CS 261 Fall 2017

Mike Lam, Professor



Binary Arithmetic

Binary Arithmetic

- Topics
 - Basic addition
 - Overflow
 - Multiplication & Division

Basic addition

- Binary and hex addition are fundamentally the same as decimal addition
 - Add digit-by-digit, using a carry as necessary
 - Result generally requires more bits than the two operands



Figure 2.21 Integer addition. With a 4-bit word size, the sum could require 5 bits.

Basic addition

- Binary and hex addition are fundamentally the same as decimal addition
 - Add digit-by-digit, using a carry as necessary
 - Result generally requires more bits than the two operands



Figure 2.21 Integer addition. With a 4-bit word size, the sum could require 5 bits.

Overflow

- Unsigned addition is subject to overflow
 - Caused by truncation to integer size



(assume a 16-bit integer)



Figure 2.23 Unsigned addition. With a 4-bit word size, addition is performed modulo 16.

Overflow

- Two's complement addition is identical to unsigned mechanically
 - Subject to both positive and negative overflow
 - Overflows if carry-in and carry-out differ for sign bit



Multiplication & Division

- Like addition, fundamentally the same as base 10
 - Actually, it's even simpler!
 - Same regardless of encoding
- Special case: multiply by powers of 2 (shift left)

- Division is expensive!
 - Special case: divide by powers of two (shift right)

101 (5) x<u>11</u> (3) 101 <u>101</u> **1111** (15)

Binary fractions

- Now we can store integers
 - But what about general real numbers?
- Extend positional binary integers to store fractions
 - Designate a certain number of bits for the fractional part
 - These bits represent negative powers of two
 - (Just like fractional digits in decimal fractions!)



4 + 1 + 0.5 + 0.125 = **5.625**

Another problem

- For scientific applications, we want to be able to store a wide *range* of values
 - From the scale of galaxies down to the scale of atoms
- Doing this with fixed-precision numbers is difficult
 - Even signed 64-bit integers
 - Perhaps allocate half for whole number, half for fraction
 - Range: ~2 x 10⁻⁹ through ~2 x 10⁹

Floating-point numbers

- Scientific notation to the rescue!
 - Traditionally, we write large (or small) numbers as $x \cdot 10^{e}$
 - This is how floating-point representations work
 - Store exponent and fractional parts (the significand) separately
 - The decimal point "floats" on the number line
 - Position of point is based on the exponent
 - Many nuances and caveats!