# CSC 3210 Computer organization and programming
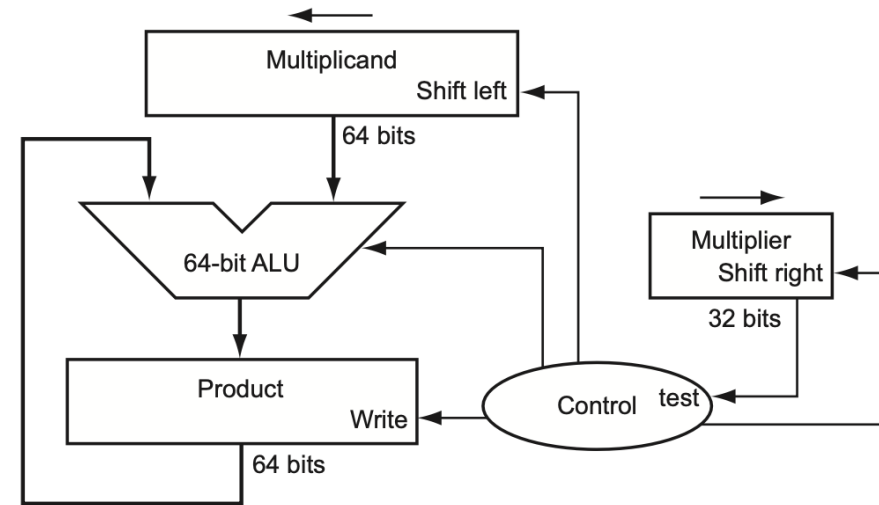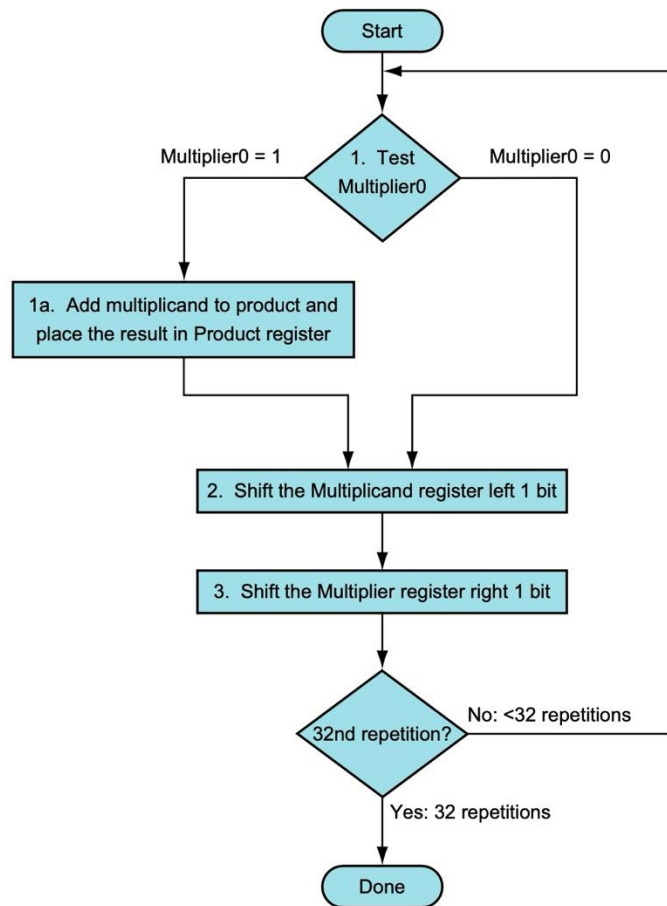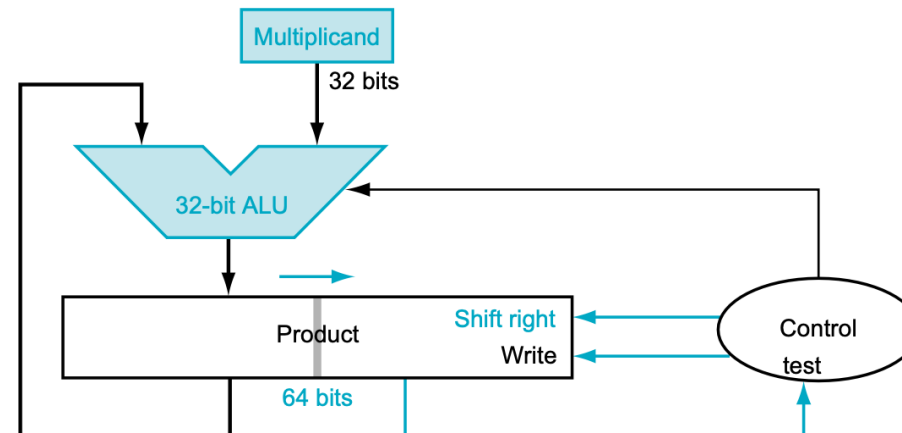
## Chapter 3 & Quiz3

Chunlan Gao

# Multiplication-Review



no    muli
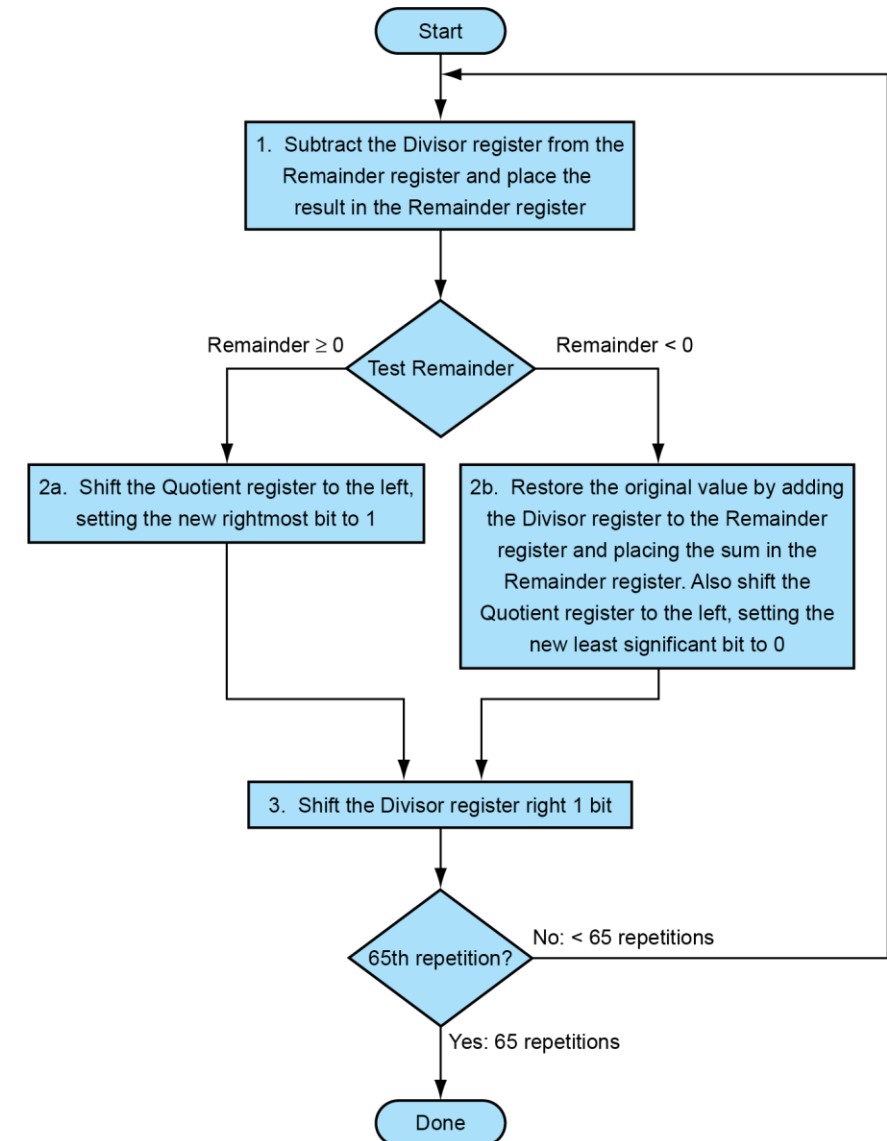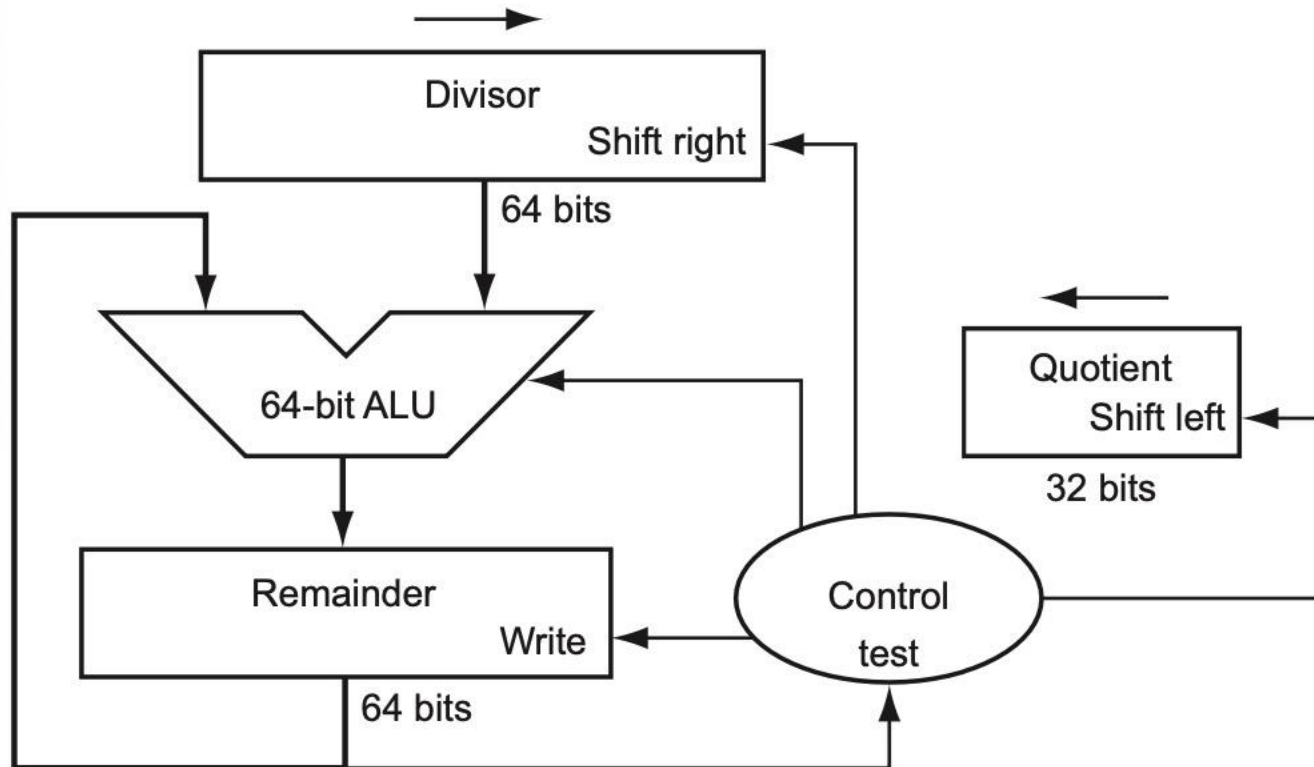
immediate 12bits, not big enough.

# Division

- Check for 0 divisor
- Long division approach
  - If divisor ≤ dividend bits
    - 1 bit in quotient, subtract
  - Otherwise
    - 0 bit in quotient, bring down next dividend bit
- Restoring division
  - Do the subtract, and if remainder goes < 0, add divisor back
- Signed division
  - Divide using absolute values
  - Adjust sign of quotient and remainder as required

$$1001_{ten} \quad \text{Quotient}$$

$$\text{Divisor } 1000_{ten} \overline{)1001010_{ten}} \quad \text{Dividend}$$

$$-1000$$

$$10$$

$$101$$

$$1010$$

$$-1000$$

$$10_{ten} \quad \text{Remainder}$$

Dividend =Quotient* Divisor+ Remainder

# Division Hardware

# Example

division_remainder_left.docx

division_noremainder_left.docx

# Example-cont

| Iteration | Step | Quotient | Divisor | Remainder |
|---|---|---|---|---|
| 0 | Initial values | 0000 | 0010 0000 | 0000 0111 |
| 1 | 1: Rem = Rem – Div | 0000 | 0010 0000 | ①110 0111 |
| | 2b: Rem < 0 ⟹ +Div, SLL Q, Q0 = 0 | 0000 | 0010 0000 | 0000 0111 |
| | 3: Shift Div right | 0000 | 0001 0000 | 0000 0111 |
| 2 | 1: Rem = Rem – Div | 0000 | 0001 0000 | ①111 0111 |
| | 2b: Rem < 0 ⟹ +Div, SLL Q, Q0 = 0 | 0000 | 0001 0000 | 0000 0111 |
| | 3: Shift Div right | 0000 | 0000 1000 | 0000 0111 |
| 3 | 1: Rem = Rem – Div | 0000 | 0000 1000 | ①111 1111 |
| | 2b: Rem < 0 ⟹ +Div, SLL Q, Q0 = 0 | 0000 | 0000 1000 | 0000 0111 |
| | 3: Shift Div right | 0000 | 0000 0100 | 0000 0111 |
| 4 | 1: Rem = Rem – Div | 0000 | 0000 0100 | ⓪000 0011 |
| | 2a: Rem ≥ 0 ⟹ SLL Q, Q0 = 1 | 0001 | 0000 0100 | 0000 0011 |
| | 3: Shift Div right | 0001 | 0000 0010 | 0000 0011 |
| 5 | 1: Rem = Rem – Div | 0001 | 0000 0010 | ⓪000 0001 |
| | 2a: Rem ≥ 0 ⟹ SLL Q, Q0 = 1 | 0011 | 0000 0010 | 0000 0001 |
| | 3: Shift Div right | 0011 | 0000 0001 | 0000 0001 |

**FIGURE 3.10  Division example using the algorithm in Figure 3.9.** The bit examined to determine the next step is circled in color.
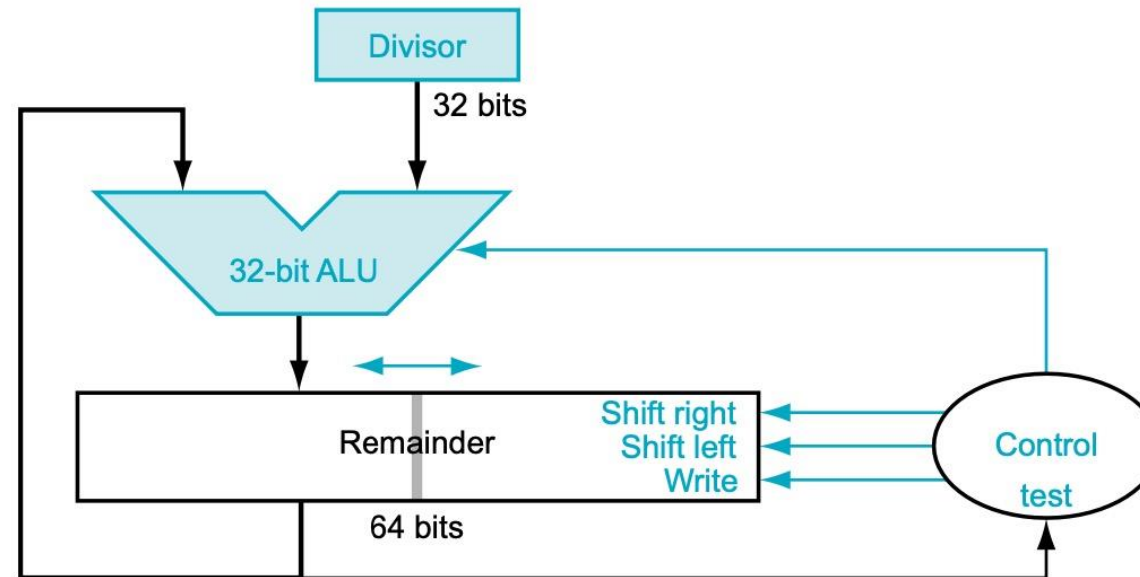
# Optimized Divider



**FIGURE 3.11 An improved version of the division hardware.** The Divisor register, ALU, and Quotient register are all 32 bits wide. Compared to Figure 3.8, the ALU and Divisor registers are halved and the remainder is shifted left. This version also combines the Quotient register with the right half of the Remainder register. As in Figure 3.5, the Remainder register has grown to 65 bits to make sure the carry out of the adder is not lost.

One cycle per partial-remainder subtraction Looks a lot like a multiplier!
Same hardware can be used for both

# Faster Division

- Can't use parallel hardware as in multiplier
  - Subtraction is conditional on sign of remainder

- Faster dividers generate multiple quotient bits per step
  - Still require multiple steps

# RISC-V Division

- Four instructions:
  - div, rem: signed divide, remainder
  - divu, remu: unsigned divide, remainder

- Overflow and division-by-zero don't produce errors
  - Just return defined results
  - Faster for the common case of no error

# Floating Point

- Representation for non-integral numbers
  - Including very small and very large numbers

- Like scientific notation
  - $-2.34\,0753 \times 10^{56}$     → normalized
  - $+0.00256779 \times 10^{-4}$
  - $+987.0245664 \times 10^{9}$     → unnormalized

- In binary
  - $\pm 1.xxxxxxx_2 \times 2^{yyyy}$

- Types `float` and `double` in C

# Floating-Point Representation

| 31 | 30 | 29 | 28 | 27 | 26 | 25 | 24 | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|---|---|---|---|---|---|---|---|---|
| s | exponent | | | | | | | | fraction | | | | | | | | | | | | | | | | | | | | | | |

1 bit      8 bits      23 bits

Floating-point numbers are of the form
$(-1)^S \times F \times 2^E$

F involves the value in the fraction field and E involves the value in the exponent field

Exponent can be positive or negative we need to check the sign number, To reduce the steps, we use a bias.

exponent = real exponent +127

# Floating Point Standard

- Defined by IEEE Std 754-1985
- Developed in response to divergence of representations
  - Portability issues for scientific code
- Now almost universally adopted
- Two representations
  - Single precision (32-bit)
  - Double precision (64-bit)

# IEEE Floating-Point Format

| Type(precision) | Sign | Exponent | Fraction |
|---|---|---|---|
| Single | 1bit | 8 bits | 23 bits |
| Double | 1bit | 23 bits | 52 bits |

$$(-1)^S \times (1 + \text{Fraction}) \times 2^E \qquad (-1)^S \times (1 + (s1 \times 2^{-1}) + (s2 \times 2^{-2}) + (s3 \times 2^{-3}) + (s4 \times 2^{-4}) + \ldots) \times 2^E$$

- S: sign bit (0 $\Rightarrow$ non-negative, 1 $\Rightarrow$ negative)
- Normalize significand: 1.0 ≤ |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 = 1203

# 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
  - $\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$(ten)
- Largest value
  - exponent: 11111110
    $\Rightarrow$ actual exponent = 254 – 127 = +127
  - Fraction: 111…11 $\Rightarrow$ significand $\approx$ 2.0
  - $\pm 2.0 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$

geometric sequence：

$$S_n = \frac{a_1\,(1 - q^n)}{1 - q}\,(q \neq 1)$$

$$Sn = \frac{1/2\,\left(1 - \frac{1}{2}\wedge n\right)}{1 - 1/2} = 1$$

# Double-Precision Range

- Exponents 0000…00 and 1111…11 reserved

- Smallest value
  - Exponent: 00000000001
    $\Rightarrow$ 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 Precision

- Relative precision
    - all fraction bits are significant
    - Single: approx $2^{-23}$
        - Equivalent to $23 \times \log_{10}2 \approx 23 \times 0.3 \approx 6$ decimal digits of precision
    - Double: approx $2^{-52}$
        - Equivalent to $52 \times \log_{10}2 \approx 52 \times 0.3 \approx 16$ decimal digits of precision

# Floating-Point Example

- Represent –0.75
  - $-0.75 = (-1)^1 \times 1.1_2 \times 2^{-1}$
  - S = 1
  - Fraction = $1000...00_2$
  - Exponent = –1 + Bias
    - Single: $-1 + 127 = 126 = 01111110_2$
    - Double: $-1 + 1023 = 1022 = 01111111110_2$
- Single: 101111110100...00
- Double: 1011111111101000...00

# Floating-Point Example

- What number is represented by the **single-precision** float

    1 10000001 01000...00

    - S = 1
    - Fraction = $01000...00_2$
    - Exponent = $10000001_2$ = 129

- $x = (-1)^1 \times (1 + 01_2) \times 2^{(129 - 127)}$

    $= (-1) \times 1.25 \times 2^2$

    $= -5.0$

# Denormal Numbers

- Exponent = 000...0 $\Rightarrow$ hidden bit is 0

$$x = (-1)^S \times (0 + \text{Fraction}) \times 2^{-\text{Bias}}$$

- Smaller than normal numbers
  - allow for gradual underflow, with diminishing precision
- Denormal with fraction = 000...0

- $$x = (-1)^S \times (0 + 0) \times 2^{-\text{Bias}} = \pm 0.0$$

Two representations of 0.0!

# Infinities and NaNs

- Exponent = 111…1, Fraction = 000…0
  - ±Infinity
  - Can be used in subsequent calculations, avoiding need for overflow check

- Exponent = 111…1, Fraction ≠ 000…0
  - Not-a-Number (NaN)
  - Indicates illegal or undefined result
    - e.g., 0.0 / 0.0
  - Can be used in subsequent calculations