

# Chapter 3

## Arithmetic for Computers

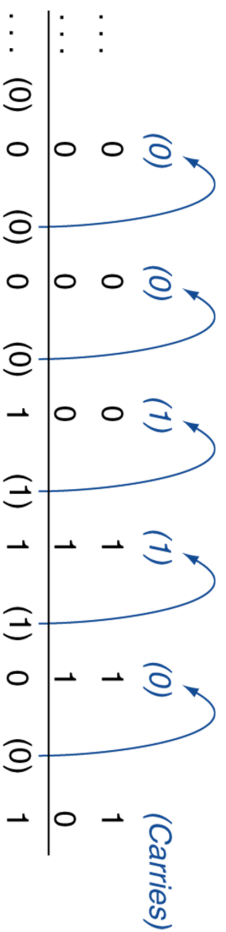
### Arithmetic for Computers

§3.1 Introduction

- 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



# Dealing with Overflow

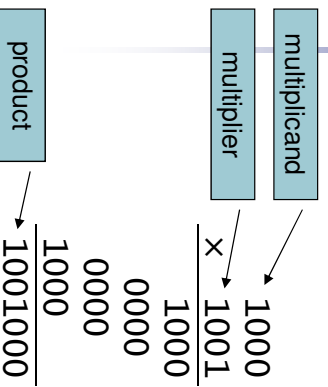
- Some languages (e.g., C) ignore overflow
  - Use MIPS addu, addui, subu instructions
- Other languages (e.g., Ada, Fortran) require raising an exception
  - Use MIPS add, addi, sub instructions
  - On overflow, invoke exception handler
  - Save PC in exception program counter (EPC) register
  - Jump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



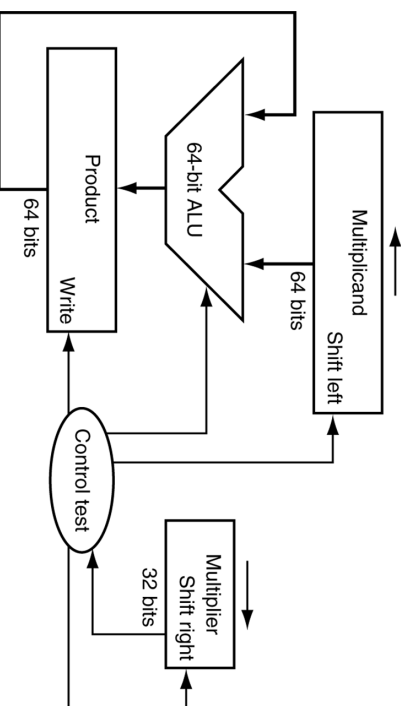
# Multiplication

- Start with long-multiplication approach

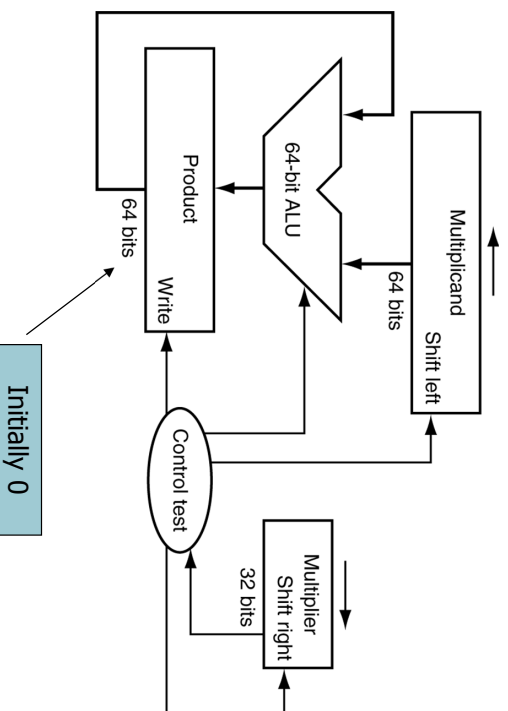
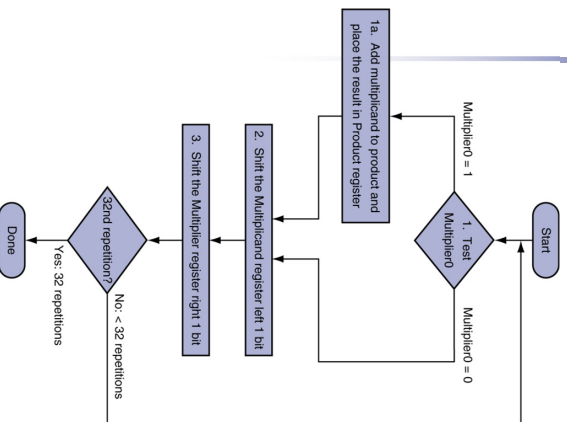
§3.3 Multiplication



Length of product is the sum of operand lengths



# Multiplication Hardware



# 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
  - `mflo rd` / `mftlo 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



# 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



# Floating Point

§3.5 Floating Point

- Representation for non-integral numbers
  - Including very small and very large numbers
- Like scientific notation
  - $-2.34 \times 10^{56}$  → normalized
  - $+0.002 \times 10^{-4}$  → not normalized
  - $+987.02 \times 10^9$  → not normalized
- In binary
  - $\pm 1.XXXXXXXXX_2 \times 2^{YYYY}$
- 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                      single: 23 bits  
double: 11 bits                    double: 52 bits



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

- 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
    - ⇒ actual exponent =  $1 - 127 = -126$
  - Fraction: 000...00 ⇒ significand = 1.0
  - $\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$
- Largest value
  - exponent: 11111110
    - ⇒ actual exponent =  $254 - 127 = +127$
  - Fraction: 111...11 ⇒ significand  $\approx 2.0$
  - $\pm 2.0 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$



# Double-Precision Range

- Exponents 0000...00 and 1111...11 reserved
- Smallest value
  - Exponent: 00000000001
    - ⇒ actual exponent =  $1 - 1023 = -1022$
  - Fraction: 000...00 ⇒ significand = 1.0
  - $\pm 1.0 \times 2^{-1022} \approx \pm 2.2 \times 10^{-308}$
- Largest value
  - Exponent: 11111111110
    - ⇒ actual exponent =  $2046 - 1023 = +1023$
  - Fraction: 111...11 ⇒ 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 =  $1000\dots00_2$
  - Exponent =  $-1 + \text{Bias}$ 
    - Single:  $-1 + 127 = 126 = 01111110_2$
    - Double:  $-1 + 1023 = 1022 = 011111111110_2$
  - Single:  $1011111101000\dots00$
  - Double:  $1011111111101000\dots00$





## Floating-Point Example

- What number is represented by the single-precision float

11000000101000...00

- S = 1
- Fraction = 01000...00<sub>2</sub>
- Exponent = 10000001<sub>2</sub> = 129
- X = (-1)<sup>1</sup> × (1 + 01<sub>2</sub>) × 2<sup>(129 - 127)</sup>  
= (-1) × 1.25 × 2<sup>2</sup>  
= -5.0



## Floating-Point Addition

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



# 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

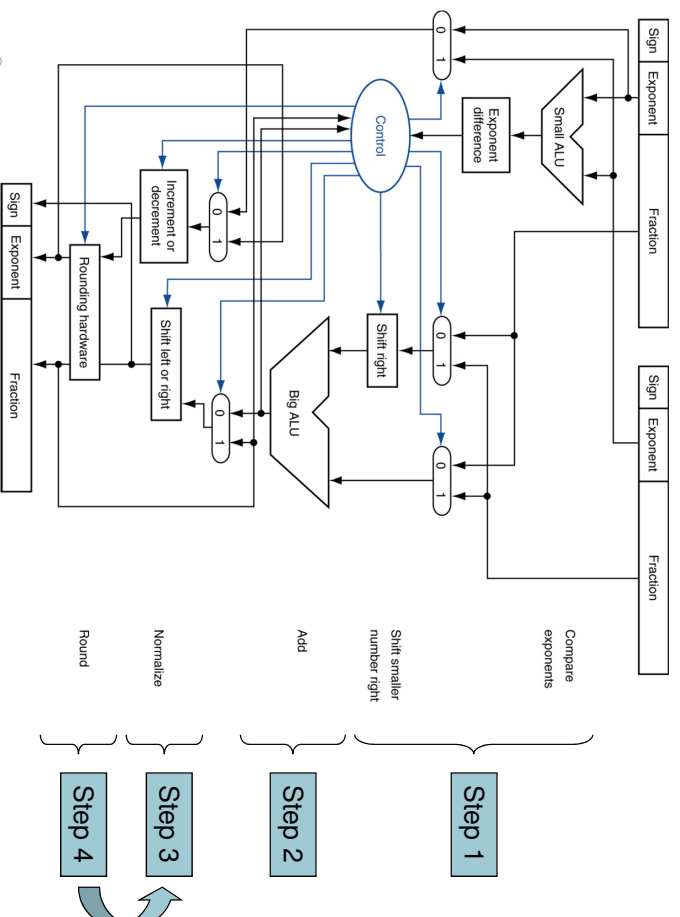


# FP Adder Hardware

- 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



# FP Adder Hardware



# 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



# FP Instructions in MIPS

- FP hardware is coprocessor 1
  - Adjunct processor that extends the ISA
- Separate FP registers
  - 32 single-precision: \$f0, \$f1, ..., \$f31
  - Paired for double-precision: \$f0/\$f1, \$f2/\$f3, ...
    - Release 2 of MIPS ISA supports 32 × 64-bit FP reg's
- FP instructions operate only on FP registers
  - Programs generally don't do integer ops on FP data, or vice versa
  - More registers with minimal code-size impact
- FP load and store instructions
  - `lwc1, ldc1, swc1, sdc1`
    - e.g., `ldc1 $f8, 32($sp)`



# FP Instructions in MIPS

- Single-precision arithmetic
  - `add.s, sub.s, mul.s, div.s`
    - e.g., `add.s $f0, $f1, $f6`
- Double-precision arithmetic
  - `add.d, sub.d, mul.d, div.d`
    - e.g., `mul.d $f4, $f4, $f6`
- Single- and double-precision comparison
  - `c.XX.s, c.XX.d` (`XX` is `eq, lt, le, ...`)
  - Sets or clears FP condition-code bit
    - e.g. `c.lt.s $f3, $f4`
- Branch on FP condition code true or false
  - `bc1t, bc1f`
    - e.g., `bc1t TargetLabel`



# FP Example: °F to °C

- C code:

```
f1oat f2c (f1oat fahr) {  
    return ((5.0/9.0)*(fahr - 32.0));  
}
```
- fahr in \$f12, result in \$f0, literals in global memory space
- Compiled MIPS code:

```
f2c: lwc1  $f16, const5($gp)  
     lwc2  $f18, const9($gp)  
     div.s $f16, $f16, $f18  
     lwc1  $f18, const32($gp)  
     sub.s $f18, $f18, $f18  
     mul.s $f0, $f16, $f18  
     jr   $ra
```



## 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
	1.00E+00	0.00E+00

§3.6 Parallelism and Computer Arithmetic: Associativity

- Need to validate parallel programs under varying degrees of parallelism



# x86 FP Architecture

- Originally based on 8087 FP coprocessor
  - 8 × 80-bit extended-precision registers
  - Used as a push-down stack
  - Registers indexed from TOS: ST(0), ST(1), ...
- FP values are 32-bit or 64 in memory
  - Converted on load/store of memory operand
  - Integer operands can also be converted on load/store
- Very difficult to generate and optimize code
  - Result: poor FP performance

§3.7 Real Stuff: Floating Point in the x86



# x86 FP Instructions

Data transfer	Arithmetic	Compare	Transcendental
FIELD mem/ST(i) FISTP mem/ST(i) FLDPI FLDI FLDZ	FIADD mem/ST(i) FISUBRP mem/ST(i) FIMULP mem/ST(i) FIDIVRP mem/ST(i) FSQRT FABS FRNDINT	FIUCOMP FSTSW AX/mem	FPATAN F2XMI FCOS FPATAN FPREM FPSIN FYL2X

- Optional variations
  - I: integer operand
  - P: pop operand from stack
  - R: reverse operand order
  - But not all combinations allowed



## Concluding Remarks

- ISAs support arithmetic
  - Signed and unsigned integers
  - Floating-point approximation to reals
- Bounded range and precision
  - Operations can overflow and underflow
- MIPS ISA
  - Core instructions: 54 most frequently used
    - 100% of SPECINT, 97% of SPECFP
  - Other instructions: less frequent

