

COMPSCI 210 S1T

Computer Systems

Floating Point Numbers Arithmetic

Floating Point Operations

- ◆ Exponents and fraction must be handled separately & differently.
- ◆ Consider the Floating point numbers:
 - $X = M_x \cdot 2^{E_x}$, $Y = M_y \cdot 2^{E_y}$
- ◆ Basic Operations
 - Addition: $X + Y = (M_x \cdot 2^{E_x} + M_y \cdot 2^{E_y}) \cdot 2^{\min(E_x, E_y)}$, $E_x <= E_y$
 - Subtraction: $X - Y = (M_x \cdot 2^{E_x} - M_y \cdot 2^{E_y}) \cdot 2^{\min(E_x, E_y)}$, $E_x <= E_y$
- ◆ Procedures for addition/subtraction:
 - Adjust exponents and align mantissa
 - The exponent of the operands must be made equal for addition and subtraction.
 - If $E_y > E_x$ Right shift M_x to form $M_x \cdot 2^{E_y - E_x}$
 - If $E_x > E_y$ Right shift M_y to form $M_y \cdot 2^{E_x - E_y}$
 - Add or subtract mantissa
 - Normalize the result
 - Left shift result, decrement result exponent (e.g., if result is $0.001xx\dots$) or
 - Right shift result, increment result exponent (e.g., if result is $10.1xx\dots$)
 - Check result
 - If larger than maximum exponent allowed return exponent overflow
 - If smaller than minimum exponent allowed return exponent underflow
 - If result mantissa is 0, may need to set the exponent to zero to return a zero.

Agenda & Reading

◆ Agenda:

- Floating Point operations
 - ◆ Addition
 - ◆ Subtraction
 - ◆ Multiplication
 - ◆ Division

Example on Decimal Value

◆ $3.25 \cdot 10^2 + 2.63 \cdot 10^{-1}$

◆ Steps:

- 1) Align decimal points:
 - ◆ Start by adjusting the smaller exponent to be equal to the larger exponent, and modify the fraction accordingly.
 - ◆ Take $2.63 \cdot 10^{-1}$
 - Original value = $2.63 \cdot 10^{-1}$ (right-shift mantissa-> increase exponent)
 - Shift 1 place = $0.263 \cdot 10^0$
 - Shift 2 places = $0.0263 \cdot 10^1$
 - Shift 3 places = $0.00263 \cdot 10^2$
- 2) Add the numbers
 - ◆
$$\begin{array}{r} 3.25 \cdot 10^2 \\ + 0.00263 \cdot 10^2 \\ \hline 3.2563 \cdot 10^2 \end{array}$$
- 3) Normalize the result
 - ◆ No need to change. It is already normalized.
- 4) Check result
 - ◆ OK. Answer = $3.2563 \cdot 10^2$

Example on Binary value

◆ $0.101 \cdot 2^2 + 0.111 \cdot 2^4$

◆ Steps:

- 1) Align decimal points:
 - Start by adjusting the smaller exponent to be equal to the larger exponent, and modify the fraction accordingly.
 - Take $0.101 \cdot 2^2$
 - Original value = $0.101 \cdot 2^2$
 - Shifted 1 place = $0.101 \cdot 2^2 \Rightarrow 0.0101 \cdot 2^3$
 - Shifted 2 places = $0.0101 \cdot 2^3 \Rightarrow 0.00101 \cdot 2^4$
 - 2) Add the numbers
 - $0.111 \cdot 2^4$
 - $+ 0.00101 \cdot 2^4$
 - $1.00001 \cdot 2^4$
 - 3) Normalize the result
 - No need to change. It is normalized.
 - 4) Check result
 - OK. Answer = $1.0011 \cdot 2^4$

COMPSCI 210

08

5

Floating Point Addition

◆ $1.25 + 0.25$

- $0.25 = 0.01111101 000\ldots 000$
- $1.25 = 0.01111111 010\ldots 000$

◆ Steps:

- Adjust exponents and align mantissa
 - Start by adjusting the smaller exponent to be equal to the larger exponent
 - Take 0.25 ($0.01111101 000\ldots 000$) (with smaller exponent)
 - Original Value: E:01111101 M:000...000
 - Shifted 1 place: E:01111110 M:100...000 (Note: "1" is the hidden bit)
 - Shifted 2 places: E:01111111 M:010...000
- Add mantissa bits
 - $0.01111111 1.010\ldots 000$
 - $+ 0.01111111 0.010\ldots 000$
 - $0.01111111 1.100\ldots 000$
- Normalize result: **The hidden bit**
 - No need to change. It is normalized.
- Check result
 - OK. Answer = $0.01111111 100\ldots 000$

COMPSCI 210

08

6

Floating Point Addition

◆ $6.5 + 1.75$

- $6.5 = 0.10000001 1010\ldots 000$
- $1.75 = 0.01111111 110\ldots 000$

◆ Steps:

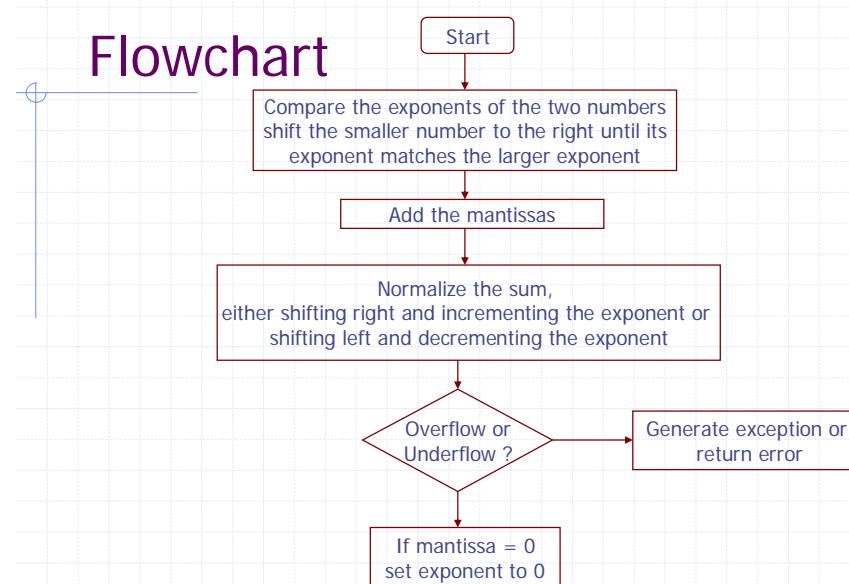
- Adjust exponents and align mantissa
 - Start by adjusting the smaller exponent to be equal to the larger exponent
 - Take 1.75 ($0.01111111 110\ldots 000$) (with smaller exponent)
 - Original Value: E:01111111 M:110...000
 - Shifted 1 place: E:10000000 M:1110...000 (Note: "1" is the hidden bit)
 - Shifted 2 places: E:10000001 M:0110...000
- Add mantissa bits
 - $0.10000001 1.1010\ldots 000$
 - $+ 0.10000001 0.01110\ldots 000$
 - $0.10000001 10.0001000\ldots 000$
- Normalize result:
 - Shift the mantissa right by 1 bit increase the exponent by 1
 - E: 10000001 +1 = 10000010, mantissa = 1.0000100..0
- Check result:
 - OK. Answer = $0.10000010 0000100\ldots 000$

COMPSCI 210

08

7

Flowchart



COMPSCI 210

08

8

Floating Point Subtraction

◆ 1.25 - 0.25

- 0.25 = 0 01111101 000...000
- 1.25 = 0 01111111 010...000

◆ Steps:

- Adjust exponents and align mantissa
 - ♦ Start by adjusting the smaller exponent to be equal to the larger exponent
 - ♦ Take 0.25 (0 01111101 000...000) (with smaller exponent)
 - Shifted 1 place: E:01111110 M:100...000 (Note: "1" is the hidden bit)
 - Shifted 2 places: E:01111111 M:010...000
- Subtract mantissa bits
 - ♦ 0 01111111 1.010...000
 - ♦ -0 01111111 0.010...000
 - ♦ 0 01111111 1.000...000
- Normalize result:
 - ♦ No need to change. It is normalized.
- Check result
 - ♦ OK. Answer = 0 01111111 000...000

COMPSCI 210

08

9

Addition Notes

◆ Truncation Errors

- If the exponents differ by more than 24, the smaller number will be shifted right entirely out of the mantissa field, producing a zero mantissa.
- The sum will then equal the larger number.
- Such truncation errors occur when the numbers differ by a factor of more than 2^{24} , which is approximately 1.6×10^7 .
- Thus, the precision of IEEE single precision floating point arithmetic is approximately 7 decimal digits.

◆ Alternatively, floating point subtraction is achieved simply by inverting the sign bit and performing addition of signed mantissas.

- Negative mantissas are handled by first converting to 2's complement and then performing the addition.
- After the addition is performed, the result is converted back to sign-magnitude form.

◆ When adding numbers of opposite sign, cancellation may occur, resulting in a sum which is arbitrarily small, or even zero if the numbers are equal in magnitude.

COMPSCI 210

08

11

Floating Point Subtraction

◆ 6.5 - 1.75

- 6.5 = 0 10000001 1010...000
- 1.75 = 0 01111111 110...000

◆ Steps:

- Adjust exponents and align mantissa
 - ♦ Start by adjusting the smaller exponent to be equal to the larger exponent
 - ♦ Take 1.75 (0 01111111 110...000) (with smaller exponent)
 - Shifted 1 place: E:10000000 M:1110...000 (Note: "1" is the hidden bit)
 - Shifted 2 places: E:10000001 M:0110...000
- Subtract mantissa bits
 - ♦ 0 10000001 1.10100...000
 - ♦ -0 10000001 0.01110...000
 - ♦ 0 10000001 1.001100...000
- Normalize result:
 - ♦ No need to change.
- Check result:
 - ♦ OK. Answer = 0 10000001 00110...00

COMPSCI 210

08

10

Multiplication

◆ Floating Point Operations:

- $X = (-1)^{S_x} M_x \times 2^{E_x}, Y = (-1)^{S_y} M_y \times 2^{E_y}$
- Multiplicatoin: $X * Y = (M_x * M_y) \times 2^{E_x+E_y}$

◆ Procedures for Multiplicatoin:

- Check Zeros
 - ♦ If one or both operands is equal to zero, return the result as zero.
- Compute the sign of the result $S_x \text{ XOR } S_y$
- Multiply mantissa
 - ♦ $M_x * M_y$
 - ♦ Round the result to the allowed number of mantissa bits
- Add exponents
 - ♦ biased exponent (E_x) + biased exponent (E_y) - bias
- Normalize the result
 - ♦ Left shift result, decrement result exponent (e.g., if result is 0.001xx...) or
 - ♦ Right shift result, increment result exponent (e.g., if result is 10.1xx...)
- Check result
 - ♦ If larger than maximum exponent allowed return exponent overflow
 - ♦ If smaller than minimum exponent allowed return exponent underflow

COMPSCI 210

08

12

Decimal Multiplication/Division

◆ Example 1: $3.0 * 10^1 * 4.5 * 10^2$

◆ Steps:

- Sign = 0 XOR 0 = 0
- Multiply mantissa (don't forget the hidden bit)
 - $3.0 * 4.5 = 13.5$
- Add exponents
 - $1 + 2 = 3$
- Normalize the result
 - $13.5 * 10^3$, Shift 1 place -> increase exponent = $1.35 * 10^4$
- Check result
 - OK. Answer = $1.35 * 10^4$

◆ Example 2: $4.0 * 10^3 * 0.5 * 10^2$

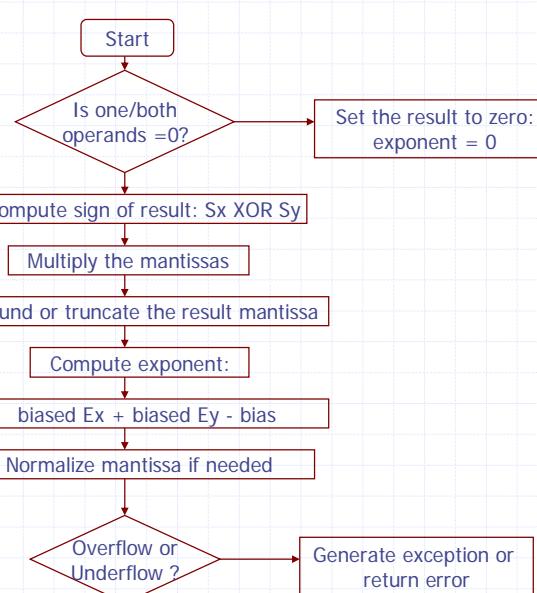
- Sign = 0 XOR 0 = 0
- Divide mantissa (don't forget the hidden bit)
 - $4.0 / 0.5 = 8.0$
- Subtract exponents
 - $3 - 2 = 1$
- Normalize the result
 - $8.0 * 10^1$ No change.
- Check result
 - OK. Answer = $8.0 * 10^1$

COMPSCI 210

08

13

Flowchart



COMPSCI 210

08

15

Floating Point Multiplication

◆ $-18 * 9.5$

- $-18 = 1 10000011 0010..000$
- $9.5 = 0 10000010 0011..000$

◆ Steps:

- Sign = 0 XOR 1 = 1
- Multiply mantissa (don't forget the hidden bit)

$$\begin{array}{r}
 1.0010 \\
 * 1.0011 \\
 \hline
 10010 \\
 10010 \\
 00000 \\
 00000 \\
 10010 \\
 101010110 = 1.01010110
 \end{array}$$

- Add exponents
 - $10000011 + 10000010 - 01111111 = 100000110$

- Normalize the result
 - It is normalized. No change

- Check result
 - OK. Answer = $1 10000110 01010110..0$

COMPSCI 210

08

14

Division

◆ Floating Point Operations:

- $X = (-1)^{S_x} M_x * 2^{E_x}, Y = (-1)^{S_y} M_y * 2^{E_y}$
- Division: $X / Y = (M_x / M_y) * 2^{E_x - E_y}$

◆ Procedures for Division:

- Check Zeros
 - If both operands is equal to zero, return the result as NaN.
 - If Y is equal to zero, return the result as infinity.
- Compute the sign of the result $S_x \text{ XOR } S_y$
- Divide mantissa
 - M_x/M_y
 - Round the result to the allowed number of mantissa bits
- Subtract exponents
 - Division: biased exponent (E_x) - biased exponent (E_y) + bias
- Normalize the result
 - Left shift result, decrement result exponent (e.g., if result is $0.001xx\dots$) or
 - Right shift result, increment result exponent (e.g., if result is $10.1xx\dots$)
- Check result
 - If larger than maximum exponent allowed return exponent overflow
 - If smaller than minimum exponent allowed return exponent underflow

COMPSCI 210

08

16

Floating Point Division

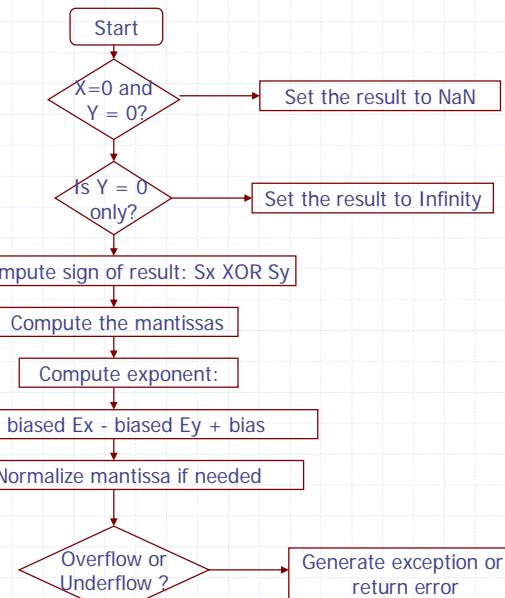
- ◆ $3.75 / 1.5$
 - $3.75 = 0\ 10000000\ 1110...000$
 - $1.5 = 0\ 01111111\ 100...000$
- ◆ Steps:
 - Sign = 0 XOR 0 = 0
 - Divide mantissa (don't forget the hidden bit)
 - ♦ $\begin{array}{r} \underline{101} \\ 11) \underline{1111} \\ \underline{11} \\ \underline{11} \\ \underline{0} = \\ = 101 \Rightarrow 1.01 \end{array}$
 - Subtract exponents
 - ♦ $10000000 - 01111111 + 01111111 = 1000\ 0000$
 - Normalize the result
 - ♦ It is normalized. No change
 - Check result
 - ♦ OK. Answer = 0 10000000 0100...000

COMPSCI 210

08

17

Flowchart



COMPSCI 210

08

18

Conversion Example:

- ◆ What is the IEEE floating point representation of 6.5_{10} ? Hence, what is the representation for 52?
- Step 1: 6.5_{10} :
 - ♦ Sign = 0
 - ♦ $6.5_{10} \Rightarrow$ Binary number = 110.1
 - ♦ Normalization (shift radix point to left by 2 places) $\Rightarrow 1.101 \times 2^2$
 - ♦ Exponent = $127+2= 129 = 10000001$
 - ♦ Mantissa = 1010...0
 - ♦ Answer = 0 10000001 1010...0 = 40D00000
- Step 2: $52/6.5 = 8 = 2^3$ i.e. $6.5 \times 2^3 = 52$
 - ♦ Sign bit : unchanged
 - ♦ Mantissa : unchanged
 - ♦ Exponent : exponent from step 1 + 3 = $(129 + 3) = 10000100$
 - ♦ Answer = 0 10000100 1010...0 = 42500000

COMPSCI 210

08

19

Conversion Example (con't)

- ◆ What is the IEEE floating point representation of 1.25_{10} ? Hence, what is the representation for 1.75 and 1.875?
- Step 1: $1.25 = 3FA00000 = 0\ 01111111\ 010...0$
- Step 2: 1.75
 - ♦ $1.75 = 1.25 * 2^0 + 0.5 * 2^0$
 - ♦ Sign bit: unchanged
 - ♦ Exponent bit: unchanged
 - ♦ Mantissa = + 0.5
 - ♦ Answer = 0 01111111 110...0 = 3FE00000
- Step 3:
 - ♦ $1.875 = 1.25 * 2^0 + 0.5 * 2^0 + 0.125 * 2^0$
 - ♦ Sign bit: unchanged
 - ♦ Exponent bit: unchanged
 - ♦ Mantissa = + 0.5 + 0.125
 - ♦ Answer = 0 01111111 1110...0 = 3FF00000
- ◆ Note: exponent = $2^0 = 1$ only
 - If exponent is not equal to 1, you need to take this into account

COMPSCI 210

08

20

Conversion Example (con't)

What is the IEEE floating point representation of 0.625_{10} ?
Hence, what is the representation for 0.875_{10} ?

■ Step 1: $0.625 = 3F200000 = 0\ 01111110\ 010...0$

■ Step 2: 0.875

$$\diamond 0.875 = 0.625 + 0.25$$

$$\diamond = 0.625 + 0.5 * 2^{-1}$$

$$\diamond = 1.25 * 2^{-1} + 0.5 * 2^{-1}$$

♦ Sign bit: same

♦ Exponent bit: same

♦ Mantissa = + 0.5

$$\diamond \text{Answer} = 0\ 01111110\ 110...0 = 3F600000$$