

CPE 431/531

Chapter 3 – Arithmetic for Computers

Dr. Rhonda Kay Gaede



3.1 Introduction

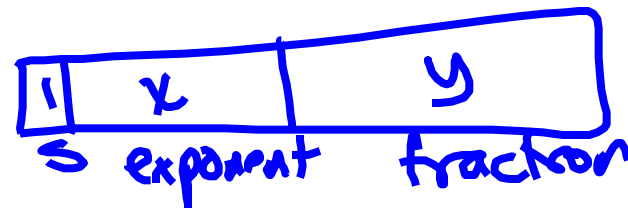
Bits are bits, what is important is how they are interpreted.

- You may have an instruction.
- You may have a signed number (integer).
- You may have an unsigned number (integer).
- You may have a floating-point numbers.

3.5 Floating Point - Basics

- Floating-point numbers are represented in scientific notation
 1.38×10^E
- Floating-point numbers use normalized representation.
- In general, floating-point numbers are of the form
 $(-1)^S \times 1.F \times 2^E$
 F fraction
 E exponent
- There is a tradeoff between range and precision
 - More y bits gives you more precision
 - More x bits gives you more range
- IEEE defines two types of floating-point numbers
 - Single Precision
 - Double Precision

Quad Precision
 Half



3.5 Floating Point – More of the Story

IEEE 754 Floating Point Standard

Adding a bias to the exponent simplifies sorting

The leading one is implicit

Representation expanded

Example: Represent -0.75 in single and double precision

16 bit
S = 1

0.75₁₀ = 0.112
①1 × 2⁻¹

8 exponent -1 + 128 = 0111 1111
bias

23 fraction 10 0

1 0111 1111 10 0
0xBFC0 0000

8 bits of exponent

-128 to +127

127 0111 1111
1 only choice -128 1000 0000
-1 1111 1111
+128 0000 0000
+127 1111 1111

3.5 Floating Point – More Examples

Example: What decimal number is represented by this single precision float? 0x4493 AB00



$S=0$ number is positive

$EXP = 137$

$- BIAS \quad \underline{-128}$
9

$$FRACTION = 2^{-3} + 2^{-6} + 2^{-7} + 2^{-8} + 2^{-10} + 2^{-12} + 2^{-14} + 2^{-15}$$

$$Number = (-1)^S (1.FRACTION) \times 2^{EXP}$$

590.

single
8
double
11
 2^{n-1}

3.6 Subword Parallelism

Graphics systems originally used 8 bits to represent color and 8 bits to represent position.

Support for sound led to 16 bits of information.

Subword items have been supported for a long time in data transfer.

Graphics processing called for arithmetic on subword items.

Often the same operation is performed on groups (vectors) of data.

128 bit adders can handle (Data Level Parallelism)

<u>1</u>	<u>128</u> bit operands
<u>2</u>	<u>64</u> bit operands
<u>4</u>	<u>32</u> bit operands
<u>8</u>	<u>16</u> bit operands
<u>16</u>	<u>8</u> bit operands

2008 Standard FP
adds
half precision 16 bits
quad precision 128 bits

3.9 Fallacies and Pitfalls

$$x + (y + z) = (x + y) + z$$

Pitfall: Floating-point addition is not associative.

Because floating-point numbers are approximations of real numbers and because computer arithmetic has limited precision, associativity does not hold for floating-point numbers.

$$x = -1.5_{10} \times 10^{38}, y = 1.5_{10} \times 10^{38}, z = 1.0$$

$$(x + y) + z = 0 + 1 = 1$$

$$x + (y + z) = 0 + 0 = 0$$

$$x + y + z = 1.0$$

Fallacy: Parallel execution strategies that work for integer data types also work for floating point data types.

Results may be credible but not identical.