# 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

**Chapter 3 — Arithmetic for Computers — 41**

# 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

**Chapter 3 — Arithmetic for Computers — 42**

# 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

**Chapter 3 — Arithmetic for Computers — 43**

# 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

**Chapter 3 — Arithmetic for Computers — 44**

# x86 FP Instructions

| Data transfer | Arithmetic | Compare | Transcendental |
|---|---|---|---|
| FILD  mem/ST(i) | FIADDP  mem/ST(i) | FICOMP | FPATAN |
| FISTP mem/ST(i) | FISUBRP mem/ST(i) | FIUCOMP | F2XMI |
| FLDPI | FIMULP  mem/ST(i) | FSTSW AX/mem | FCOS |
| FLD1 | FIDIVRP mem/ST(i) | | FPTAN |
| FLDZ | FSQRT | | FPREM |
| | FABS | | FPSIN |
| | FRNDINT | | FYL2X |

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

**Chapter 3 — Arithmetic for Computers — 45**

# Streaming SIMD Extension 2 (SSE2)

- Adds 4 × 128-bit registers
  - Extended to 8 registers in AMD64/EM64T
- Can be used for multiple FP operands
  - 2 × 64-bit double precision
  - 4 × 32-bit double precision
  - Instructions operate on them simultaneously
    - Single-Instruction Multiple-Data

**Chapter 3 — Arithmetic for Computers — 46**

# Right Shift and Division

§3.8 Fallacies and Pitfalls

- 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. $11111011_2 >>> 2 = 00111110_2 = +62$

**Chapter 3 — Arithmetic for Computers — 47**

# Who Cares About FP Accuracy?

- Important for scientific code
  - But for everyday consumer use?
    - "My bank balance is out by 0.0002¢!" ☹
- The Intel Pentium FDIV bug
  - The market expects accuracy
  - See Colwell, *The Pentium Chronicles*

**Chapter 3 — Arithmetic for Computers — 48**

# Concluding Remarks

§3.9 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

**Chapter 3 — Arithmetic for Computers — 49**