unsigned numbers to tiny floating point format

Example Numbers

| | |
|-----|----------|
| 128 | 10000000 |
| 15  | 00001101 |
| 17  | 00010001 |
| 19  | 00010011 |
| 138 | 10001010 |
| 63  | 00111111 |

# Normalize

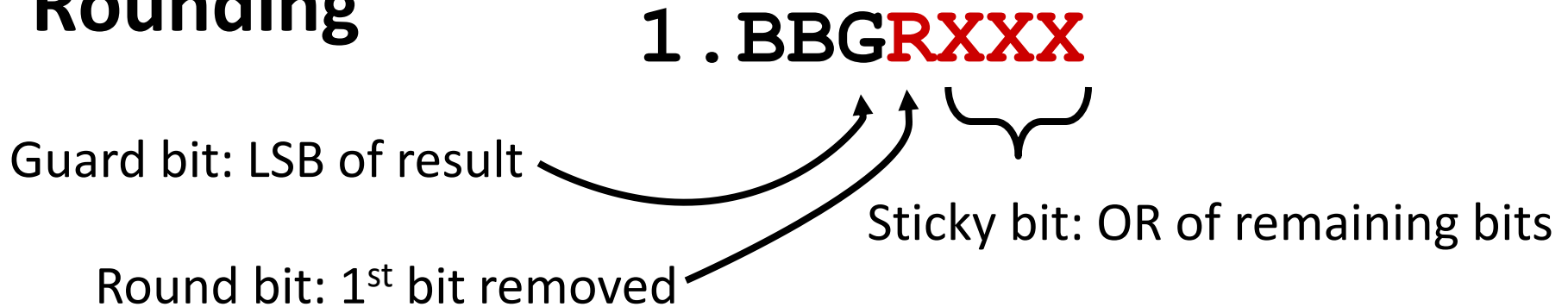| s | exp | frac |
|---|-----|------|
| 1 | 4-bits | 3-bits |

- **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

$$\texttt{1.BBG}\color{red}{\texttt{RXXX}}$$

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 |

# Interesting Numbers {single,double}

| Description | exp | frac | Numeric Value |
|---|---|---|---|
| **Zero** | 00…00 | 00…00 | 0.0 |
| **Smallest Pos. Denorm.** | 00…00 | 00…01 | $2^{-\{23,52\}} \times 2^{-\{126,1022\}}$ |

- Single $\approx 1.4 \times 10^{-45}$
- Double $\approx 4.9 \times 10^{-324}$

| | | | |
|---|---|---|---|
| **Largest Denormalized** | 00…00 | 11…11 | $(1.0 - \varepsilon) \times 2^{-\{126,1022\}}$ |

- Single $\approx 1.18 \times 10^{-38}$
- Double $\approx 2.2 \times 10^{-308}$

| | | | |
|---|---|---|---|
| **Smallest Pos. Normalized** | 00…01 | 00…00 | $1.0 \times 2^{-\{126,1022\}}$ |

- Just larger than largest denormalized

| | | | |
|---|---|---|---|
| **One** | 01…11 | 00…00 | 1.0 |
| **Largest Normalized** | 11…10 | 11…11 | $(2.0 - \varepsilon) \times 2^{\{127,1023\}}$ |

- Single $\approx 3.4 \times 10^{38}$
- Double $\approx 1.8 \times 10^{308}$

# 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

Get binary points lined up



- **Exact Result: $(-1)^s\ M\ 2^E$**
  - Sign s, significand M:
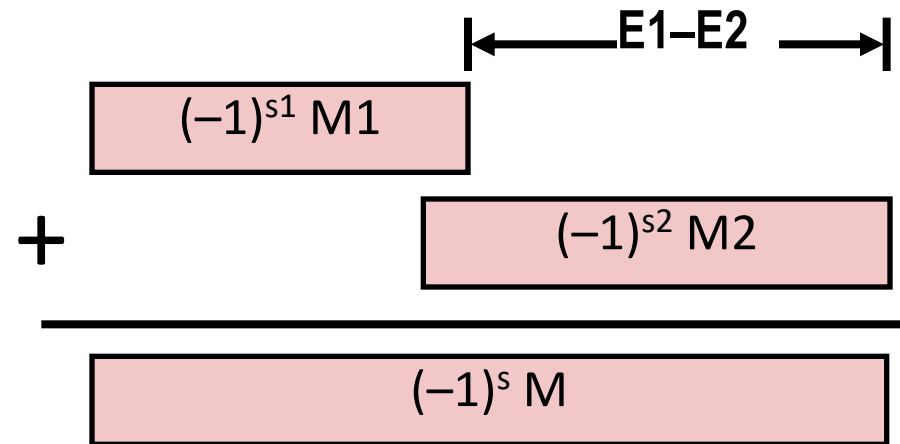    - Result of signed align & add
  - Exponent E:     E1

- **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

# 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?     Yes
  - Every element has additive inverse?     Almost
    - Yes, except for infinities & NaNs

- **Monotonicity**     Almost
  - $a \geq b \Rightarrow a+c \geq b+c$?
    - 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?                        No
    - Possibility of overflow, inexactness of rounding
    - Ex: `(1e20*1e20)*1e-20=` `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 = NaN`
                                                          Almost
- **Monotonicity**
  - $a \geq b$ & $c \geq 0 \Rightarrow a * c \geq b * c$?
    - Except for infinities & NaNs

# 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

# 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`

# Summary

- **IEEE Floating Point has clear mathematical properties**
- **Represents numbers of form M x $2^E$**
- **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