

# **Chapter 3**

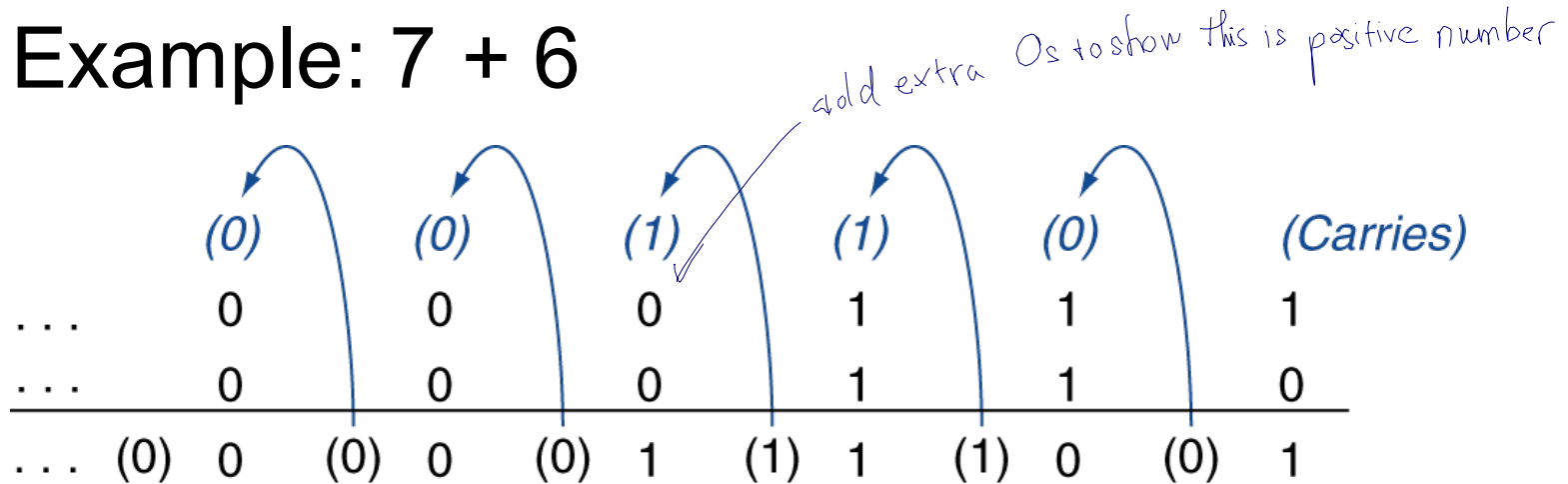
## **Arithmetic for Computers**

# Arithmetic for Computers

- Operations on integers
  - Addition and subtraction
  - Multiplication and division
  - Dealing with overflow
- Floating-point real numbers
  - Representation and operations

# Integer Addition

## ■ Example: $7 + 6$



## ■ Overflow if result out of range

- Adding positive (+) and negative (−) operands, no overflow
- Adding two positive (+ +) operands
  - Overflow if result sign is 1 (negative)
- Adding two negative (− −) operands
  - Overflow if result sign is 0 (positive)

$$\begin{array}{r} 2/6 \\ 03 \\ \hline 110 \\ \hline 6 = 110 \end{array}$$

# Integer Subtraction

- Add negation of second operand
- Example:  $7 - 6 = 7 + (-6)$

$$\begin{array}{r} +7: \quad 0000 \ 0000 \ \dots \ 0000 \ 0111 \\ -6: \quad 1111 \ 1111 \ \dots \ 1111 \ 1010 \\ \hline +1: \quad 0000 \ 0000 \ \dots \ 0000 \ 0001 \end{array}$$

*drop the overflow sign abs(7) > abs(-6).*

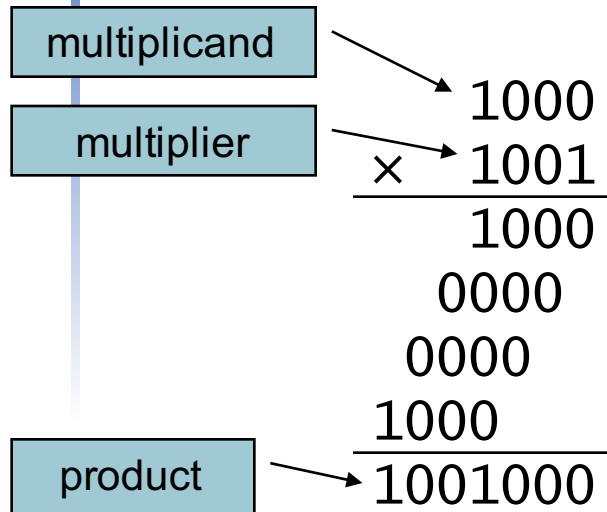
- Overflow if result out of range
  - Subtracting two positive (+ +) or two negative (– –) operands, no overflow
  - Subtracting positive (+) from negative (–) operand
    - Overflow if result sign is 0 (positive)
  - Subtracting negative (–) from positive (+) operand
    - Overflow if result sign is 1 (negative)

# Dealing with Overflow

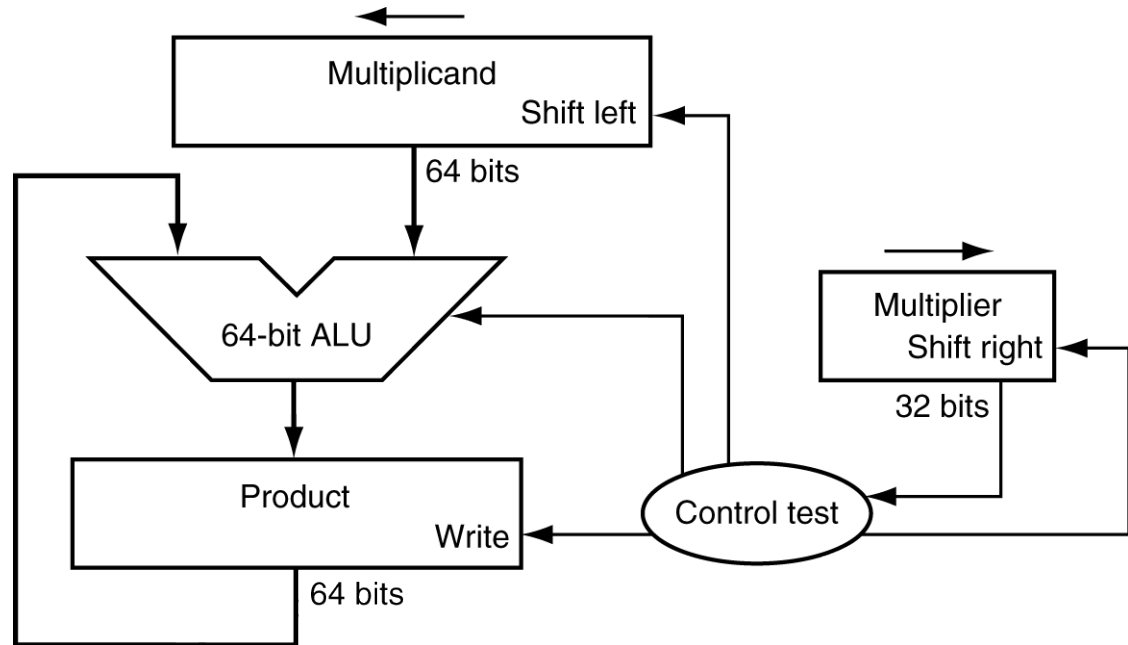
- Some languages (e.g., C) ignore overflow
- Other languages (e.g., Ada, Fortran) has overflow
- MIPS has overflow, require raising an exception on overflow
  - Use MIPS add, addi, sub instructions
  - On overflow, invoke exception handler
    - Save PC in exception program counter (EPC) register
    - Jump to predefined handler address
    - mfc0 (move from coprocessor reg) instruction can retrieve EPC value, to return after corrective action

# Multiplication

- Start with long-multiplication approach



Length of product is the sum of operand lengths



# Multiplication Hardware

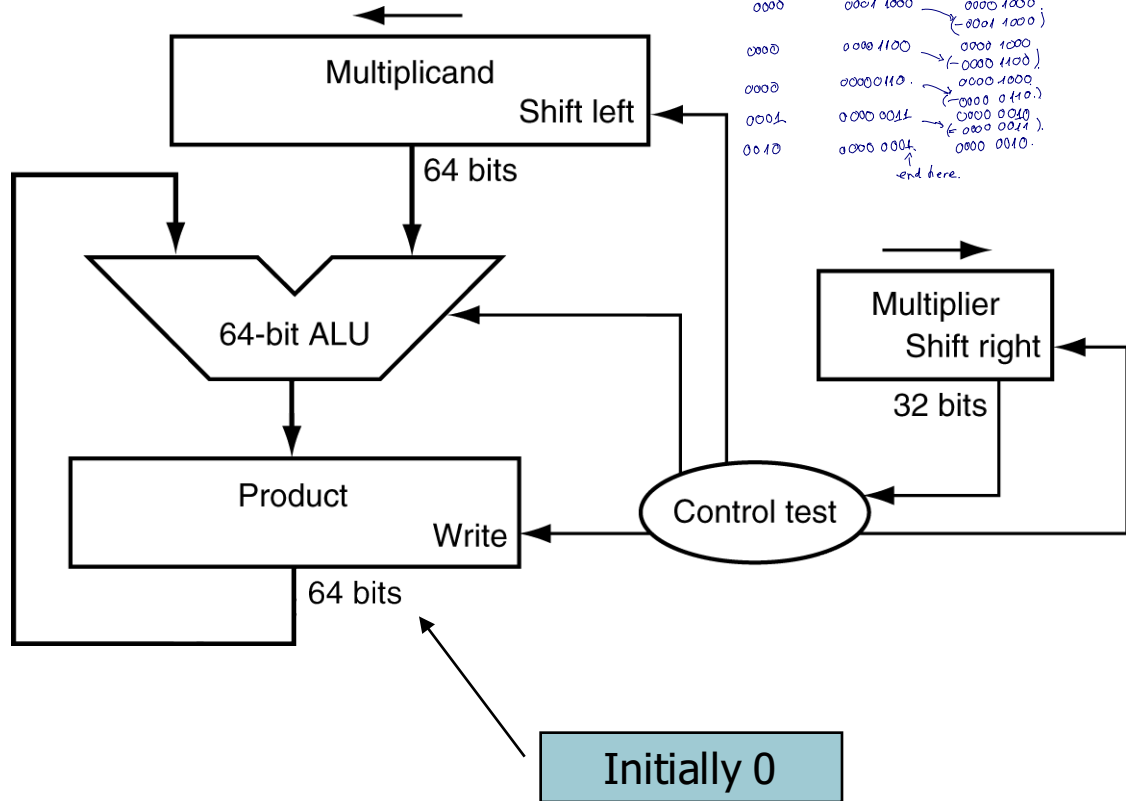
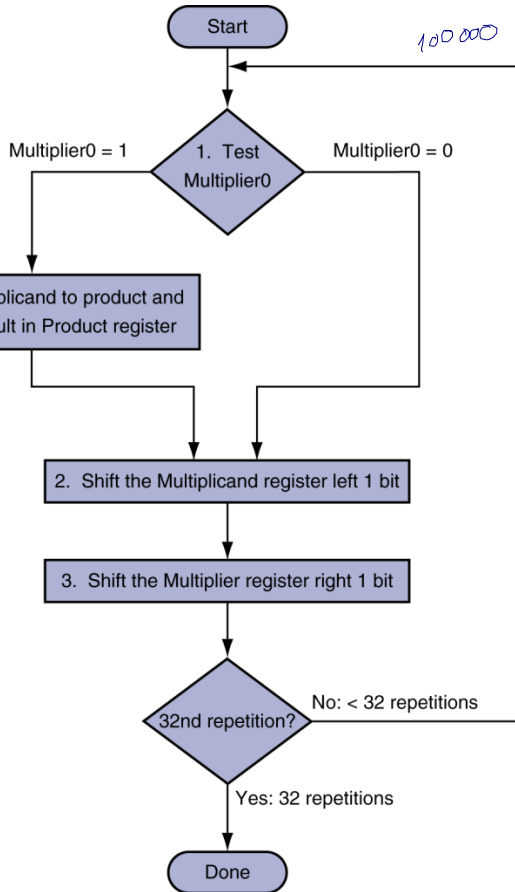
Multiplicand	Multiplier	Product
1000	1011	$  \begin{array}{r}  0000 \\  + 1000 \\  \hline  1000  \end{array}  $
10000	0001	$  \begin{array}{r}  1000 \\  + 10000 \\  \hline  11000  \end{array}  $
100000	0000	$  \begin{array}{r}  11000 \\  + 100000 \\  \hline  111000  \end{array}  $

$8 \times 5 = 24$   
 Multiplicand      multiplier      Product  
 $\begin{array}{r} 0011 \\ \rightarrow 0001 \\ \rightarrow 0000 \end{array}$ 
 $\begin{array}{r} 0001000 \\ 0001000 \\ 0001000 \end{array}$ 
 $\begin{array}{r} 00000000 \\ + 00010000 \\ + 00010000 \\ + 00010000 \\ \hline 00011000 \end{array}$ 
 $\begin{array}{r} 8 \\ \times 5 \\ \hline 16 \end{array}$   
 ← final product

8 ÷ 3

Quotient	Divisor	Remainder
0000	0011 0000	0000 1000 (-0001 0000)
0000	0001 1000	0000 1000 (-0001 1000)
0000	0000 1100	0000 1000 (-0000 1100)
0000	0000 0110	0000 1000 (-0000 1000)
0000	0000 0011	0000 0100 (-0000 0100)
0010	0000 0001	0000 0010 (-0000 0011)
	0000 0000	0000 0010

↑  
end here.



# MIPS Multiplication

- Two 32-bit registers for product
    - HI: most-significant 32 bits
    - LO: least-significant 32-bits
  - Instructions
    - `mult rs, rt` / `multu rs, rt`
      - 64-bit product in HI/LO ← *special registers not accessible by name but rather by mfhil or mflh.*
    - `mfhi rd` / `mflo rd`
      - Move from HI/LO to rd
      - Can test HI value to see if product overflows 32 bits
    - `mul rd, rs, rt`
      - Least-significant 32 bits of product → rd
- similar to mult but only use least-significant 32 bit (i.e LO)*



# Division

Dividend: 8    Divisor: 3    Quotient: ?    Remainder: ?

Step 1:  $8 \div 3$  (left)

Quotient	Divisor	Remainder
0000	0011 0000	0000 1000
		<del>0011 0000</del>
0000	0001 1000	0000 1000
		<del>0001 1000</del>
0000	0000 1100	0000 1000
		<del>0000 1100</del>
0000	0000 0110	0000 1000
		<del>0000 0110</del>

# minus the divisor - (negative, raise previous value.)

# Shift Divisor 1 bit to the right

Quotient 1    left

# — Divisor    right

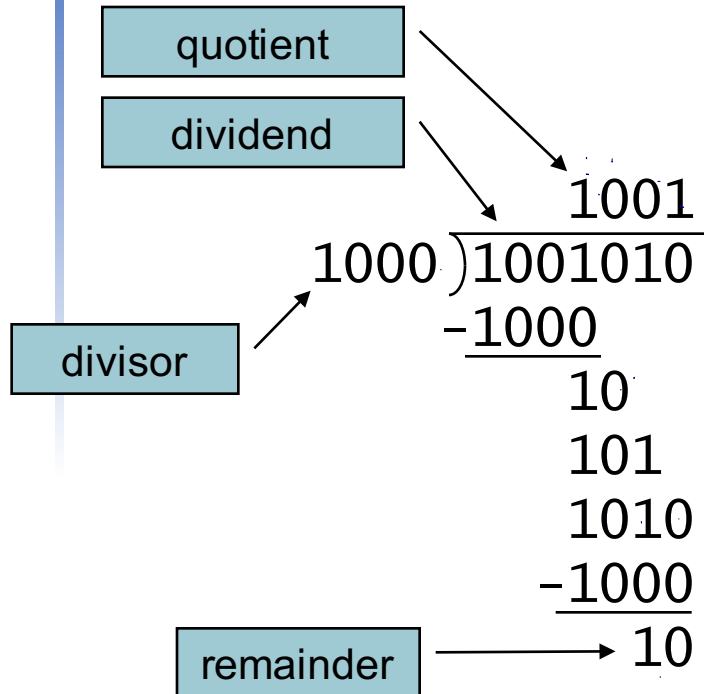
# Quotient —    left

0001    0000 0111    0000 0000 0 ← keep. # add quotients with 1

0001    0000 0000    0000 0000

or add with 1

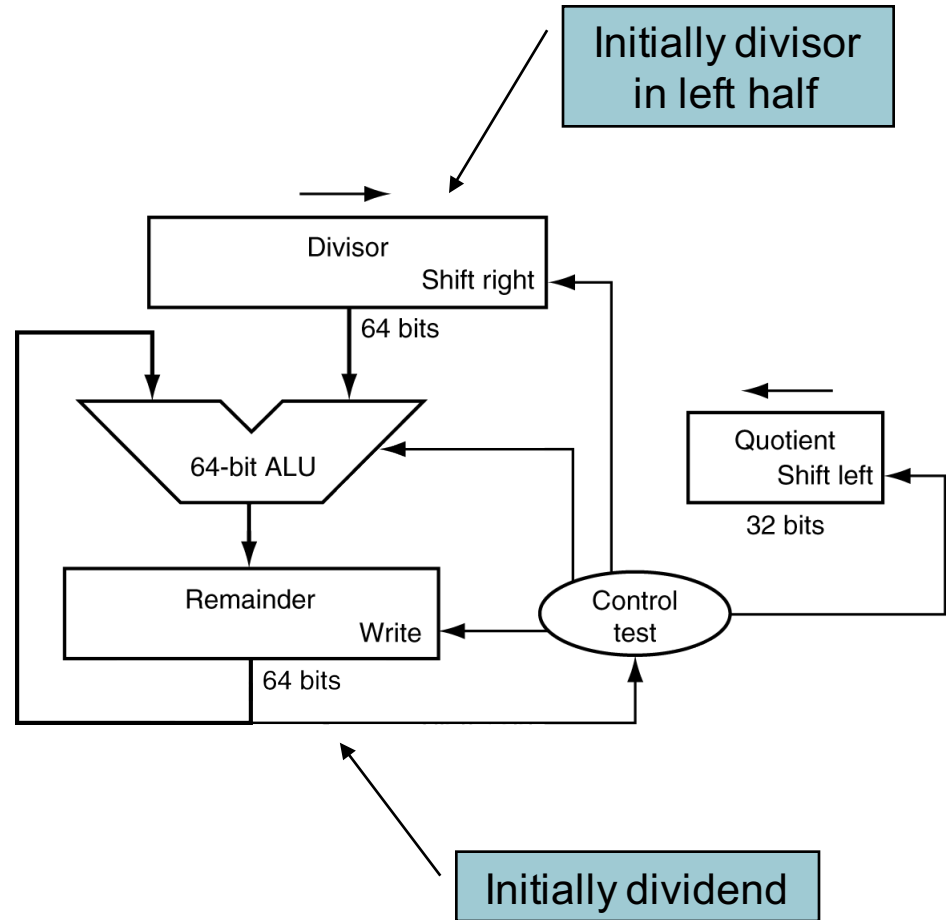
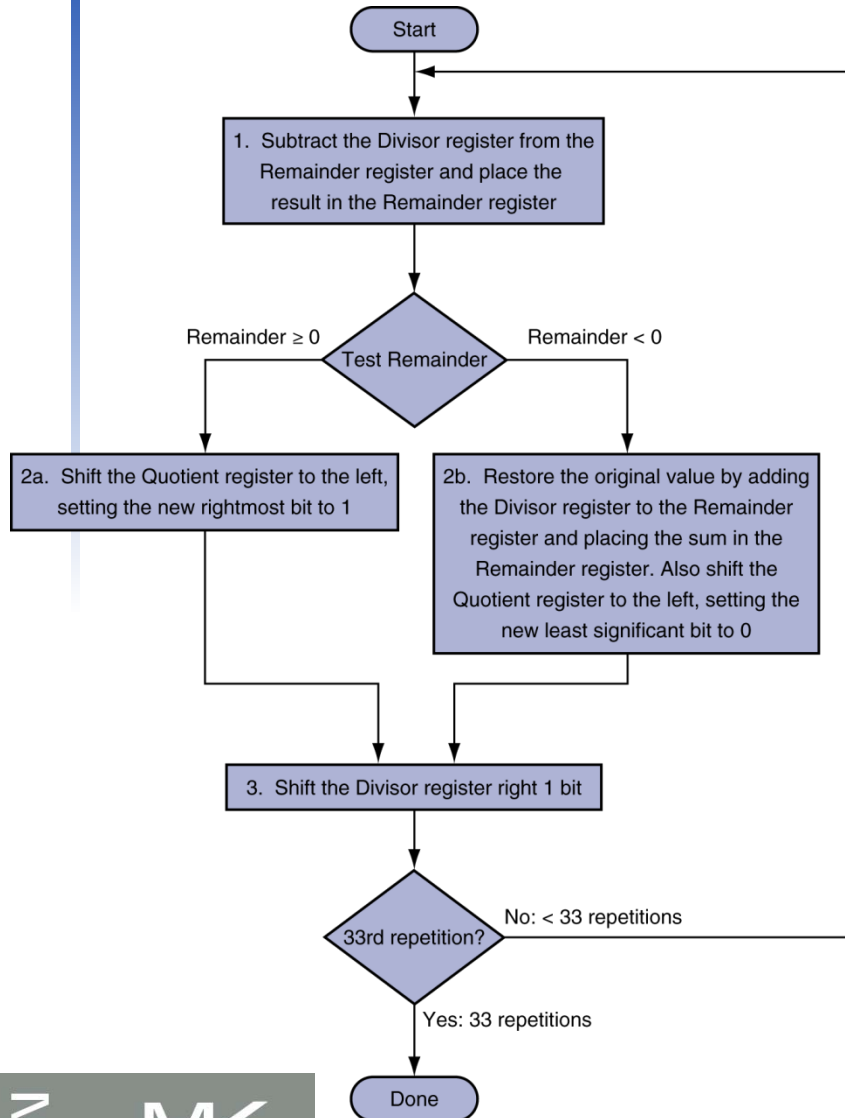
Check for 0 d



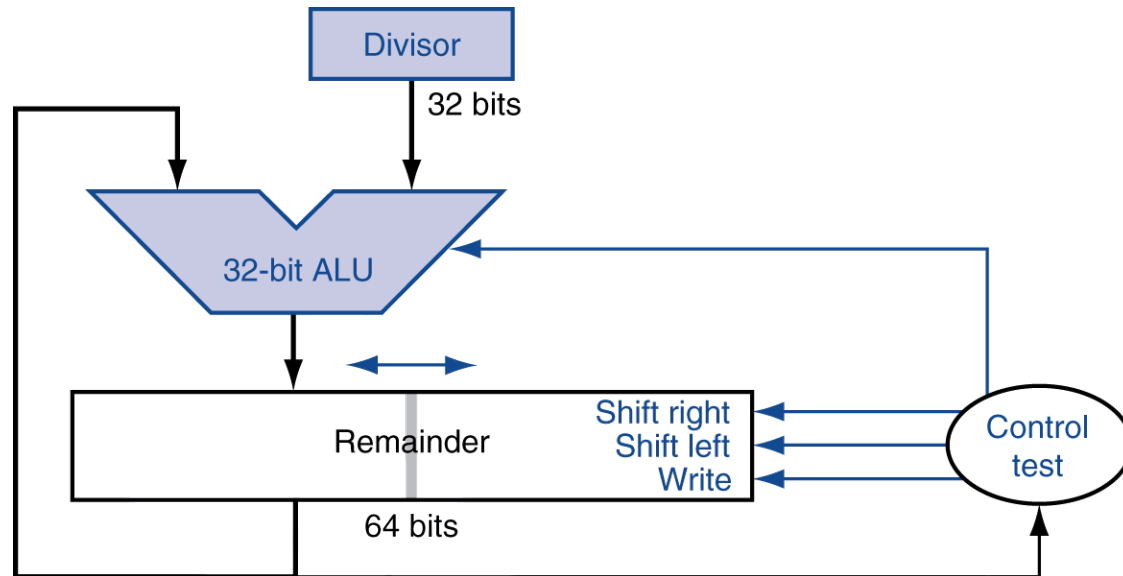
$n$ -bit operands yield  $n$ -bit quotient and remainder

- Check for 0 divisor
- Long division approach
  - If divisor  $\leq$  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

# Division Hardware



# Optimized Divider



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

# MIPS Division

- Use HI/LO registers for result
  - HI: 32-bit remainder
  - LO: 32-bit quotient
- Instructions
  - `div rs, rt` / `divu rs, rt`

# Floating Point

- Representation for non-integral numbers
  - Including very small and very large numbers
- Like scientific notation
  - $-1.34 \times 10^{56}$  ← normalized
  - $+0.002 \times 10^{-4}$  ← not normalized
  - $+987.02 \times 10^9$  ← not normalized
- In binary
  - $\pm 1.xxxxxxx_2 \times 2^{yyyy}$
- Types `float` and `double` in C

# 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

single: 8 bits

double: 11 bits

single: 23 bits

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 1 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
- Largest value
  - exponent: 11111110  
 $\Rightarrow$  actual exponent =  $254 - 127 = +127$
  - Fraction: 111...11  $\Rightarrow$  significand  $\approx 2.0$



# 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
- Largest value
  - Exponent: 11111111110  
 $\Rightarrow$  actual exponent =  $2046 - 1023 = +1023$
  - Fraction: 111...11  $\Rightarrow$  significand  $\approx 2.0$

# Floating-Point Precision

- Relative precision
  - all fraction bits are significant
  - Single: approx  $2^{-23}$
  - Double: approx  $2^{-52}$

# Floating-Point Example

- Represent  $-0.75$ 
  - $-0.75 = (-1)^1 \times 1.1_2 \times 2^{-1}$
  - $S = 1$
  - Fraction =  $1000\dots00_2$
  - Exponent =  $-1 + \text{Bias}$ 
    - Single:  $-1 + 127 = 126 = 01111110_2$
    - Double:  $-1 + 1023 = 1022 = 01111111110_2$
- Single:  $10111111101000\dots00$
- Double:  $101111111111101000\dots00$

# Floating-Point Example

- What number is represented by the single-precision float

11000000101000...00

- $S = 1$

- Fraction =  $01000...00_2$

- Exponent =  $10000001_2 = 129$

- $$\begin{aligned}x &= (-1)^1 \times (1 + 01_2) \times 2^{(129 - 127)} \\&= (-1) \times 1.25 \times 2^2 \\&= -5.0\end{aligned}$$

# 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
- Denormal with fraction = 000...0

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

Two representations  
of 0.0!



# Floating-Point Example

1. Use the MIPS algorithm in this chapter to get the binary sum for these two decimal numbers, verify your result is correct.  $5_{10} + 9_{10} = ?_2$
2. Use the MIPS algorithm in this chapter to get the binary difference for these two decimal numbers, verify your result is correct.  $5_{10} - 9_{10} = ?_2$
3. Use the MIPS algorithm in this chapter to get the binary product for these two decimal numbers, verify your result is correct.  $5_{10} \times 9_{10} = ?_2$
4. Use the MIPS algorithm in this chapter to get the binary quotient and binary remainder for these two decimal numbers, verify your result is correct.  $9_{10} / 2_{10} = ?_2$