

Integer Addition

- Example: 7 + 6

| | | | | | | |
|-----|-----|-----|-----|-----|-----|-----------|
| ... | (0) | (0) | (1) | (1) | (0) | (Carries) |
| ... | 0 | 0 | 0 | 1 | 1 | 1 |
| ... | 0 | 0 | 0 | 1 | 1 | 0 |
| ... | (0) | (0) | (0) | (1) | (1) | (0) |
| ... | (0) | (0) | (0) | (1) | (1) | (0) |

- 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

Chapter 3 — Arithmetic for Computers — 3

Integer Subtraction

- Add negation of second operand $x - (y) = x + (-y)$
- 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

Chapter 3 — Arithmetic for Computers — 4

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

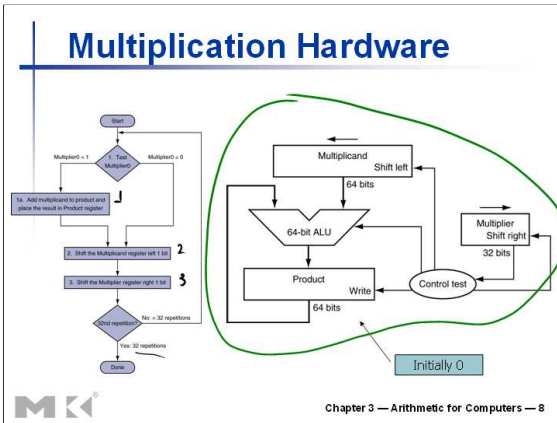
Chapter 3 — Arithmetic for Computers — 5

Multiplication

- Start with long-multiplication approach

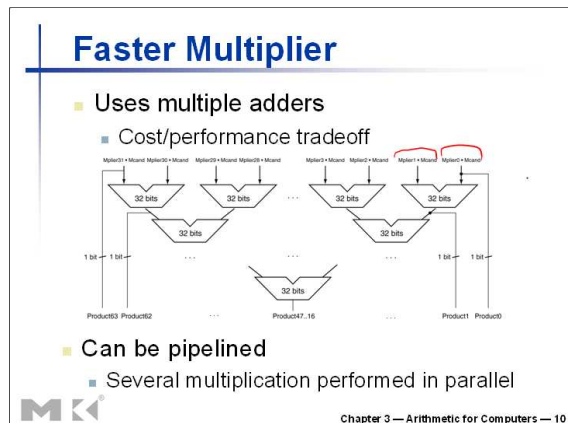
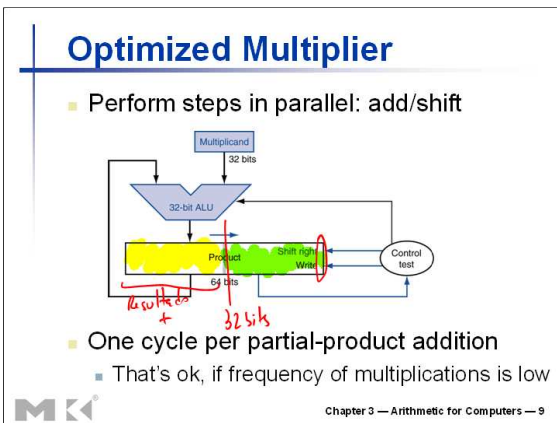
| | | |
|--------------|---|---------|
| multiplicand | → | 1000 |
| multiplier | x | 1001 |
| product | → | 1001000 |

Chapter 3 — Arithmetic for Computers — 7

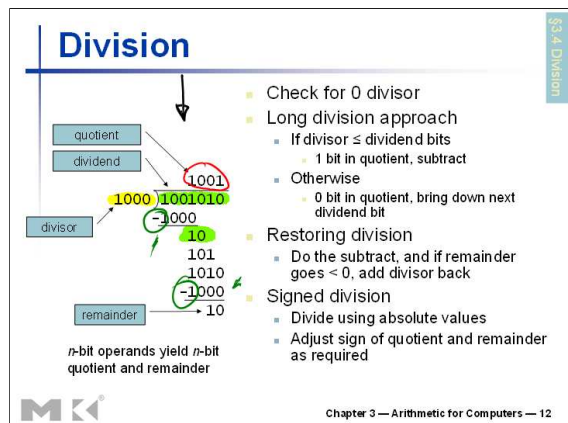


| Step | Pass | Multiplier | Multiplicand | Product |
|------|---------|------------|--------------|-----------|
| 0 | Initial | 0011 | 0000 0010 | 0000 0000 |
| 1 | 1 | 0011 | 0000 0010 | 0000 0010 |
| 2,3 | 2,3 | 0001 | 0000 0100 | 0000 0100 |
| 2 | 1 | 0001 | 0000 0100 | 0000 0110 |
| 2,3 | 2,3 | 0000 | 0000 1000 | 0000 0110 |
| 3 | 2,3 | 0000 | 0001 0000 | 0000 0110 |
| 4 | 2,3 | 0000 | 0010 0000 | 0000 0110 |

4 6



- ## 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
 - `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
- Chapter 3 — Arithmetic for Computers — 11



Division Hardware

Chapter 3 — Arithmetic for Computers — 13

| Iter | Passo | Quotient | Divisor | Remainder |
|------|----------|----------|-----------|-----------|
| 0 | Início | 0000 | 0010 0000 | 0000 0111 |
| 1 | 1: R=R-D | 0000 | 0010 0000 | 0110 0111 |
| | 2: R<0 | 0000 | 0010 0000 | 0000 0111 |
| | 3: D>> | 0000 | 0001 0000 | 0000 0111 |
| 2 | 1: R=R-D | 0000 | 0001 0000 | 0111 0111 |
| | 2: R<0 | 0000 | 0001 0000 | 0000 0111 |
| | 3: D>> | 0000 | 0000 1000 | 0000 0111 |
| 3 | 1: R=R-D | 0000 | 0000 1000 | 1111 1111 |
| | 2: R<0 | 0000 | 0000 1000 | 0000 0111 |
| | 3: D>> | 0000 | 0000 0100 | 0000 0111 |
| 4 | 1: R=R-D | 0000 | 0000 0100 | 0000 0011 |
| | 2: R>0 | 0001 | 0000 0100 | 0000 0011 |
| | 3: D>> | 0001 | 0000 0010 | 0000 0011 |
| 5 | 1: R=R-D | 0001 | 0000 0010 | 0000 0001 |
| | 2: R>0 | 0011 | 0000 0010 | 0000 0001 |
| | 3: D>> | 0011 | 0000 0001 | 0000 0001 |

Optimized Divider

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

Chapter 3 — Arithmetic for Computers — 14

| Iter | Passo | Divisor | Remainder |
|------|--------|---------|-----------|
| 0 | Início | 0010 | 0000 0111 |
| 1 | R-D | 0010 | 0110 1110 |
| | R<0 | 0010 | 0001 1100 |
| 2 | R-D | 0010 | 0111 1100 |
| | R<0 | 0010 | 0011 1000 |
| 3 | R-D | 0010 | 0001 1000 |
| | R>> | 0010 | 0011 0001 |
| 4 | R-D | 0010 | 0001 0001 |
| | R>> | 0010 | 0010 0011 |

Exemplo usando HW otimizado

0001

0011

R Q

Faster Division

- Can't use parallel hardware as in multiplier
 - Subtraction is conditional on sign of remainder
- Faster dividers (e.g. SRT division) generate multiple quotient bits per step
 - Still require multiple steps

Chapter 3 — Arithmetic for Computers — 15

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

Chapter 3 — Arithmetic for Computers — 16

Floating Point

- Representation for non-integral numbers
 - Including very small and very large numbers
- Like scientific notation
 - -2.34×10^{56} ← normalized
 - $+0.002 \times 10^{-4}$ ← not normalized
 - $+987.02 \times 10^9$ ← not normalized
- In binary
 - $\pm 1.xxxxxx_2 \times 2^{yyyy}$
- Types float and double in C

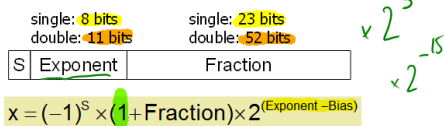


Floating Point Standard

- Defined by IEEE Std 754-1985
- Developed in response to divergence of representations
 - Portability issues for scientific code 5.0
- Now almost universally adopted 3
- Two representations
 - Single precision (32-bit)
 - Double precision (64-bit)



IEEE Floating-Point Format



- S: sign bit (0 ⇒ non-negative, 1 ⇒ 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 1023
 - 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 ≈ 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: 0000000001 ⇒ 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: 1111111110 ⇒ actual exponent = $2046 - 1023 = +1023$
 - Fraction: 111...11 ⇒ significand ≈ 2.0
 - $\pm 2.0 \times 2^{+1023} \approx \pm 4.8 \times 10^{+308}$



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 = $-1 + \text{Bias}$
 - Single: $-1 + 127 = 126 = 01111110_2$
 - Double: $-1 + 1023 = 1022 = 01111111110_2$
- Single: $1011111101000...00$
- Double: $1011111111101000...00$



Floating-Point Example

- What number is represented by the single-precision float $11000000101000...00$
 - $S = 1$
 - 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 = $000...0$
 - \pm Infinity
 - Can be used in subsequent calculations, avoiding need for overflow check
- Exponent = $111...1$, Fraction $\neq 000...0$
 - Not-a-Number (NaN)
 - Indicates illegal or undefined result
 - e.g., $0.0 / 0.0$
 - Can be used in subsequent calculations



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
 - 1.0015×10^2
- 4. Round and renormalize if necessary
 - 1.002×10^2



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
 - $1.000_2 \times 2^{-2}$, with no over/underflow
- 4. Round and renormalize if necessary
 - $1.000_2 \times 2^{-4}$ (no change) = 0.0625

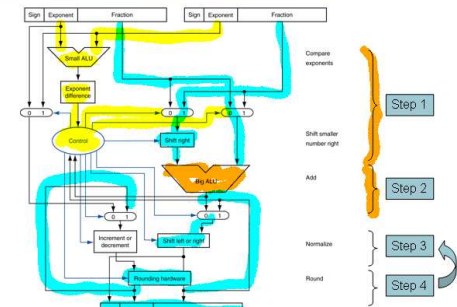


FP Adder Hardware

- 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



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
 - $1.110 \times 9.200 = 10.212 \Rightarrow 10.212 \times 10^5$
- 3. Normalize result & check for over/underflow
 - 1.0212×10^6
- 4. Round and renormalize if necessary
 - 1.021×10^6
- 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×-0.4375)
- 1. Add exponents
 - Unbiased: $-1 + -2 = -3$
 - Biased: $(-1 + 127) + (-2 + 127) = -3 + 254 - 127 = -3 + 127$
- 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



FP Instructions in MIPS

- FP hardware is coprocessor 1
 - Adjunct processor that extends the ISA
- Separate FP registers
 - 32 single-precision: $\$f0, \$f1, \dots, \$f31$
 - Paired for double-precision: $\$f0/\$f1, \$f2/\$f3, \dots$
 - 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., `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:

```
float f2c (float 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, $f12, $f18
     mul.s $f0, $f16, $f18
     jr $ra
```



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

