# Chapter 3

## Arithmetic for Computers

(Revised)
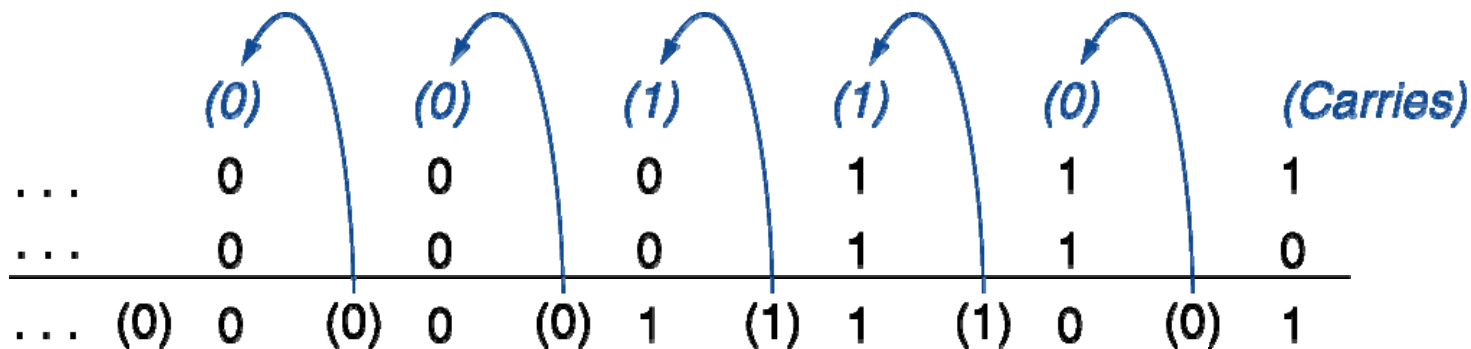
# Arithmetic for Computers

- ## 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 +ve and –ve operands, no overflow

  - Adding two +ve operands

    - Overflow if result sign is 1

  - Adding two –ve operands

    - Overflow if result sign is 0

# Integer Subtraction

- Add negation of second operand

- Example: 7 – 6 = 7 + (–6)

```
+7:        0000 0000 … 0000 0111
–6:        1111 1111 … 1111 1010
+1:        0000 0000 … 0000 0001
```

- Overflow if result out of range

  - Subtracting two +ve or two –ve operands, no overflow

  - Subtracting +ve from –ve operand
    - Overflow if result sign is 0

  - Subtracting –ve from +ve operand
    - Overflow if result sign is 1

- Consider the operations A + B, and A – B

  - Can overflow occur if B is 0 ?
  - Can overflow occur if A is 0 ?

# Effects of Overflow

- An exception (interrupt) occurs
  - Control jumps to predefined address for exception
  - Interrupted address is saved for possible resumption

- Details based on software system / language

- Don't always want to detect overflow
  — new MIPS instructions: `addu, addiu, subu`

  *note:* `addiu` *still sign-extends!*
  *note:* `sltu, sltiu` *for unsigned comparisons*

  **`addiu` sign-extends its 16-bit immediate to 32-bit, when performing addition with a 32-bit register.**
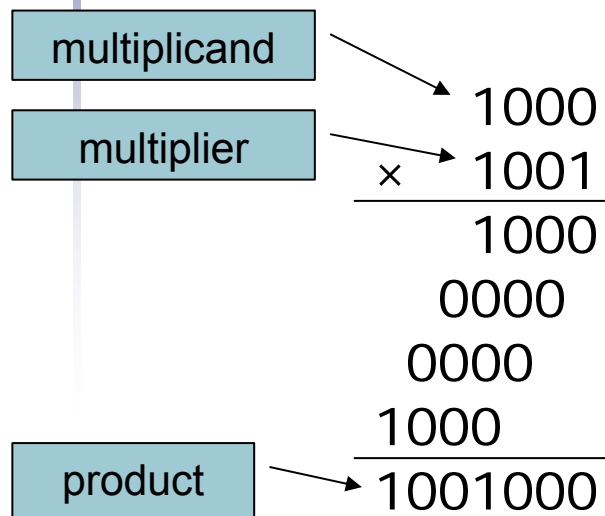
  **This is because in its design, the immediate can be negative. So `addiu` is only an `addi` counterpart that *ignores* *overflow detection*.**
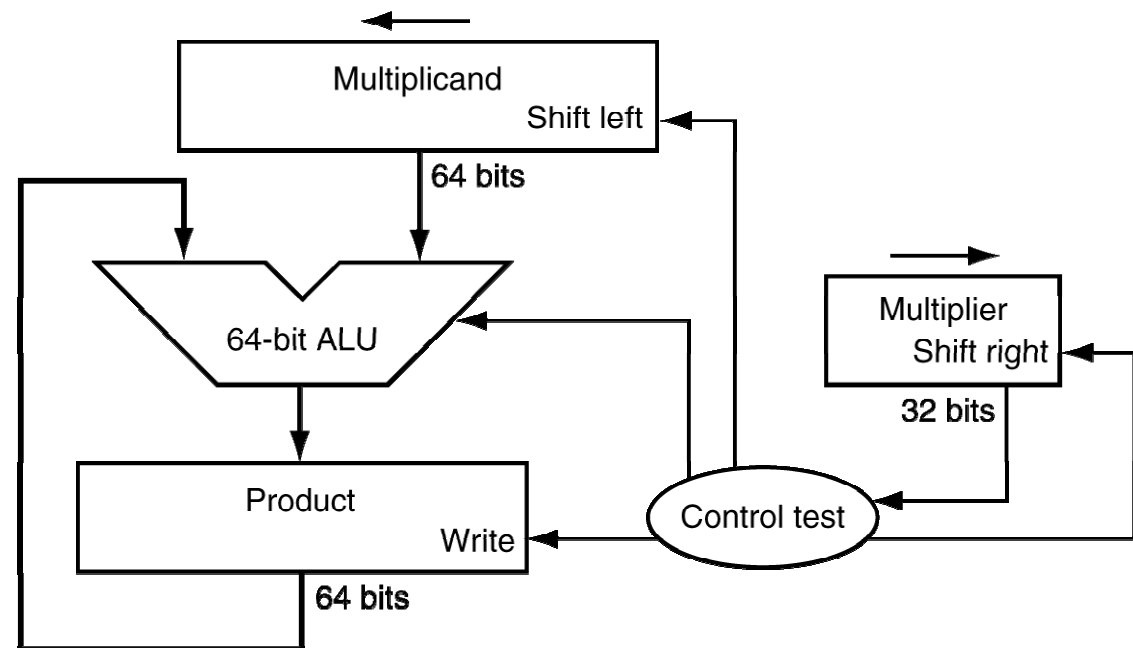
# Arithmetic for Multimedia

- Graphics and media processing operates on vectors of 8-bit and 16-bit data
    - Use 64-bit adder, with partitioned carry chain
        - Operate on 8×8-bit, 4×16-bit, or 2×32-bit vectors
    - SIMD (single-instruction, multiple-data)
- Saturating operations
    - On overflow, result is largest representable value
        - c.f. 2s-complement modulo arithmetic
    - E.g., clipping in audio, saturation in video
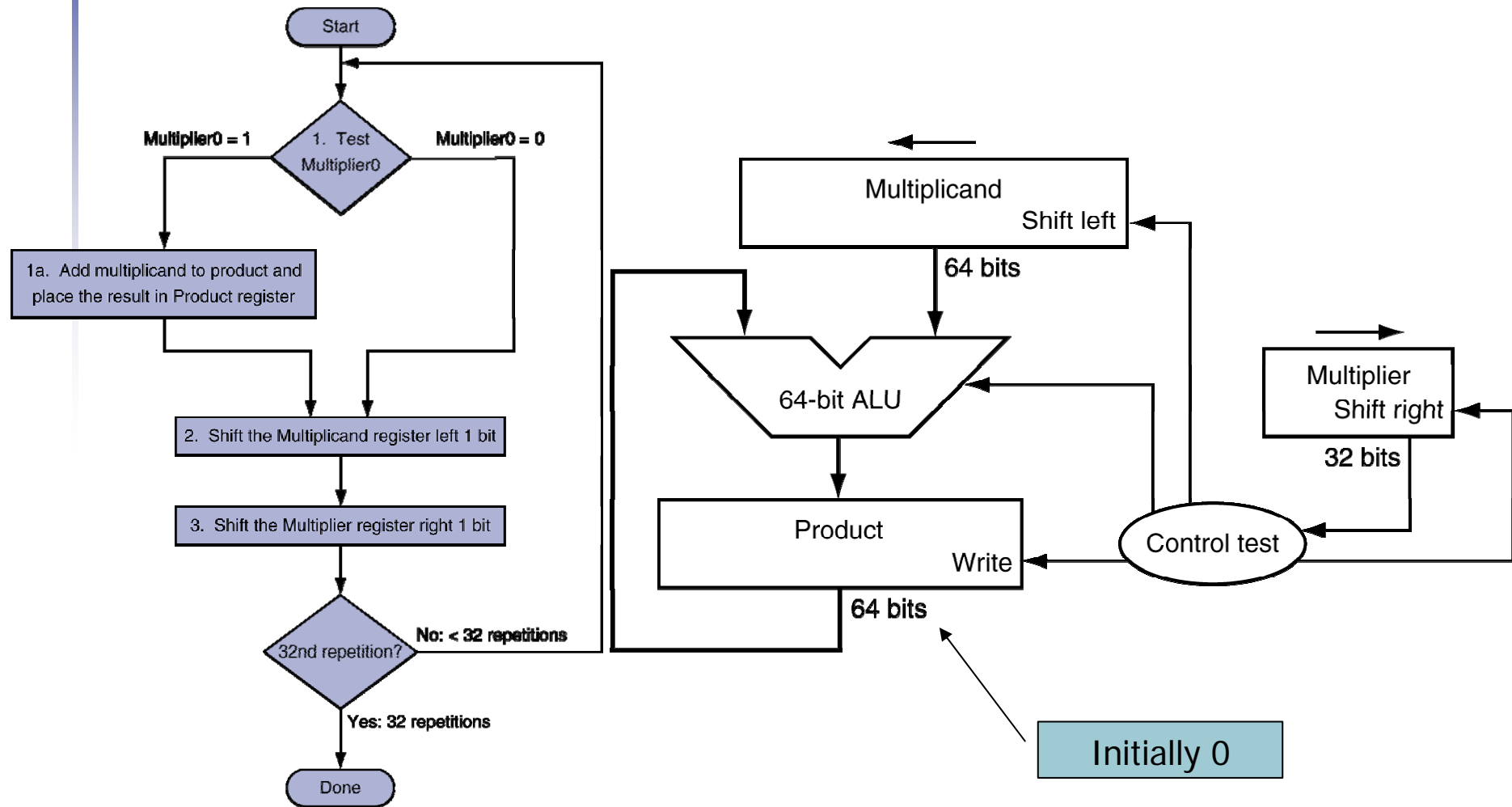
# Multiplication

- Start with long-multiplication approach

multiplicand

multiplier

```
      1000
  ×   1001
      1000
     0000
    0000
   1000
   1001000
```

product

Length of product is
the sum of operand
lengths

# Multiplication Hardware



**Start**

1. Test Multiplier0

Multiplier0 = 1 → 1a. Add multiplicand to product and place the result in Product register

Multiplier0 = 0

2. Shift the Multiplicand register left 1 bit

3. Shift the Multiplier register right 1 bit

32nd repetition?

No: < 32 repetitions

Yes: 32 repetitions

**Done**

Multiplicand — Shift left — 64 bits

64-bit ALU

Product — Write — 64 bits

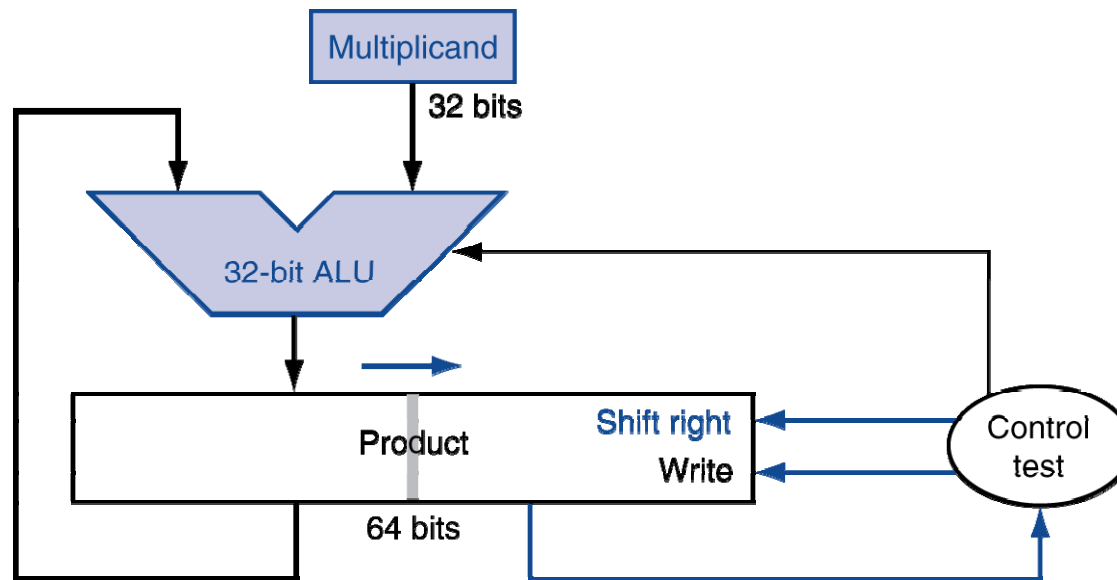Multiplier — Shift right — 32 bits

Control test

Initially 0

Example (p.234)
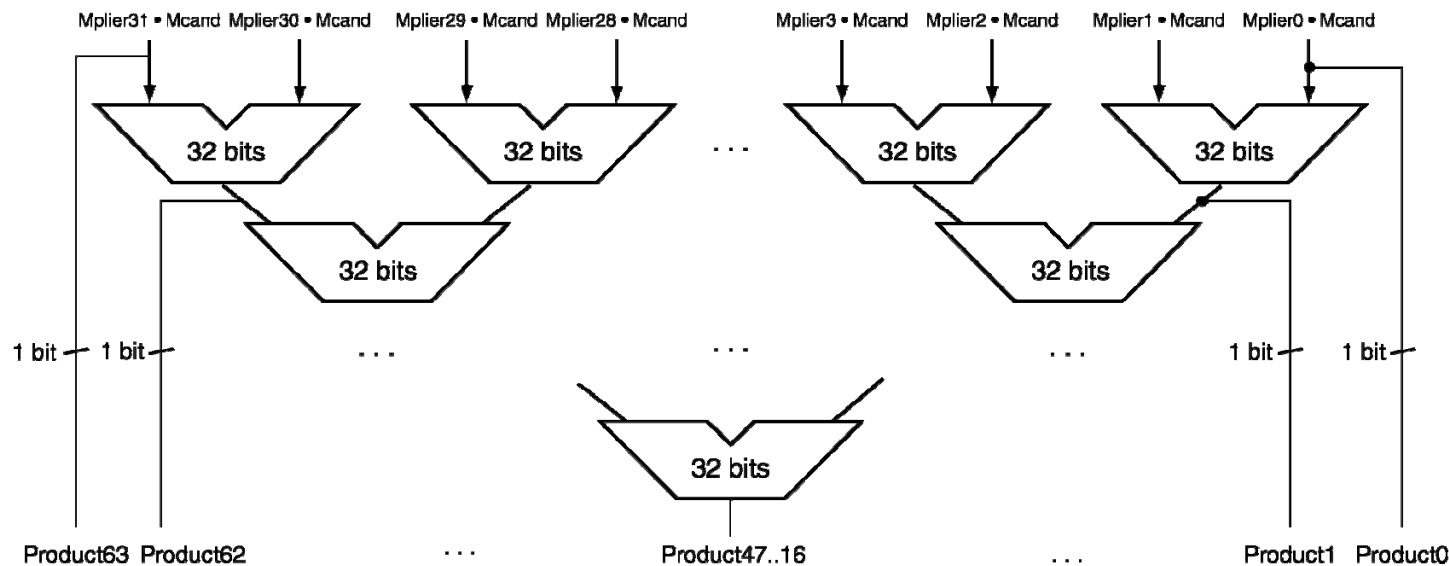
# Optimized Multiplier

- Perform steps in parallel: add/shift



- One cycle per partial-product addition
  - That's ok, if frequency of multiplications is low

# Faster Multiplier

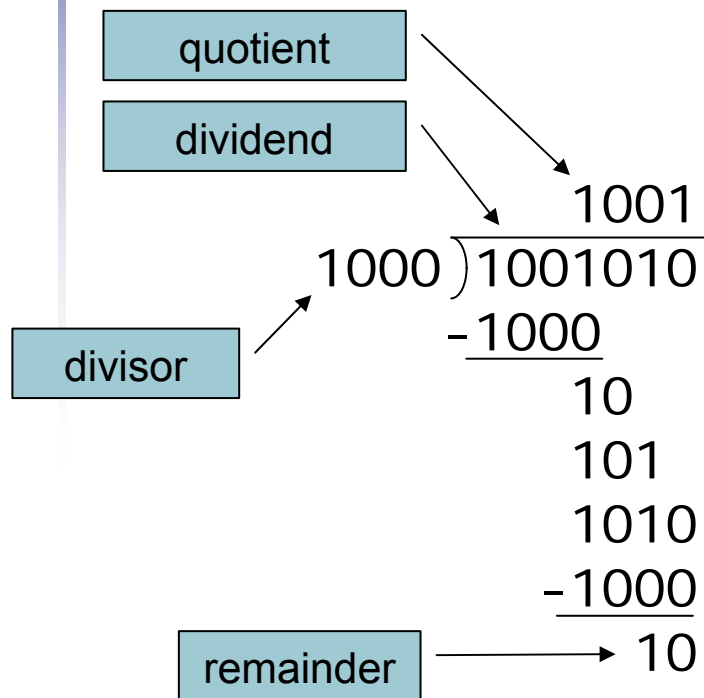- Uses multiple adders
  - Cost/performance tradeoff



- Can be pipelined
  - Several multiplication performed in parallel

# Multiply in MIPS

- MIPS provides a separate pair of 32-bit register ($Hi$ and $Lo$) to contain the 64-bit product

- MIPS instructions:
  - `mult`: *multiply*
  - `multu`: *multiply unsigned*
  - `mflo`: *move from lo*
  - `mfhi`: *move from hi*

# Division

quotient

dividend

```
            1001
1000 ) 1001010
      -1000
          10
         101
        1010
       -1000
          10
```

divisor

remainder

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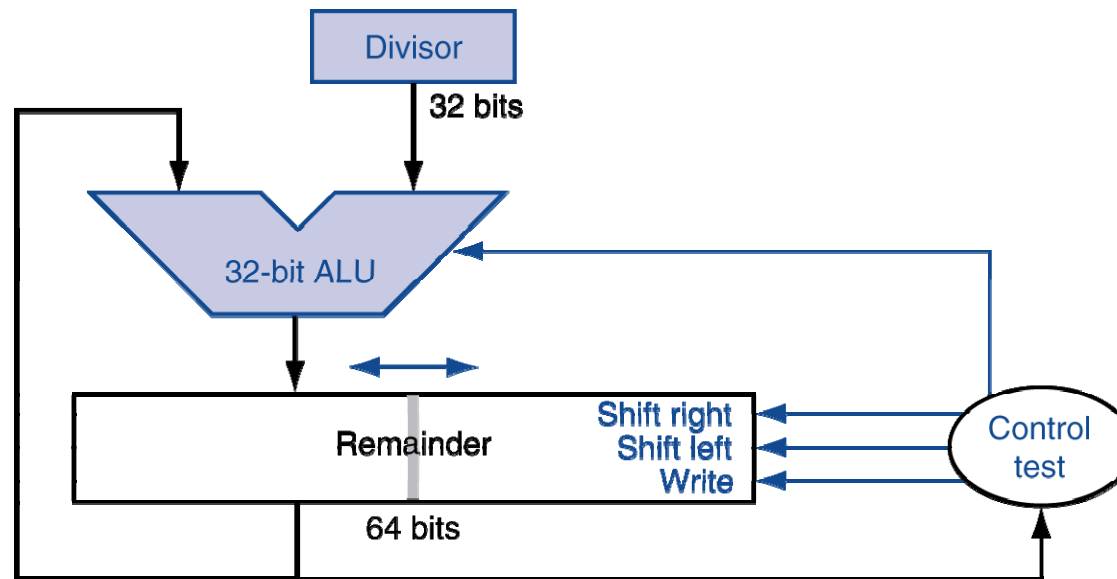
- Check for 0 divisor
- Long division approach
  - If divisor ≤ dividend bits
    - 1 bit in quotient, subtract
  - Otherwise
    - 0 bit in quotient, bring down next dividend bit
- Restoring division
  - Do the subtract, and if remainder goes < 0, add divisor back
- Signed division
  - Divide using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

# 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

Remainder          Write
64 bits

Quotient          Shift left
32 bits

Control test

Initially dividend

Example (p.240)

# Optimized Divider



- 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 (e.g. SRT devision) generate multiple quotient bits per step
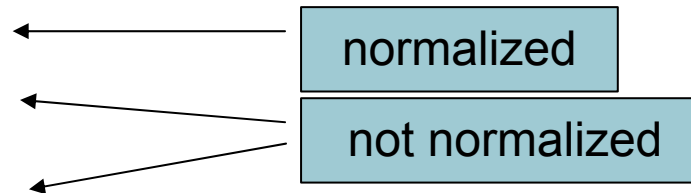  - Still require multiple steps

# Divide in MIPS

■ MIPS provides a separate pair of 32-bit register ($Hi$ and $Lo$) to contain the pair of resulting remainder & quotient

■ MIPS instructions:

- `div`: *divide*
- `divu`: *divide unsigned*
- `mflo`: *move from lo*
- `mfhi`: *move from hi*

# Floating Point

- Representation for non-integral numbers
    - Including very small and very large numbers
- Like scientific notation
    - $-2.34 \times 10^{56}$
    - $+0.002 \times 10^{-4}$
    - $+987.02 \times 10^{9}$

    normalized

    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
    - Two representations
        - Single precision (32-bit)
        - Double precision (64-bit)

# IEEE Floating-Point Format

single: 8 bits single: 23 bits
double: 11 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 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
  - $\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$

- Largest value
  - exponent: 11111110
    $\Rightarrow$ actual exponent = 254 – 127 = +127
  - Fraction: 111…11 $\Rightarrow$ significand = 1.111…11
  - $\pm 1.\,111…11 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$

# 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 $\approx$ 1. 111…11

  - $\pm 1. 111…11 \times 2^{+1023} \approx \pm 1.8 \times 10^{+308}$

# Floating-Point Example

- Represent –0.75
  - $-0.75 = (-1)^1 \times 1.1_2 \times 2^{-1}$
  - $S = 1$
  - Fraction = $1000\ldots00_2$
  - Exponent = –1 + Bias
    - Single: –1 + 127 = 126 = $01111110_2$
    - Double: –1 + 1023 = 1022 = $01111111110_2$
- Single: $1\,01111110\,1000\ldots00$
- Double: $1\,01111111110\,1000\ldots00$

# Floating-Point Example

- What number is represented by the single-precision float

  1 1000000 1 01000...00

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

- $x = (-1)^1 \times (1 + 0.01_2) \times 2^{(129 - 127)}$

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

  $= -5.0$

# Denormalized Numbers

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

$$x = (-1)^S \times (0 + \text{Fraction}) \times 2^{-126}, \text{ for single precision}$$

$$x = (-1)^S \times (0 + \text{Fraction}) \times 2^{-1022}, \text{ for double precision}$$

- Smaller than normal numbers
  - allow for gradual underflow, with diminishing precision

# 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

# IEEE 754 floating-point standard

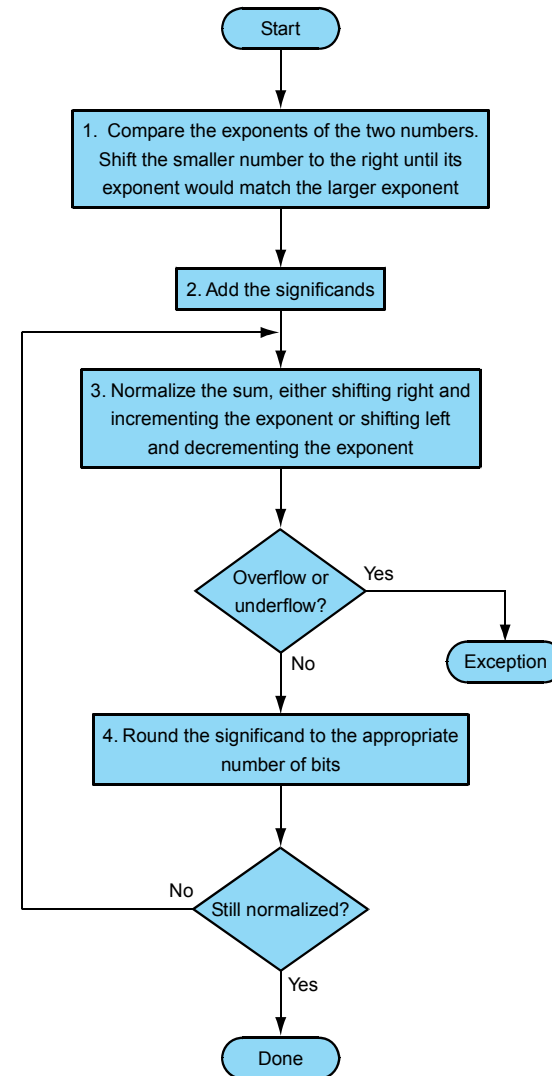| Single precision | | Double precision | | Object represented |
|---|---|---|---|---|
| Exponent | Fraction | Exponent | Fraction | |
| 0 | 0 | 0 | 0 | 0 |
| 0 | nonzero | 0 | nonzero | ± denormalized number |
| 1-254 | anything | 1-2046 | anything | ± floating-point number |
| 255 | 0 | 2047 | 0 | ± infinity |
| 255 | nonzero | 2047 | nonzero | NaN (Not a number) |

# Floating-Point Addition

- Consider a 4-digit decimal example
    - $9.999 \times 10^1 + 1.610 \times 10^{-1}$
- 1. Align decimal points
    - Shift number with smaller exponent
    - $9.999 \times 10^1 + 0.016 \times 10^1$
- 2. Add significands
    - $9.999 \times 10^1 + 0.016 \times 10^1 = 10.015 \times 10^1$
- 3. Normalize result & check for over/underflow
    - $1.0015 \times 10^2$
- 4. Round and renormalize if necessary
    - $1.002 \times 10^2$

# Floating-Point Addition

- Now consider a 4-digit binary example
  - $1.000_2 \times 2^{-1} + -1.110_2 \times 2^{-2}$ (0.5 + –0.4375)
- 1. Align binary points
  - Shift number with smaller exponent
  - $1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1}$
- 2. Add significands
  - $1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1} = 0.001_2 \times 2^{-1}$
- 3. Normalize result & check for over/underflow
  - $1.000_2 \times 2^{-4}$, with no over/underflow
- 4. Round and renormalize if necessary
  - $1.000_2 \times 2^{-4}$ (no change)  = 0.0625

# FP Adder Hardware

- Much more complex than integer adder

- Doing it in one clock cycle would take too long
  - Much longer than integer operations
  - Slower clock would penalize all instructions

- FP adder usually takes several cycles
  - Can be pipelined

Start

1. Compare the exponents of the two numbers. Shift the smaller number to the right until its exponent would match the larger exponent

2. Add the significands

3. Normalize the sum, either shifting right and incrementing the exponent or shifting left and decrementing the exponent

Overflow or underflow?

Yes → Exception

No

4. Round the significand to the appropriate number of bits

Still normalized?

No

Yes

Done

# FP Adder Hardware
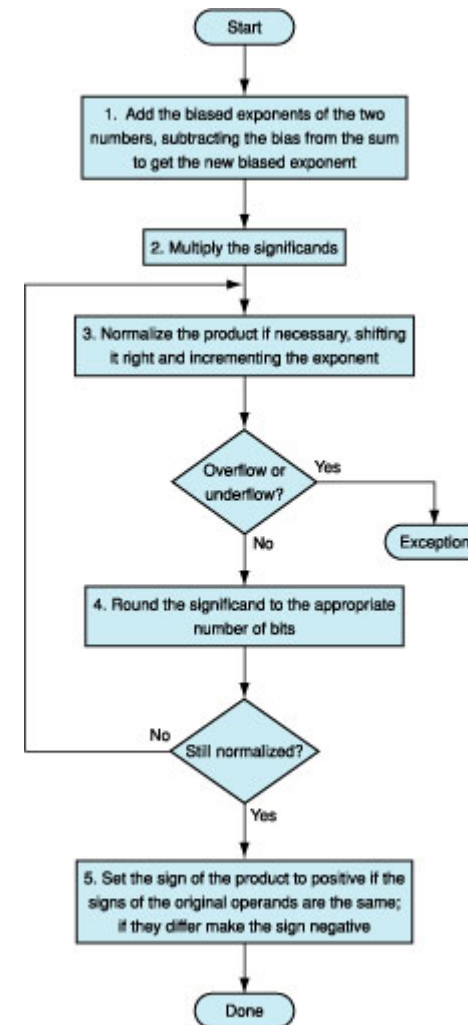
# Floating-Point Multiplication

- Consider a 4-digit decimal example
  - $1.110 \times 10^{10} \times 9.200 \times 10^{-5}$
- 1. Add exponents
  - For biased exponents, subtract bias from sum
  - New exponent = 10 + –5 = 5
- 2. Multiply significands
  - $1.110 \times 9.200 = 10.212 \implies 10.212 \times 10^{5}$
- 3. Normalize result & check for over/underflow
  - $1.0212 \times 10^{6}$
- 4. Round and renormalize if necessary
  - $1.021 \times 10^{6}$
- 5. Determine sign of result from signs of operands
  - $+1.021 \times 10^{6}$

# Floating-Point Multiplication

- Now consider a 4-digit binary example
  - $1.000_2 \times 2^{-1} \times -1.110_2 \times 2^{-2}$ ($0.5 \times -0.4375$)
- 1. Add exponents
  - Unbiased: $-1 + -2 = -3$
  - Biased: $(-1 + 127) + (-2 + 127) - 127 = -3 + 127$
- 2. Multiply significands
  - $1.000_2 \times 1.110_2 = 1.1102 \Rightarrow 1.110_2 \times 2^{-3}$
- 3. Normalize result & check for over/underflow
  - $1.110_2 \times 2^{-3}$ (no change) with no over/underflow
- 4. Round and renormalize if necessary
  - $1.110_2 \times 2^{-3}$ (no change)
- 5. Determine sign: +ve $\times$ –ve $\Rightarrow$ –ve
  - $-1.110_2 \times 2^{-3} = -0.21875$

# FP Arithmetic Hardware

- FP multiplier is of similar complexity to FP adder
  - But uses a multiplier for significands instead of an adder
- FP arithmetic hardware usually does
  - Addition, subtraction, multiplication, division, reciprocal, square-root
  - FP $\leftrightarrow$ integer conversion
- Operations usually takes several cycles
  - Can be pipelined

Start

1. Add the biased exponents of the two numbers, subtracting the bias from the sum to get the new biased exponent

2. Multiply the significands

3. Normalize the product if necessary, shifting it right and incrementing the exponent

Overflow or underflow? — Yes → Exception

No

4. Round the significand to the appropriate number of bits

Still normalized? — No

Yes

5. Set the sign of the product to positive if the signs of the original operands are the same; if they differ make the sign negative

Done

# Accurate Arithmetic

- IEEE Std 754 specifies additional rounding control
  - Extra bits of precision (guard, round, sticky)
  - Choice of rounding modes
  - Allows programmer to fine-tune numerical behavior of a computation
- Not all FP units implement all options
  - Most programming languages and FP libraries just use defaults
- Trade-off between hardware complexity, performance, and market requirements

# Interpretation of Data

**The BIG Picture**

- Bits have no inherent meaning
  - Interpretation depends on the instructions applied

- Computer representations of numbers
  - Finite range and precision
  - Need to account for this in programs

# Associativity

- Parallel programs may interleave operations in unexpected orders

  - Assumptions of associativity may fail

| | | (x+y)+z | x+(y+z) |
|---|---|---|---|
| x | -1.50E+38 | | -1.50E+38 |
| y | 1.50E+38 | 0.00E+00 | |
| z | 1.0 | 1.0 | 1.50E+38 |
| | | 1.00E+00 | 0.00E+00 |

- Need to validate parallel programs under varying degrees of parallelism

# Right Shift and Division

- Left shift by $i$ places multiplies an integer by $2^i$

- Right shift divides by $2^i$?
  - Only for unsigned integers

- For signed integers
  - Arithmetic right shift: replicate the sign bit
  - e.g., –5 / 4
    - $11111011_2 >> 2 = 11111110_2 = -2$
    - Rounds toward $-\infty$
  - c.f. $11111011_2 >> 2 = 00111110_2 = +62$

# Concluding Remarks

- ISAs support arithmetic
  - Signed and unsigned integers
  - Floating-point approximation to reals

- Bounded range and precision
  - Operations can overflow and underflow