

CPS104
Computer Organization and Programming
Lecture 3 : Data representations

Robert Wagner

<http://www.cs.duke.edu/~raw/cps104/Lectures.html>

CPS104 DR.1 ©RW Fall 2000

Overview of Today's Lecture:

- Introduction: Decimal, Binary, Octal and Hexadecimal numbers.
- Storage types: Byte, Word, Double-word.
- ASCII Characters.
- 2's complement numbers.
- Floating-point numbers.

Sections 3,7, 4.1-4.3, 4.8 In the textbook

CPS104 DR.2 ©RW Fall 2000

Binary , Octal and Hexidecimal numbers

- Decimal numbers are used by people! Most computers use Binary numbers and Binary Arithmetic!!
- Computers can input and output decimal numbers but they convert them to internal binary number representation.
- Binary numbers use only two different digits: {0,1}
 - * Example: $1200_{10} = 0000\ 0100\ 1011\ 0000_2$
- Octal numbers use 8 digits: {0 - 7}
 - * Example: $1200_{10} = 04260_8$
- Hexadecimal numbers use 16 digits: {0-9, A-F}
 - * Example: $1200_{10} = 0X04B0_{16}$

CPS104 DR.3 ©RW Fall 2000

Binary representations for Integers Sign-Magnitude

- Advantages:
 - * Simple extension of unsigned numbers.
 - * Same number of positive and negative numbers.
- Disadvantages:
 - * Two representations for 0: $0=000000$; $-0=100000$.
 - * Algorithm (circuit) for addition depends on the arguments' signs.

CPS104 DR.13

©RW Fall 2000

Binary representations for Integers 2's Complement

- Advantages:
 - * Only one representation for 0: $0 = 000000$
 - * Addition algorithm independent of sign bits.
- Disadvantages:
 - * One more negative number than positive :
Example: $100000_2 = -32_{10}$; but 32_{10} could not be represented as a 6-bit 2's complement number.

CPS104 DR.14

©RW Fall 2000

Binary Arithmetic

- Addition
 - * As in DECIMAL arithmetic, add low order column of digits up, forming a digits sum DS
 - * Record the sum-digit for this column as $DS \bmod 2$
 - * Include a CARRY of $\text{floor}(DS/2)$ in the next column's sum
 - * Example: $18_{10} + 14_{10} = 32_{10}$

| | | |
|---------------|------------|------------------------------|
| | 010010_2 | |
| → $18_{10} =$ | 010010_2 | $+001110_2$ |
| → $14_{10} =$ | 001110_2 | <u>100000_2</u> |

CPS104 DR.15

©RW Fall 2000

Binary Arithmetic 2

- Shifts
 - * A LEFT shift of K bit positions MULTIPLIES the number by 2^K
 - $18_{10} \ll 2$ is $010010_2 \ll 2 = 01001000_2 = 64_{10} + 8_{10} = 72 = 18 \cdot 4$
 - * A RIGHT shift of K bit positions DIVIDES the number by 2^K dropping the remainder (low-order digits)
 - $18_{10} \gg 2$ is $010010_2 \gg 2 = 0100_2 = 4$
 - * The REMAINDER left when a number is divided by 2^K is the LOW-ORDER K BITS
 - $18_{10} \% 4$ is $010010_2 \% 2^2 = 0100_2 = 4$

CPS104 DR.16

©RW Fall 2000

2's Complement (cont.) Complementing (negating)

- To negate a number do:
 - * A: complement each digit (Change 0 to 1, and 1 to 0)
 - * B: add 1.
 - * Example: $14_{10} = 001110_2$; $-14_{10} = 110001_2 + 1 = 110010_2$
- To add signed numbers use regular addition but disregard carry out.
 - * Example: $18_{10} - 14_{10} = 18_{10} + (-14_{10}) = 4_{10}$

$$\begin{array}{r} 010010_2 \\ +110010_2 \\ \hline 000100_2 \end{array}$$

CPS104 DR.17

©RW Fall 2000

Why 2's Complement Addition Works (1)

- Examples for 5-bit 2's complement integers

$$\begin{array}{l} 4_{10} + (-5_{10}) = -1 \\ 00100 + (-00101) \\ 00100 + 11011 = 11111 = -(00001) = -1 \end{array}$$

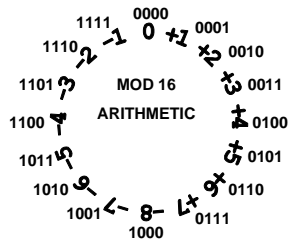
$$\begin{array}{l} (-4_{10}) + (-5_{10}) = -9_{10} \\ (-00100) + (-00101) \\ 11100 + 11011 = 11011 = -(01001) = -9_{10} \end{array}$$

CPS104 DR.18

©RW Fall 2000

How 2's Complement Works (4 bit)

○



CPS104 DR.19

©RW Fall 2000

Why 2's Complement Addition Works (2)

○ N-bit 2's complement arithmetic is exact arithmetic mod 2^N

- * Suppose J, K are small positive integers, $J > K$.
- * $R(-J) = \sim J + 1 = 111\dots111$ (N 1's) $\sim J + 1 = 2^N - 1 - J + 1$.
- $R(-J) = 2^N - J$, $R(-K) = 2^N - K$, $R(J) = J$, $R(K) = K$

$$R(-J) + R(-K) = 2^N + (2^N - J - K) = (2^N - J - K) \text{ mod } 2^N$$

$$R(J) + R(-K) = 2^N + (J - K) = (J - K) \text{ mod } 2^N$$

(Because $J - K$ is greater than 0)

$$R(-J) + R(K) = (2^N - J + K) = (2^N - (J - K)) \text{ mod } 2^N$$

(Because $-J + K$ is less than 0)

$$R(J) + R(K) = (J + K) = (J + K) \text{ mod } 2^N$$

CPS104 DR.20

©RW Fall 2000

2's Complement (cont.)

Example: A = 0x0ABC; B = 0x0FEB.

Compute: A + B and A - B in 16-bit 2s complement arithmetic.

CPS104 DR.21

©RW Fall 2000

**2's Complement (cont.)
Precision Extension**

- To extend precision do: sign bit extension:

Example:

* $14_{10} = 001110_2$ in 6-bit representation.
 $14_{10} = 00000001110_2$ in 12-bit representation

* $-14_{10} = 110010_2$ in 6-bit representation
 $-14_{10} = 111111110010_2$ in 12-bit representation.

CPS104 DR.22

©RW Fall 2000

Changing base representation

- Suppose you can do arithmetic in base Z, and want to convert a number N from base Z to base B

- Write $N = q \times B + r$, where

$$q = \sum_{i=1}^k N_i \times B^{i-1}, \quad r = N_0$$

- * R is the low-order digit of N in base B, and q can be converted to get the remaining digits

- * Like shifting the base-point:

$$N = N_k N_{k-1} \dots N_1 . N_0 \times B$$

- * For fractions, MULTIPLY by B; the integer part is the high-order digit, and the fraction holds the other digits:

$$N \times B = N_k . N_{k-1} \dots N_1 N_0$$

- To convert TO base Z, just evaluate the sum:

$$N = \sum_{i=0}^k N_i \times B^i$$

CPS104 DR.23

©RW Fall 2000

Limits on 2's Complement Arithmetic

- N-bit 2's complement represents a finite RANGE of #s

- * Numbers outside this range cannot be represented

- * Arithmetic operation can produce unrepresentable #s

- * Only values of I satisfying eqn below can be represented:

$$-2^{N-1} \leq I \leq 2^{N-1} - 1$$

- When operation gives an out-of-range answer A, machine gives A', where $R(A') = R(A) \text{ mod } 2^N$

- Examples (6-bit 2's complement)

$$\begin{array}{r} 100010 = -30 \\ + 100001 = -31 \\ \hline 000011 = 3 \end{array} \qquad \begin{array}{r} 001101 = 13 \\ + 011000 = 24 \\ \hline 100101 = -27 \end{array}$$

CPS104 DR.24

©RW Fall 2000

2's Complement with Implicit Binary Point

- 2's complement numbers need not represent integers -- each can have a "hidden" binary point, and thus represent a number with an integral part, AND a fractional part:

$$\begin{array}{r} 010010_2 = 18_{10} \\ +001011 = 11 \\ \hline 011101 = 29 \end{array} \qquad \begin{array}{r} 010.010_2 = 2.25_{10} \\ +001.011 = 1.375 \\ \hline 011.101 = 3.625 \end{array}$$

- Addition (IF the binary points are aligned) produces the same sequence of bits, regardless of the position of the "."
- This property underlies floating point addition
 - * Uses a "significand", F, with a hidden, implicit binary point, multiplied by 2^E , where the number representation includes fields that give values for "E" and "F".

CPS104 DR.25

©RW Fall 2000

Floating Point Representation

Numbers are represented by:

$$X = (-1)^s \times 2^{E-127} \times 1.M$$

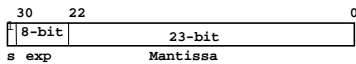
S := 1-bit field ; Sign bit

E := 8-bit field; Exponent: Biased integer, $0 \leq E \leq 255$.

M := 23-bit field; Mantissa: Normalized fraction with hidden 1.

F := 1.M = Significand

Single precision floating point number



CPS104 DR.26

©RW Fall 2000

Floating Point Representation (cont.)

- The mantissa represents a fraction using binary notation:

$$M = .s_1, s_2, s_3 \dots = 1.0 + s_1 \cdot 2^{-1} + s_2 \cdot 2^{-2} + s_3 \cdot 2^{-3} + \dots$$

- Example: $X = -0.75_{10}$ in single precision:

$$-0.75_{10} = -0.11_2 = (-1) \times 1.1_2 \times 2^{-1} = (-1) \times 1.1_2 \times 2^{126-127}$$

$$S = 1 ; \text{Exp} = 126_{10} = 0111\ 1110_2 ;$$

$$M = 100\ 0000\ 0000\ 0000\ 0000_2$$

$$X = \boxed{1\ 0111\ 1110\ 100\ 0000\ 0000\ 0000\ 0000\ 0000}_2$$

CPS104 DR.27

©RW Fall 2000

Floating Point Representation
(cont.)

Example:
What floating-point number is : 0xC1580000 ?

CPS104 DR.28

©RW Fall 2000

Floating Point Representation
(cont.)

Example:
What floating-point number is : 0xC1580000 ?
1100 0001 0101 1000 0000 0000 0000 0000

| | | |
|---|------------|------------------------------|
| S | Exp 8-bit | Mantissa 23-bit |
| 1 | 100 0001 0 | 101 1000 0000 0000 0000 0000 |

- $E=128+2-127=3$ $-1.1011_2 \times 2^3 = -1101.1_2 = -13.5$

CPS104 DR.29

©RW Fall 2000

Floating Point Arithmetic

- To add FP #s, Equalize Exponents, then add Adjusted Significands.
 - * Extend M's, by including the hidden leading "1": gives F's.
 - * Shift F_A of #, A, with SMALLER exponent right
 - Each 1-bit right-shift adds 1 to E_A
 - Shift until $E_A = E_B$
 - * Add shifted F_A to F_B to get F_S observing sign rules
 - * Re-normalize (shift F_S left, reducing E_s , till leading 1 disappears)
 - * Significant bits may be lost on the RIGHT. Approximate answer. "Weight" of lost bits is (hopefully) small.
- Numbers out of range produce "special" representations, using E field of all 1's, or all 0's.
 - * NaN's, and 0 See TABLE, end of Sect. 4.8, p 301

CPS104 DR.30

©RW Fall 2000

FP addition example

A = 0xC1D80000, B=0xC2E80000
 0xC1D80000 = 1 | 100 0001 1 | 101 1000 0000 ...
 0xC2E80000 = 1 | 100 0010 1 | 110 1000 0000 ...
 $F_A = 1.1011000$
 Adjust F_A : 5-3=2 shifts: new $F_A = 0.011011000$

$F_B = 1.1101000, E_B=100\ 0010\ 1$
 shifted $F_A = 0.0110110$
 $F_S = 10.00111110, E_S=100\ 0010\ 1$
 Re-normalize $F_S = 1.00011111, E_S=100\ 0011\ 0$
 Result = 1 | 100 0011 0 | 000 1111 1 = 0xC30F8000

CPS104 DR.31

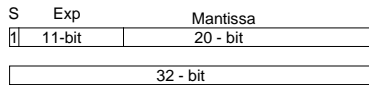
©RW Fall 2000

Floating Point Representation
(cont.)

- o **Double Precision Floating point:**

64-bit representation: 1-bit sign, 11-bit (biased) exponent;
 52-bit mantissa (with hidden 1).

Double precision floating point number



$$X = (-1)^S \times 1.M \times 2^{(E-1023)}$$

CPS104 DR.32

©RW Fall 2000
