# Division Hardware



- Start
  1. Subtract the Divisor register from the Remainder register and place the result in the Remainder register
  - Test Remainder
    - Remainder ≥ 0
    - Remainder < 0
  - 2a. Shift the Quotient register to the left, setting the new rightmost bit to 1
  - 2b. Restore the original value by adding the Divisor register to the Remainder register and placing the sum in the Remainder register. Also shift the Quotient register to the left, setting the new least significant bit to 0
  3. Shift the Divisor register right 1 bit
  - 33rd repetition?
    - No: < 33 repetitions
    - Yes: 33 repetitions
  - Done

Initially divisor in left half

Divisor — Shift right — 64 bits

64-bit ALU

Quotient — Shift left — 32 bits

Remainder — Write — 64 bits

Control test

Initially dividend

CSULB

# Optimized Divider



Divisor — 32 bits

32-bit ALU

Remainder — Shift right / Shift left / Write — 64 bits

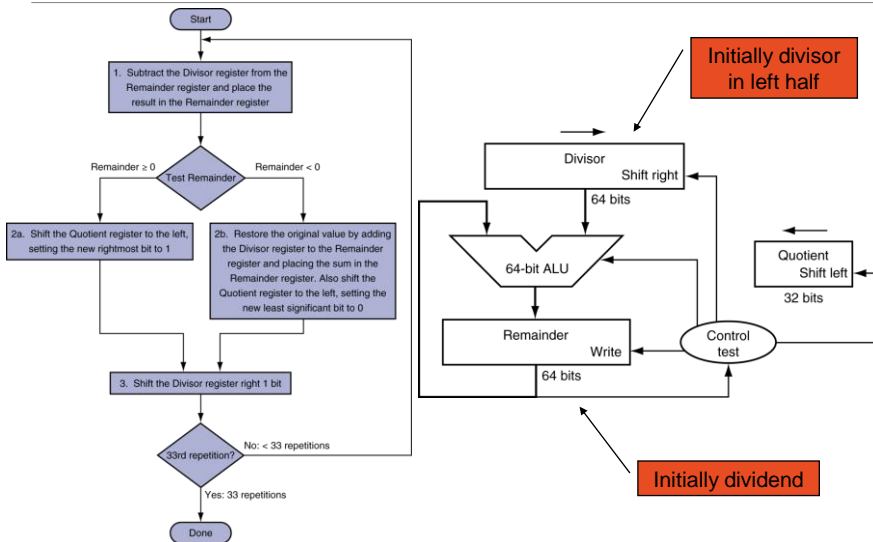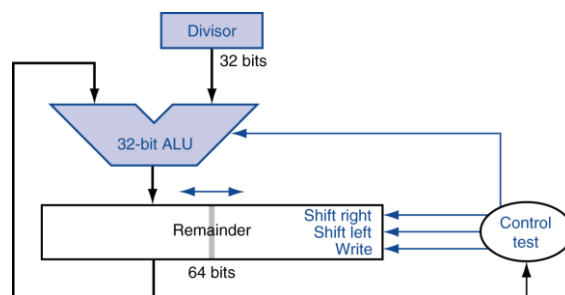Control test

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

CSULB

# MIPS Division

- Use HI/LO registers for result
  - HI: 32-bit remainder
  - LO: 32-bit quotient

- Instructions
- `div rs, rt  /  divu rs, rt`
- No overflow or divide-by-0 checking
  - Software must perform checks if required
- Use `mfhi, mflo` to access result

CSULB

# Floating Point

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

- Like scientific notation
- –2.34 × $10^{56}$  ← normalized
- +0.002 × $10^{-4}$ ← not normalized  → +2.0 × $10^{-7}$
- +987.02 × $10^{9}$ ←  → +9.8702 × $10^{7}$

- In binary
- ±1.$xxxxxxx_2$ × $2^{yyyy}$

- Types `float` and `double` in C

CSULB

# 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)

CSULB

# IEEE Floating-Point Format

| | single: 8 bits<br>double: 11 bits | single: 23 bits<br>double: 52 bits |
|---|---|---|
| S | Exponent | Fraction |

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

▪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 = 1263  1023

CSULB

3

# Single-Precision Range

- Exponents <u>00000000</u> and 11111111 reserved

- Smallest absolute value
  - Exponent:
    $\Rightarrow$ actual exponent = $00 \sim 1 \Rightarrow -126$
  - Fraction: 000…00 $\Rightarrow$ significand =
    $\pm 1.0 \times 2^{-126}$

$1 \times 2^{1-127}$

| S | Exponent | Fraction |
|---|----------|----------|

- Largest absolute value
  - Exponent: $1111 \sim 0$
    $\Rightarrow$ actual exponent = $254 - 127 = 127$
  - Fraction: 111…11 $\Rightarrow$ significand $\approx$
    $\pm 1.999\ldots \times 2^{127}$

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

# 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
  - $\pm 1.0 \times 2^{-1022} \approx \pm 2.2 \times 10^{-308}$

| S | Exponent | Fraction |
|---|----------|----------|

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

- Largest value
  - Exponent: 11111111110
    $\Rightarrow$ actual exponent = $2046 - 1023 = +1023$
  - Fraction: 111…11 $\Rightarrow$ significand $\approx$ 2.0
  - $\pm 2.0 \times 2^{+1023} \approx \pm 1.8 \times 10^{+308}$

# Denormal Numbers

▪Exponent = 000...0 ⇒ hidden bit is 0

- Denormal with fraction = 000...0

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

Two representations of 0.0!

| Sign | Exponent (e) | Fraction (f) | Value |
|------|------|------|------|
| 0 | 00···00 | 00···00 | +0 |
| 1 | 00···00 | 00···00 | −0 |

CSULB

# Denormal Numbers

▪Exponent = 000...0 ⇒ hidden bit is 0

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

- Smaller than normal numbers

  $0.0000\cdots1 \times 2^{-127}$
  23 bits

  - allow for gradual underflow, with diminishing precision

| Sign | Exponent (e) | Fraction (f) | Value |
|------|------|------|------|
| 0 | 00···00 | 00···01 ⋮ 11···11 | Positive Denormalized Real $0.f \times 2^{(-b+1)}$ |
| 1 | 00···00 | 00···01 ⋮ 11···11 | Negative Denormalized Real $-0.f \times 2^{(-b+1)}$ |

CSULB

5

# 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, 0.0 * ∞
  - Can be used in subsequent calculations

CSULB

# Floating-Point Precision

$$1.9999 \times \boxed{10^3}$$
$$= 1999.9 \Rightarrow 3 \text{ digits}$$

- Relative precision
  - all fraction bits are significant
  - Single: approx $\boxed{2^{23}}$
    - Equivalent to $\underline{\log 2^{23} = 23 \cdot \log 2 \approx 6}$ decimal digits of precision
  - Double: approx $2^{52}$
    - Equivalent to $\underline{\log 2^{52} \approx 16}$ decimal digits of precision

CSULB

6

# Floating-Point Example

- What number is represented by the single-precision float

  1 10000001 01000…00

  $1.\boxed{01}0000\sim 0$ (23)

  - S = $1$
  - Fraction = $010\sim 0 = 0100\sim 0_2$
  - Exponent = $1000001 = 129$
- x = $(-1)^1 \times (1 + 01_2) \times 2^{129-127}$

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

  $1.01_2 = 1.25_{10}$
  $+ 1\times 2^0 + 0\times 2^{-1} + 1\times 2^{-2}$

CSULB

# Floating-Point Example

- Represent –0.75
  - –0.75 = $-0.5 - 0.25 = -0.11_2 = -1.1 \times 2^{-1}$
  - S = $1$
  - Fraction = $10\sim 0$
  - Exponent = $x - 127 = -1$, $x = 126$
    - Single:
    - Double:
- Single:
- Double: 0

$\frac{9}{28}$

CSULB