

University of Southern California

Viterbi School of Engineering

EE352

Computer Organization and Architecture

Fixed Point Systems

References:

- 1) Textbook
- 2) Mark Redekopp's slide series

Shahin Nazarian

Spring 2010

Data Representation

- In C/C++ variables can be of different types and sizes
 - Integer Types (signed and unsigned)

C Type	Bytes	Bits	MIPS Name
[unsigned] char	1	8	byte
[unsigned] short [int]	2	16	half-word
[unsigned] long [int]	4	32	word
[unsigned] long long [int]	8	64	double-word

- Floating Point Types

C Type	Bytes	Bits	MIPS Name
float	4	32	single
double	8	64	double

Binary Representation Systems

- **Integer Systems**

- **Unsigned**

- Unsigned (Normal) binary

- **Signed**

- Signed Magnitude
 - 2's complement
 - *Excess-N**
 - *1's complement**

- **Floating Point**

- For very large and small (fractional) numbers

- **Codes**

- **Text**

- ASCII / Unicode

- **Decimal Codes**

- BCD (Binary Coded Decimal) / (8421 Code)

* = Not fully covered in this class

Shahin Nazarian/EE352/Spring10

Number Systems

- Number systems consist of
 1. A base (radix) r
 2. r coefficients $[0 \text{ to } r-1]$
- Human System: Decimal (Base 10): 0,1,2,3,4,5,6,7,8,9
- Computer System: **Binary** (Base 2): 0,1
- Human systems for working with computer systems (shorthand for human to read/write binary)
 - **Octal** (Base 8): 0,1,2,3,4,5,6,7
 - **Hexadecimal** (Base 16): 0-9,A,B,C,D,E,F (A thru F = 10 thru 15)

Anatomy of a Decimal Number

- A number consists of a string of explicit coefficients (digits)
- Each coefficient has an implicit place value which is a power of the base
- The value of a decimal number (a string of decimal coefficients) is the sum of each coefficient times its place value

Radix (base)

$$(934)_{10} = 9 * 10^2 + 3 * 10^1 + 4 * 10^0 = 934$$

Explicit coefficients

Implicit place values

$$(3.52)_{10} = 3 * 10^0 + 5 * 10^{-1} + 2 * 10^{-2} = 3.52$$

Positional Number Systems (Unsigned)

- A number in base r has place values/weights that are the powers of the base
- Denote the coefficients as: a_i

Left-most digit = Most
Significant Digit
(MSD)

Right-most digit =
Least Significant Digit
(LSD)

$$\dots \quad \frac{a_3}{r^3} \quad \frac{a_2}{r^2} \quad \frac{a_1}{r^1} \quad \frac{a_0}{r^0} \quad . \quad \frac{a_{-1}}{r^{-1}} \quad \frac{a_{-2}}{r^{-2}} \quad \dots$$

$$N_r = \sum_i (a_i * r^i) = D_{10}$$

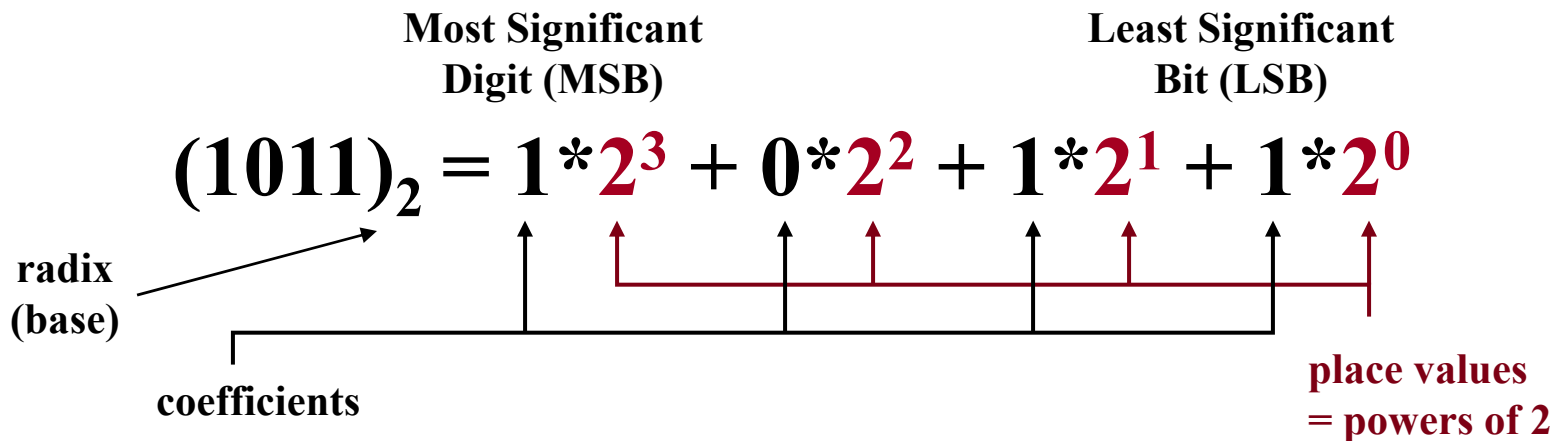
Examples

$$\begin{aligned}(746)_8 &= 7*8^2 + 4*8^1 + 6*8^0 \\ &= 448 + 32 + 16 = 486_{10}\end{aligned}$$

$$\begin{aligned}(1A5)_{16} &= 1*16^2 + 10*16^1 + 5*16^0 \\ &= 256 + 160 + 5 = 421_{10}\end{aligned}$$

Anatomy of a Binary Number

- Same as decimal but now the coefficients are 1 and 0 and the place values are the powers of 2



Binary Examples

$$\begin{array}{c} \underline{1} \underline{0} \underline{0} \underline{1} . \underline{1} \\ 8 \quad 1.5 \end{array} \bigg)_2 = 8 + 1 + 0.5 = 9.5_{10}$$

$$\begin{array}{c} \underline{1} \underline{0} \underline{1} \underline{1} \underline{0} \underline{0} \underline{0} \underline{1} \\ 128 \ 32 \ 16 \quad 1 \end{array} \bigg)_2 = 128 + 32 + 16 + 1 = 177_{10}$$

Powers of 2

$$2^0 = 1$$

$$2^1 = 2$$

$$2^2 = 4$$

$$2^3 = 8$$

$$2^4 = 16$$

$$2^5 = 32$$

$$2^6 = 64$$

$$2^7 = 128$$

$$2^8 = 256$$

$$2^9 = 512$$

$$2^{10} = 1024$$

<u>1024</u>	<u>512</u>	<u>256</u>	<u>128</u>	<u>64</u>	<u>32</u>	<u>16</u>	<u>8</u>	<u>4</u>	<u>2</u>	<u>1</u>
-------------	------------	------------	------------	-----------	-----------	-----------	----------	----------	----------	----------

Conversion to Decimal Examples

- Decimal equivalent is...

... the sum of each coefficient multiplied by its implicit place value (power of the base)

$$= \sum_i (a_i * r^i) \quad [a_i = \text{coefficient}, r = \text{base}]$$

$$\begin{aligned}(11010)_2 &= 1*2^4 + 1*2^3 + 1*2^1 \\ &= 16 + 8 + 2 = (26)_{10}\end{aligned}$$

$$\begin{aligned}(6523)_8 &= 6*8^3 + 5*8^2 + 2*8^1 + 3*8^0 \\ &= 3072 + 320 + 16 + 3 = (3411)_{10}\end{aligned}$$

$$\begin{aligned}(AD2)_{16} &= 10*16^2 + 13*16^1 + 2*16^0 \\ &= 2560 + 208 + 2 = (2770)_{10}\end{aligned}$$

Unique Combinations

- Given n digits of base r , how many unique numbers can be formed? r^n
 - What is the range? $[0 \text{ to } r^n - 1]$

2-digit, decimal numbers ($r=10, n=2$)

0-9 0-9

100 combinations:
00-99

3-digit, decimal numbers ($r=10, n=3$)

____ ____ ____

1000 combinations:
000-999

4-bit, binary numbers ($r=2, n=4$)

____ ____ ____ ____
0-1 0-1 0-1 0-1

16 combinations:
0000-1111

6-bit, binary numbers ($r=2, n=6$) ____ ____ ____ ____ ____ ____

64 combinations:
000000-111111

Main Point: Given n digits of base r , r^n unique numbers can be made with the range $[0 - (r^n - 1)]$

Approximating Large Powers of 2

- Often need to find decimal approximation of a large powers of 2 like 2^{16} , 2^{32} , etc.

- Use following approximations:

- $2^{10} \approx 10^3$ (1 thousand) = 1 Kilo-
- $2^{20} \approx 10^6$ (1 million) = 1 Mega-
- $2^{30} \approx 10^9$ (1 billion) = 1 Giga-
- $2^{40} \approx 10^{12}$ (1 trillion) = 1 Tera-

- For other powers of 2, decompose into product of 2^{10} or 2^{20} or 2^{30} and a power of 2 that is less than 2^{10}

- 16-bit half word: 64K numbers
- 32-bit word: 4G numbers
- 64-bit dword: 16 million trillion numbers

$$2^{16} = 2^6 * 2^{10} \\ \approx 64 * 10^3 = 64,000$$

$$2^{24} = 2^4 * 2^{20} \\ \approx 16 * 10^6 = 16,000,000$$

$$2^{28} = 2^8 * 2^{20} \\ \approx 256 * 10^6 = 256,000,000$$

$$2^{32} = 2^2 * 2^{30} \\ \approx 4 * 10^9 = 4,000,000,000$$

Decimal to Unsigned Binary

- To convert a decimal number, x , to binary:
 - Only coefficients of 1 or 0. So simply find place values that add up to the desired values, starting with larger place values and proceeding to smaller values and place a 1 in those place values and 0 in all others

$$25_{10} = \frac{0}{32} \frac{1}{16} \frac{1}{8} \frac{0}{4} \frac{0}{2} \frac{1}{1}$$

For 25_{10} the place value 32 is too large to include so we include 16.
Including 16 means we have to make 9 left over. Include 8 and 1

Decimal to Unsigned Binary

$$73_{10} = \begin{array}{cccccccc} \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{1} \\ \hline 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \end{array}$$

$$87_{10} = \begin{array}{cccccccc} \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{1} & \mathbf{1} & \mathbf{1} \\ \hline & & & & & & & \end{array}$$

$$145_{10} = \begin{array}{cccccccc} \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} \\ \hline & & & & & & & \end{array}$$

$$0.625_{10} = \begin{array}{ccccc} \mathbf{1} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \hline .5 & .25 & .125 & .0625 & .03125 \end{array}$$

Decimal to Another Base

- To convert a decimal number, x , to base r :
 - Use the place values of base r (powers of r)
 - Starting with largest place values, fill in coefficients that sum up to desired decimal value without going over

$$75_{10} = \frac{0}{256} + \frac{4}{16} + \frac{B}{1}$$

Conversion Methods Summary

- Base $r \Rightarrow$ Base 10
 - Sum of coefficients * place values
- Base 10 \Rightarrow Base r
 - Find place values and coefficients that add up to desired value

Binary, Octal, and Hexadecimal

- Octal (base 8 = 2^3)
- 1 Octal digit ($_$)₈ can represent: 0 - 7
- 3 bits of binary ($_ _ _$)₂ can represent:
000-111 = 0 - 7
- Conclusion...
1 Octal digit = 3 bits
- Hex (base 16 = 2^4)
- 1 Hex digit ($_$)₁₆ can represent: 0-F (0-15)
- 4 bits of binary ($_ _ _ _$)₂ can represent:
0000-1111 = 0-15
- Conclusion...
1 Hex digit = 4 bits

Binary to Octal or Hex

- Make groups of 3 bits starting from radix point and working outward
- Add 0's where necessary
- Convert each group of 3 to an octal digit

101001110.110

5 1 6 6

516.6₈

- Make groups of 4 bits starting from radix point and working outward
- Add 0's where necessary
- Convert each group of 4 to an octal digit

000101001110.1100

1 4 E C

14E.C₁₆

Octal or Hex to Binary

- Expand each octal digit to a group of 3 bits

317.2_8

$\overbrace{011}\overbrace{001}\overbrace{111}.0\overbrace{10}_2$

11001111.01_2

- Expand each hex digit to a group of 4 bits

$D93.8_{16}$

$\overbrace{1101}\overbrace{1100}\overbrace{1001}.1\overbrace{000}_2$

110110010011.1_2

Hexadecimal Representation

- Since values in modern computers are many bits, we use hexadecimal as a shorthand notation (4 bits = 1 hex digit)
 - 11010010 = D2 hex
 - 0111011011001011 = 76CB hex
- Important Point: To interpret the value of a hex number, you must know what underlying binary system is assumed (unsigned, 2's comp. etc.)

Conversion Methods

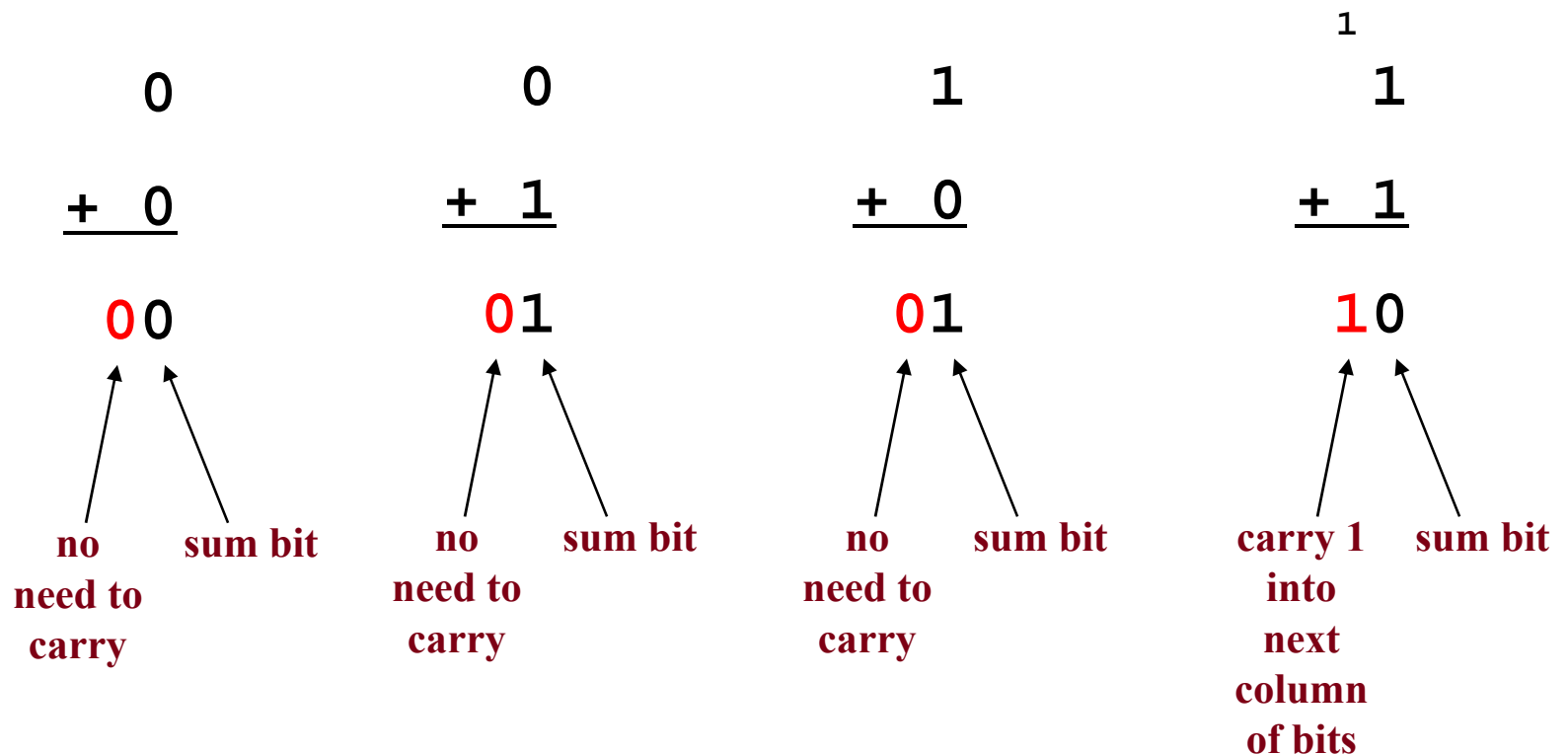
- Base $r \Rightarrow$ Base 10
 - Sum of coefficients * place values
- Base 10 \Rightarrow Base r
 - Find place values and coefficients that add up to desired value
- Binary \Leftrightarrow Octal
 - 3-bits = 1 octal digit
- Binary \Leftrightarrow Hex
 - 4-bits = 1 hex digit

Binary Arithmetic

- Can perform all arithmetic operations (+, -, *, ÷) on binary numbers
- Can use same methods as in decimal
 - Still use carries and borrows, etc.
 - Only now we carry when sum is 2 or more rather than 10 or more (decimal)
 - We borrow 2's not 10's from other columns
- Easiest method is to add bits in your head in decimal ($1+1 = 2$) then convert the answer to binary ($2_{10} = 10_2$)

Binary Addition

- In decimal addition we carry when the sum is 10 or more
- In binary addition we carry when the sum is 2 or more
- Add bits in binary to produce a sum bit and a carry bit



Binary Addition & Subtraction

$$\begin{array}{r} 1 \ 1 \ 1 \\ 0 \ 1 \ 1 \ 1 \ (7) \\ + \ 0 \ 0 \ 1 \ 1 \ (3) \\ \hline 1 \ 0 \ 1 \ 0 \ (10) \\ \text{8} \quad \text{4} \quad \text{2} \quad \text{1} \end{array}$$

$$\begin{array}{r} 0 \quad 0 \\ \cancel{1} \ 10 \ \cancel{1} \ 10 \ (10) \\ - \ 0 \ 1 \ 0 \ 1 \ (5) \\ \hline 0 \ 1 \ 0 \ 1 \ (5) \\ \text{8} \quad \text{4} \quad \text{2} \quad \text{1} \end{array}$$

Binary Addition & Subtraction

$$\begin{array}{r} 1 \ 1 \ 1 \\ 0 \ 1 \ 1 \ 1 \ (7) \\ + \ 0 \ 0 \ 1 \ 1 \ (3) \\ \hline 1 \ 0 \ 1 \ 0 \ (10) \\ \text{8} \quad \text{4} \quad \text{2} \quad \text{1} \end{array}$$

$$\begin{array}{r} 0 \quad 0 \\ \cancel{1} \ 10 \ \cancel{1} \ 10 \ (10) \\ - \ 0 \ 1 \ 0 \ 1 \ (5) \\ \hline 0 \ 1 \ 0 \ 1 \ (5) \\ \text{8} \quad \text{4} \quad \text{2} \quad \text{1} \end{array}$$

Binary Addition

$$\begin{array}{r} 110 \\ 0110 (6) \\ + 0111 (7) \\ \hline 1101 (13) \\ \text{8 4 2 1} \end{array}$$

Binary Addition

$$\begin{array}{r} 0 \\ 0110 (6) \\ + 0111 (7) \\ \hline 1101 (13) \end{array}$$

$$\begin{array}{r} 0 \\ + 1 \\ \hline 01 \end{array}$$

carry bit sum bit

Binary Addition

$$\begin{array}{r} \textcolor{red}{1}0 \\ 0110 \text{ (6)} \\ + 0111 \text{ (7)} \\ \hline 1101 \text{ (13)} \end{array}$$

$$\begin{array}{r} 0 \\ 1 \\ + 1 \\ \hline \textcolor{red}{1}0 \end{array}$$

carry bit sum bit

Binary Addition

$$\begin{array}{r} \textcolor{red}{1}10 \\ 0110 \text{ (6)} \\ + 0111 \text{ (7)} \\ \hline 1101 \text{ (13)} \end{array}$$

$$\begin{array}{r} 1 \\ 1 \\ + 1 \\ \hline \textcolor{red}{1}1 \end{array}$$

carry bit sum bit

$1+1+1 = 3_{10} = 11_2$

Binary Addition

$$\begin{array}{r} 110 \\ 0110 (6) \\ + 0111 (7) \\ \hline 1101 (13) \end{array}$$

$$\begin{array}{r} 1 \\ 0 \\ + 0 \\ \hline 01 \end{array}$$

carry bit sum bit

Binary Subtraction

- If you can't perform subtraction in one column borrow from higher order columns
- When you borrow you are borrowing a 2, not a 10 as in decimal

$$\begin{array}{r} 1010 \text{ (10)} \\ - 0111 \text{ (7)} \\ \hline 0011 \text{ (3)} \end{array}$$

Binary Subtraction

$$\begin{array}{r} 1010 \text{ (10)} \\ - 0111 \text{ (7)} \\ \hline 0011 \text{ (3)} \end{array}$$

Can't perform $0 - 1$, so we must borrow

$$\begin{array}{r} \\ \cancel{} \\ - \\ \hline \end{array}$$

We borrow the 1 (which is worth 2) from the next column and now we can perform $10 - 1 = 9$

And now we go to the next column

Binary Subtraction

$$\begin{array}{r} 10\cancel{1}0 \text{ (10)} \\ - 0111 \text{ (7)} \\ \hline 0011 \text{ (3)} \end{array}$$

The diagram illustrates a binary subtraction problem. The top number is 1010 (decimal 10), and the bottom number is 0111 (decimal 7). A red oval highlights the third column from the right, where the top digit is 1 and the bottom digit is 1. The result of the subtraction is 0011 (decimal 3).

Can't perform 0 – 1, so we must borrow

Binary Subtraction

We get the borrow from this column and work back to the current column

$$\begin{array}{r} \text{0 1 10} \\ \cancel{1}\cancel{1}0 \cancel{1}10 \text{ (10)} \\ - \quad 0 \ 1 \ 1 \ 1 \text{ (7)} \\ \hline 0 \ 0 \ 1 \ 1 \text{ (3)} \end{array}$$

We can't borrow from the column next to us (it is 0 as well) so we must try to borrow from the next column and then work our way back to the current column where we perform $10 - 1 = 1$

Binary Subtraction

$$\begin{array}{r} 0110 \\ \cancel{1}\cancel{1}\cancel{0} \text{ (10)} \\ - 0111 \text{ (7)} \\ \hline 0011 \text{ (3)} \end{array}$$

We can perform $1 - 1 = 0$

Binary Subtraction

$$\begin{array}{r}
 \begin{array}{cccc}
 0 & 1 & 1 & 0 \\
 \cancel{1} & \emptyset & \cancel{1} & 10 \text{ (10)}
 \end{array} \\
 - \begin{array}{cccc}
 0 & 1 & 1 & 1 \text{ (7)}
 \end{array} \\
 \hline
 \begin{array}{cccc}
 0 & 0 & 1 & 1 \text{ (3)}
 \end{array}
 \end{array}$$

We can perform $0 - 0 = 0$

Binary Multiplication

- Like decimal multiplication, find each partial product and shift them, then sum them up
- Multiplying two n -bit numbers yields at most a $2*n$ -bit product

$$\begin{array}{r} 0110(6) \\ * 0101(5) \\ \hline 0110 \\ 0000 \\ 0110 \\ + 0000 \\ \hline 0011110 \end{array} \leftarrow \text{Sum of the partial products}$$

Binary Multiplication

$$\begin{array}{r} (6) \\ * (5) \\ \hline \end{array} \leftarrow \text{First partial product}$$

Binary Multiplication

$$\begin{array}{r} 0110(6) \\ * 0101(5) \\ \hline 0110 \\ 0000 \\ 0110 \\ 0000 \end{array}$$

← Fourth partial product

Binary Multiplication

$$\begin{array}{r} 0110(6) \\ * 0101(5) \\ \hline 0110 \\ 0000 \\ 0110 \\ + 0000 \\ \hline 0011110 \end{array}$$

← Sum up the partial products

Binary Division

- Use the same long division techniques as in decimal

$$\begin{array}{r}
 \begin{array}{c} (2)_1 \\ 0 \end{array} \quad 10 \overline{) \begin{array}{cccc} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ \hline -1 & 0 & & \\ \hline & 0 & 1 & \\ & -0 & 0 & \\ \hline & & 1 & 1 \\ & & -1 & 0 \\ \hline & & & 0 & 1 \end{array} } \begin{array}{l} r.1 \\ (5 \text{ } r.1)_{10} \\ (11)_{10} \end{array}
 \end{array}$$

Binary Division

10 (2) goes into 1, 0 times. Since it doesn't, bring in the next bit

$$\begin{array}{r} 0 \\ 10 \overline{) 1011} \end{array}$$

10 (2) goes into 10, 1 time. Multiply, subtract, and bring down the next bit

[illegible]

Binary Division

10 (2) goes into 01, 0 times. Multiply, subtract, and bring down the next bit

$$\begin{array}{r} 00 \\ 10 \overline{) 1011} \\ \underline{-10} \\ 01 \\ \underline{-00} \\ 11 \end{array}$$

Binary Division

10 (2) goes into 11, 1 time. Multiply and subtract. The remainder is 1

$$\begin{array}{r} 0 0 r.1 \\ 10 \overline{) 1 1 1} \\ \underline{-1 0} \\ 1 \\ \underline{-0 0} \\ 1 \\ \\ \\ \end{array}$$

Hexadecimal Arithmetic

- Same style of operations
 - Carry when sum is 16 or more, etc.

$$\begin{array}{r} 1 \ 1 \\ 4 \ D_{16} \\ + \ B \ 5_{16} \\ \hline 1 \ 0 \ 2_{16} \end{array}$$

$$13+5 = 18_{10} = \underline{1} \underline{2}_{16}$$

16 1

$$1+4+11 = 16_{10} = \underline{1} \underline{0}_{16}$$

16 1

One More Example

$$\begin{array}{r} 0 \\ 0 \quad \cancel{1} \quad 10 \quad 0 \quad (4) \\ - \quad \underline{0 \quad 1 \quad 1 \quad 0} \quad (6) \\ \hline \phantom{\cancel{1}} \quad (-2) \end{array}$$

- What about negative numbers?
 - Current unsigned binary system can only represent positive numbers
 - Given n -bits \Rightarrow Range: $[0 \text{ to } 2^n - 1]$

Signed Magnitude

2's Complement System

SIGNED SYSTEMS

Binary Representation Systems

- Integer Systems
 - Unsigned
 - Unsigned (Normal) binary
 - Signed
 - Signed Magnitude
 - 2's complement
 - *1's complement**
 - *Excess-N**
- Floating Point
 - For very large and small (fractional) numbers
- Codes
 - Text
 - ASCII / Unicode
 - Decimal Codes
 - BCD (Binary Coded Decimal) / (8421 Code)

* = Not covered in this class

Unsigned and Signed

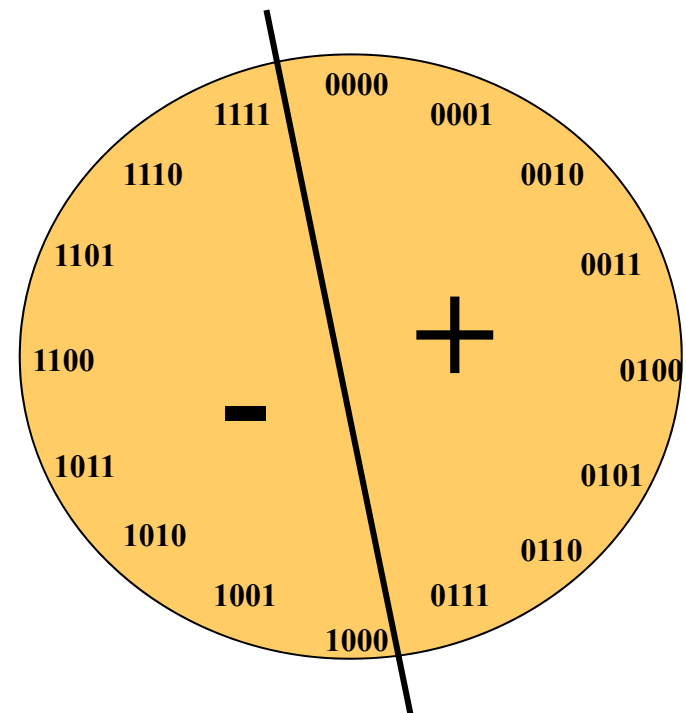
- Normal (unsigned) binary can only represent positive numbers
 - All place values are positive
- To represent negative numbers we must use a modified binary representation that takes into account sign (pos. or neg.)
 - We call these *signed* representations

Signed Number Representation

- 2 Primary Systems
 - Signed Magnitude
 - *Two's Complement (most widely used for integer representation)*

Signed numbers

- All systems used to represent negative numbers split the possible binary combinations in half (half for positive numbers / half for negative numbers)
- In both signed magnitude and 2's complement, positive and negative numbers are separated using the MSB
 - **MSB=1 means negative**
 - **MSB=0 means positive**



Signed Magnitude System

- Use binary place values but now MSB represents the sign (1 if negative, 0 if positive)

4-bit
Unsigned

Bit 3	Bit 2	Bit 1	Bit 0
8	4	2	1



0 to 15

4-bit Signed
Magnitude

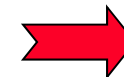
Bit 3	Bit 2	Bit 1	Bit 0
+/-	4	2	1



-7 to +7

8-bit Signed
Magnitude

Bit 7	Bit 6	Bit 5	Bit 4	Bit 3	Bit 2	Bit 1	Bit 0
+/-	64	32	16	8	4	2	1



-127 to +127

Signed Magnitude Examples

**4-bit Signed
Magnitude**

$$\begin{array}{cccc} \underline{1} & \underline{1} & \underline{0} & \underline{1} \\ +/- & 4 & 2 & 1 \end{array} = -5$$

$$\begin{array}{cccc} \underline{0} & \underline{0} & \underline{1} & \underline{1} \\ +/- & 4 & 2 & 1 \end{array} = +3$$

$$\begin{array}{cccc} \underline{1} & \underline{1} & \underline{1} & \underline{1} \\ +/- & 4 & 2 & 1 \end{array} = -7$$

Notice that +3 in signed magnitude is the same as in the unsigned system

**8-bit Signed
Magnitude**

$$\begin{array}{cccccccc} \underline{1} & \underline{0} & \underline{0} & \underline{1} & \underline{0} & \underline{0} & \underline{1} & \underline{1} \\ +/- & 64 & 32 & 16 & 8 & 4 & 2 & 1 \end{array} = -19$$

$$\begin{array}{cccccccc} \underline{0} & \underline{0} & \underline{0} & \underline{1} & \underline{1} & \underline{0} & \underline{0} & \underline{1} \\ +/- & 64 & 32 & 16 & 8 & 4 & 2 & 1 \end{array} = +25$$

Important: Positive numbers have the same representation in signed magnitude as in normal unsigned binary

Signed Magnitude Range

- Given n bits...
 - MSB is sign
 - Other $n-1$ bits = normal unsigned place values
 - Range with $n-1$ unsigned bits = $[0 \text{ to } 2^{n-1} - 1]$

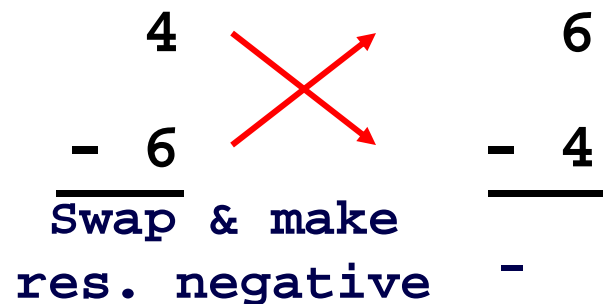
Range with n -bits of Signed Magnitude
 $[-2^{n-1} - 1 \text{ to } +2^{n-1} - 1]$

Disadvantages of Signed Magnitude

1. Wastes a combination to represent -0

$$0000 = 1000 = 0_{10}$$

2. Addition and subtraction algorithms for signed magnitude are different than unsigned binary (we'd like them to be the same to use same HW)




4 6


- 6 - 4


Swap & make
res. negative

2's Complement System

- Normal binary place values except MSB has negative weight
 - MSB of 1 = -2^{n-1}

4-bit Unsigned	Bit 3	Bit 2	Bit 1	Bit 0		0 to 15
	_____	_____	_____	_____		
	8	4	2	1		

4-bit 2's complement	Bit 3	Bit 2	Bit 1	Bit 0		-8 to +7
	_____	_____	_____	_____		
	-8	4	2	1		

8-bit 2's complement	Bit 7	Bit 6	Bit 5	Bit 4	Bit 3	Bit 2	Bit 1	Bit 0		-128 to +127
	_____	_____	_____	_____	_____	_____	_____	_____		
	-128	64	32	16	8	4	2	1		

2's Complement Examples

4-bit
2's complement

$$\begin{array}{cccc} \underline{1} & \underline{0} & \underline{1} & \underline{1} & = -5 \\ -8 & 4 & 2 & 1 & \end{array}$$

$$\begin{array}{cccc} \underline{0} & \underline{0} & \underline{1} & \underline{1} & = +3 \\ -8 & 4 & 2 & 1 & \end{array}$$

$$\begin{array}{cccc} \underline{1} & \underline{1} & \underline{1} & \underline{1} & = -1 \\ -8 & 4 & 2 & 1 & \end{array}$$

Notice that +3 in 2's comp. is the same as in the unsigned system

8-bit
2's complement

$$\begin{array}{cccccccc} \underline{1} & \underline{0} & \underline{0} & \underline{0} & \underline{0} & \underline{0} & \underline{0} & \underline{1} & = -127 \\ -128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 & \end{array}$$

$$\begin{array}{cccccccc} \underline{0} & \underline{0} & \underline{0} & \underline{1} & \underline{1} & \underline{0} & \underline{0} & \underline{1} & = +25 \\ -128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 & \end{array}$$

Important: Positive numbers have the same representation in 2's complement as in normal unsigned binary

2's Complement Range

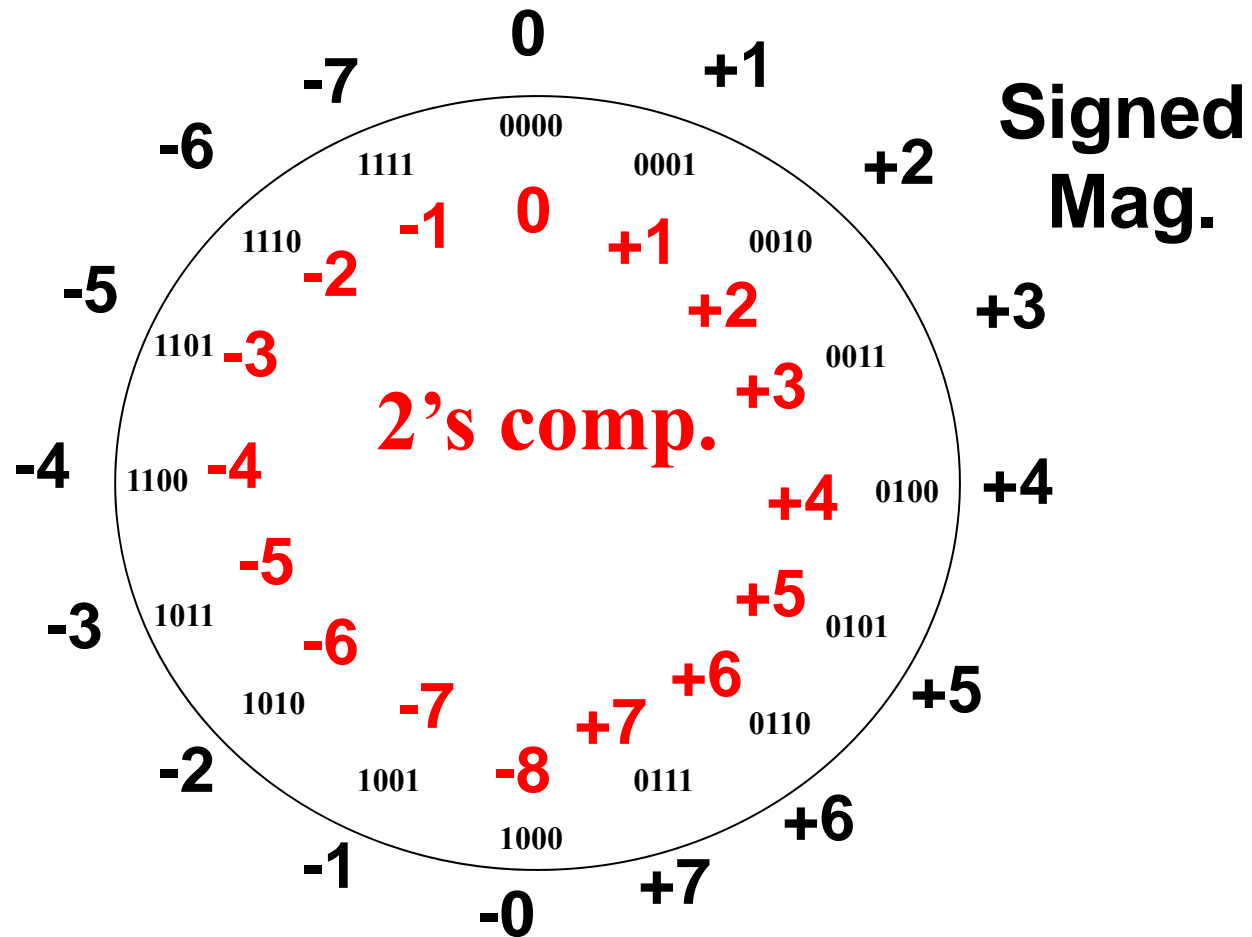
- Given n bits...
 - Max positive value = 011...11
 - Includes all n-1 positive place values
 - Max negative value = 100...00
 - Includes only the negative MSB place value

Range with n-bits of 2's complement

$$[-2^{n-1} \text{ to } +2^{n-1}-1]$$

- Side note - What decimal value is 111...11?
 - -1_{10}

Comparison of Systems



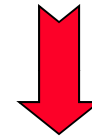
Taking the Negative

- Given a number in signed magnitude or 2's complement how do we find its negative (i.e. $-1 * X$)
 - Signed Magnitude: Flip the sign bit
 $-0110 = +6 \Rightarrow 1110 = -6$
 - 2's complement: "Take the 2's complement"
 $-0110 = +6 \Rightarrow -6 = 1010$
 - Operation defined as:
 1. Flip/invert/not all the bits (1's complement)
 2. Add 1 and drop any carry (i.e. finish with the same # of bits as we start with)

Taking the 2's Complement

- Invert (flip) each bit (take the 1's complement)
 - 1's become 0's
 - 0's become 1's
- Add 1 (drop final carry-out, if any)

-32 16 8 4 2 1
010011



101100

+ 1

101101

Original number = +19

Bit flip is called the 1's complement of a number

Resulting number = -19

Important: Taking the 2's complement is equivalent to taking the negative (negating)

Taking the 2's Complement

1

-32	16	8	4	2	1
1	0	1	0	1	0

Original number = -22

0	1	0	1	0	1
+					1
<hr/>					
0	1	0	1	1	0

Take the 2's complement yields the negative of a number

Resulting number = +22

1	0	1	0	0	1
+					1
<hr/>					
1	0	1	0	1	0

Taking the 2's complement again yields the original number (the operation is symmetric)

Back to original = -22

2

0	0	0	0
---	---	---	---

Original # = 0

1	1	1	1
+			1
<hr/>			
0	0	0	0

Take the 2's complement

2's comp. of 0 is 0

3

1	0	0	0
---	---	---	---

Original # = -8

0	1	1	1
+			1
<hr/>			
1	0	0	0

Take the 2's complement

Negative of -8 is -8
(i.e. no positive equivalent, but this is not a huge problem)

2's Complement System Facts

- Normal binary place values but MSB has negative weight
- MSB determines sign of the number
 - 0 = positive / 1 = negative
- Special Numbers
 - 0 = All 0's (00...00)
 - -1 = All 1's (11...11)
 - Max Positive = 0 followed by all 1's (011...11)
 - Max Negative = 1 followed by all 0's (100...00)
- To take the negative of a number (e.g. -7 => +7 or +2 => -2), requires taking the complement
 - 2's complement of a # is found by flipping bits and adding 1

$$\begin{array}{rcl} 1001 & x = -7 \\ 0110 & \text{Bit flip (1's comp.)} \\ + \quad 1 & \text{Add 1} \\ \hline 0111 & -x = -(-7) = +7 \end{array}$$

Summary

- **Signed Magnitude**
 - Range $(-2^{n-1} - 1 \text{ to } +2^{n-1} - 1)$
- **2's Complement**
 - Range $(-2^{n-1} \text{ to } +2^{n-1} - 1)$

Unsigned and Signed Variables

- Unsigned variables use unsigned binary (normal power-of-2 place values) to represent numbers

$$\begin{array}{rcccccccc} 1 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & = +147 \\ \hline 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \end{array}$$

- Signed variables use the 2's complement system (Neg. MSB weight) to represent numbers

$$\begin{array}{rcccccccc} 1 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & = -109 \\ \hline -128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \end{array}$$

Zero and Sign Extension

- Extension is the process of increasing the number of bits used to represent a number without changing its value

Unsigned = Zero Extension (Always add leading 0's):

$$111011 = 00111011$$

↑
Increase a 6-bit number to 8-bit number by zero extending

2's complement = Sign Extension (Replicate sign bit):

pos. $011010 = 00\overset{\curvearrowright}{0}\overset{\curvearrowright}{0}11010$

neg. $110011 = 11\overset{\curvearrowright}{1}\overset{\curvearrowright}{1}0011$

Sign bit is just repeated as many times as necessary

Zero and Sign Truncation

- Truncation is the process of decreasing the number of bits used to represent a number without changing its value

Unsigned = Zero Truncation (Remove leading 0's):

$$\text{~~00~~111011 = 111011}$$

Decrease an 8-bit number to 6-bit number by truncating 0's. Can't remove a '1' because value is changed

2's complement = Sign Truncation (Remove copies of sign bit):

$$\text{pos. } \text{~~00~~011010 = 011010}$$

$$\text{neg. } \text{~~11~~10011 = 10011}$$

Any copies of the MSB can be removed without changing the numbers value. Be careful not to change the sign by cutting off ALL the sign bits

Translating Hexadecimal

- Hex place values (16^2 , 16^1 , 16^0) can ONLY be used if the number is positive
- If hex represents unsigned binary
 1. Apply hex place values
 - B2 hex = $11 \cdot 16^1 + 2 \cdot 16^0 = 178_{10}$
- If hex represents signed value (2's comp.)
 1. Determine the sign by looking at the underlying MSB
 2. If pos., apply hex place values (as if it were unsigned)
 3. If neg., take the 16's complement and apply hex place values to find the neg. number's magnitude

Translating Hexadecimal

- Given 6C hex
 - If it is unsigned, apply hex place values
$$6C \text{ hex} = 6 \cdot 16^1 + 12 \cdot 16^0 = 108_{10}$$
 - If it is signed...
 - Determine the sign by looking at MSD
 - 0-7 hex has a 0 in the MSB [i.e. positive]
 - 8-F hex has a 1 in the MSB [i.e. negative]
 - Thus, 6C (start with 6 which has a 0 in the MSB is positive)
 - Since it is positive, apply hex place values
$$6C \text{ hex} = 6 \cdot 16^1 + 12 \cdot 16^0 = 108_{10}$$

Translating Hexadecimal

- Given B2 hex
 - If it is unsigned
$$\text{B2 hex} = 11 * 16^1 + 2 * 16^0 = 178_{10}$$
 - If it is signed
 - Determine the sign: MSD = 'B' thus negative
[i.e. (- ____)]
 - Take the 16's complement and apply hex place values to find the neg. number's magnitude

Taking the 16's Complement

- Taking the 2's complement of a binary number yields its negative and is accomplished by finding the 1's complement (bit flip) and adding 1
- Taking the 16's complement of a hex number yields its negative and is accomplished by finding the 15's complement and adding 1
 - 15's complement is found by subtracting each digit of the hex number from F_{16}

Original value B2:	FF	
	<u>- B2</u>	Subtract each digit from F
	4D	15's comp. of B2
	<u>+ 1</u>	Add 1
16's comp. of B2:	4E	16's comp. of B2

Finding the Value of Hex Numbers

- 8A hex representing a signed (2's comp.) value
 - Step 1: Determine the sign: Neg.
 - Step 2: Take the 16's comp. to find magnitude

$$\begin{array}{r} \text{FF} \\ - 8\text{A} \\ \hline 75 \\ + 1 \\ \hline 76 \end{array}$$

- Step 3: Apply hex place values ($76_{16} = +118_{10}$)
 - Step 4: Final value: B2 hex = -118_{10}
- 7C hex representing a signed (2's comp.) value
 - Step 1: Determine the sign: Pos.
 - Step 2: Apply hex place values ($7C_{16} = +124_{10}$)
- 7C hex representing an unsigned value = $+130_{10}$

ARITHMETIC AND OVERFLOW

2's Complement Addition/Subtraction

- Addition
 - Sign of the numbers do not matter
 - Add column by column
 - Drop any final carry-out
- Subtraction
 - Any subtraction $(A-B)$ can be converted to addition $(A + -B)$ by taking the 2's complement of B
 - $(A-B)$ becomes $(A + 1's \text{ comp. of } B + 1)$
 - Drop any carry-out

2's Complement Addition

- No matter the sign of the operands just add as normal
- Drop any extra carry out

$$\begin{array}{r} 0000 \\ 0011 (3) \\ + 0010 (2) \\ \hline 0101 (5) \end{array}$$

$$\begin{array}{r} 0000 \\ 1101 (-3) \\ + 0010 (2) \\ \hline 1111 (-1) \end{array}$$

Drop final carry-out →

$$\begin{array}{r} \cancel{1}110 \\ 0011 (3) \\ + 1110 (-2) \\ \hline 0001 (1) \end{array}$$

$$\begin{array}{r} \cancel{1}100 \\ 1101 (-3) \\ + 1110 (-2) \\ \hline 1011 (-5) \end{array}$$

Unsigned and Signed Addition

- Addition process is the same for both unsigned and signed numbers
 - Add columns right to left
- Examples:

	<u>1 1</u>	<u>If unsigned</u>	<u>If signed</u>
	1001	(9)	(-7)
+ 0011	<u> </u>	(3)	(3)
	1100	(12)	(-4)

2's Complement Subtraction

- Take the 2's complement of the subtrahend and add to the original minuend
- Drop any extra carry out

$$\begin{array}{r} 0011 (+3) \\ - 0010 (+2) \\ \hline \end{array}$$

$$\begin{array}{r} 1101 (-3) \\ - 1110 (-2) \\ \hline \end{array}$$


Drop final carry-out →

$$\begin{array}{r} \cancel{1}111_ \\ 0011 \\ 1101 \text{ 1's comp. of } +2 \\ + \quad 1 \text{ Add 1} \\ \hline 0001 \end{array}$$

$$\begin{array}{r} 1_ \\ 1101 \\ 0001 \text{ 1's comp. of } -2 \\ + \quad 1 \text{ Add 1} \\ \hline 1111 \end{array}$$

Unsigned and Signed Subtraction

- Subtraction process is the same for both unsigned and signed numbers
 - Convert $A - B$ to $A + \text{Comp. of } B$
 - Drop any final carry out
- Examples:

	<u>If unsigned</u>	<u>If signed</u>				
1100	(12)	(-4)		11_1_		
- 0010	(2)	(2)		1100	A	
				1101	1's comp. of B	
				+ 1	Add 1	
				1010	(10)	(-6)
				<u>If unsigned</u>	<u>If signed</u>	

Overflow

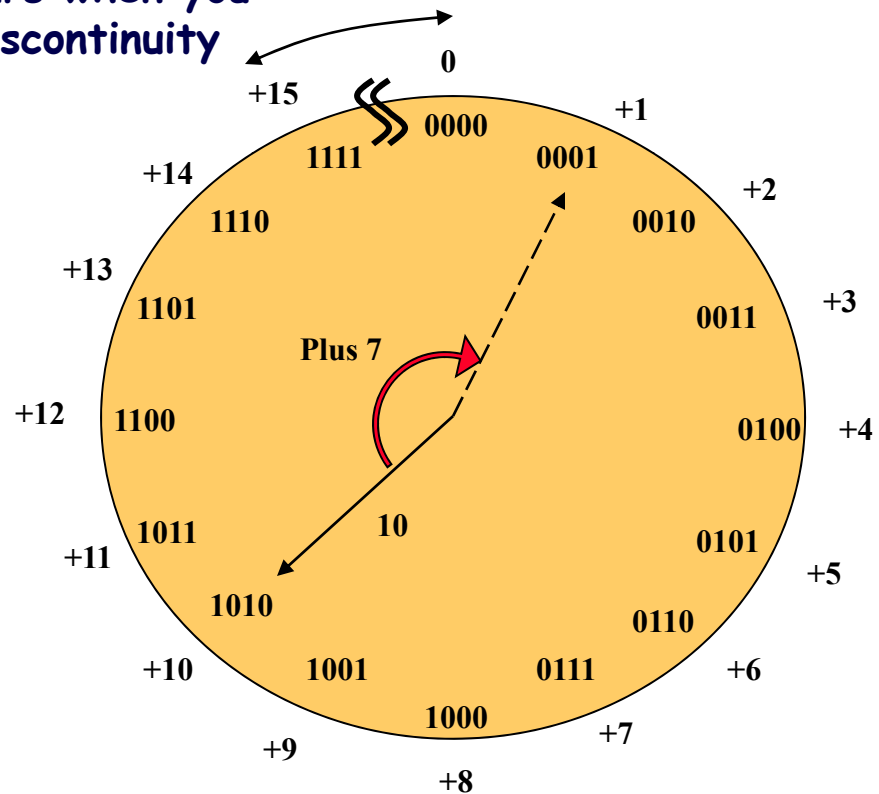
- Overflow occurs when the result of an arithmetic operation is too large to be represented with the given number of bits
- Conditions and tests to determine overflow depend on sign

Unsigned Overflow

Overflow occurs when you cross this discontinuity

$$10 + 7 = 17$$

With 4-bit *unsigned* numbers we can only represent 0 - 15.
Thus, we say overflow has occurred

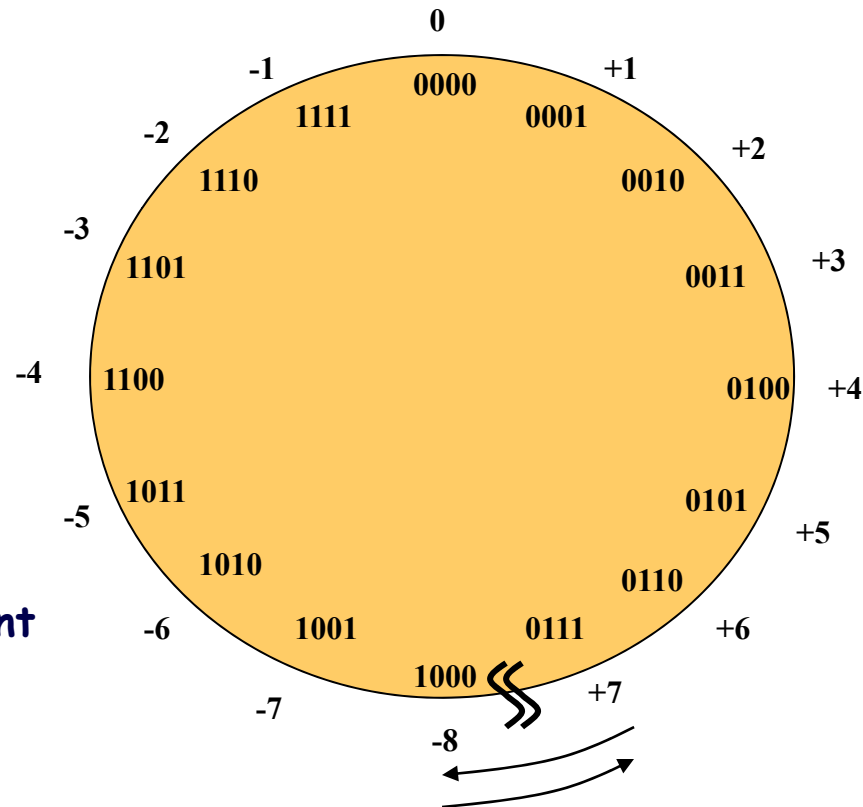


2's Complement Overflow

$$5 + 7 = +12$$

$$-6 + -4 = -10$$

With 4-bit 2's complement numbers we can only represent -8 to +7. Thus, we say overflow has occurred



Overflow occurs when you cross this discontinuity

Overflow in Addition

- Overflow occurs when the result of the addition cannot be represented with the given number of bits
- Tests for overflow:
 - Unsigned: if Cout = 1
 - Signed: if $p + p = n$ or $n + n = p$

1 1	<u>If unsigned</u>	<u>If signed</u>
1101	(13)	(-3)
+ 0100	(4)	(4)
<u>0001</u>	(17)	(+1)
	<u>Overflow</u>	<u>No Overflow</u>
	Cout = 1	n + p

0 1	<u>If unsigned</u>	<u>If signed</u>
0110	(6)	(6)
+ 0101	(5)	(5)
<u>1011</u>	(11)	(-5)
	<u>No Overflow</u>	<u>Overflow</u>
	Cout = 0	p + p = n

Overflow in Subtraction

- Overflow occurs when the result of the subtraction cannot be represented with the given number of bits
- Tests for overflow:
 - Unsigned: if Cout = 0
 - Signed: if addition is $p + p = n$ or $n + n = p$

	<u>If unsigned</u>	<u>If signed</u>
0111	(7)	(7)
- 1000	(8)	(-8)
<u> </u>	(-1)	(15)



0111_	
0111	A
0111	1's comp. of B
+ 1	Add 1
<u> </u>	
1111	(15) (-1)

Desired
Results

<u>If unsigned</u>	<u>If signed</u>
Overflow	Overflow
Cout = 0	$p + p = n$ 88

Hex Addition and Overflow

- Same rules as in binary
 - Add left to right
 - Drop any carry (carry occurs when $\text{sum} > F_{16}$)
- Same addition overflow rules
 - Unsigned: Check if final Cout = 1
 - Signed: Check signs of inputs and result

$$\begin{array}{r}
 \textcolor{blue}{1} \quad \textcolor{blue}{1} \\
 \textcolor{green}{7AC5} \\
 + \textcolor{green}{C18A} \\
 \hline
 \textcolor{green}{3C4F}
 \end{array}$$

If unsigned

Overflow
Cout = 1

If signed

No Overflow
 $p + n$

$$\begin{array}{r}
 \textcolor{blue}{0} \textcolor{blue}{1} \quad \textcolor{blue}{1} \\
 \textcolor{green}{6C12} \\
 + \textcolor{green}{549F} \\
 \hline
 \textcolor{green}{C0B1}
 \end{array}$$

If unsigned

No Overflow
Cout = 0

If signed

Overflow
 $p + p = n$

Hex Subtraction and Overflow

- Same rules as in binary
 - Convert $A - B$ to $A + \text{Comp. of } B$
 - Drop any final carry out
- Same subtraction overflow rules
 - Unsigned: Check if final Cout = 0
 - Signed: Check signs of addition inputs and result

$$\begin{array}{r}
 \text{B1ED} \\
 - \text{76FE} \\
 \hline
 \end{array}
 \xrightarrow{\text{1}}
 \begin{array}{r}
 \text{B1ED} \\
 + \text{8901} \\
 \hline
 \text{3AEF}
 \end{array}$$

If unsigned
No_Overflow
Cout = 1

If signed
Overflow
 $n + n = p$

$$\begin{array}{r}
 0001 \\
 - 0002 \\
 \hline
 \end{array}
 \xrightarrow{\text{0}}
 \begin{array}{r}
 0001 \\
 + \text{FFFD} \\
 \hline
 \text{FFFF}
 \end{array}$$

If unsigned
Overflow
Cout = 0

If signed
No Overflow
 $p + n$

Unsigned Multiplication Review

- Same rules as decimal multiplication
- Multiply each bit of Q by M shifting as you go
- An m-bit * n-bit mult. produces an m+n bit result
- Notice each partial product is a shifted copy of M or 0 (zero)

$$\begin{array}{rcl} & 1010 & \text{M (Multiplicand)} \\ * & \underline{1011} & \text{Q (Multiplier)} \\ \hline & 1010 & \\ & 1010_ & \text{PP(Partial} \\ & 0000_ & \text{Products)} \\ + & \underline{1010} & \\ \hline & 01101110 & \text{P (Product)} \end{array}$$

Signed Multiplication Techniques

- When multiplying signed (2's comp.) numbers, some new issues arise
- Must sign extend partial products (out to $2n$ bits)

**Without Sign Extension...
Wrong Answer!**

$$\begin{array}{r} 1001 = -7 \\ * 0110 = +6 \\ \hline 0000 \\ 1001_ \\ 1001_ \\ + 0000 \\ \hline 00110110 = +54 \end{array}$$

**With Sign Extension...
Correct Answer!**

$$\begin{array}{r} 1001 = -7 \\ * 0110 = +6 \\ \hline 00000000 \\ 1111001_ \\ 111001_ \\ + 00000 \\ \hline 11010110 = -42 \end{array}$$

Signed Multiplication Techniques

- Also, must worry about negative multiplier
 - MSB of multiplier has negative weight
 - If MSB=1, multiply by -1 (i.e. take 2's comp. of multiplicand)

With Sign Extension but w/o
consideration of MSB...
Wrong Answer!

$$\begin{array}{r}
 1100 = -4 \\
 * 1010 = -6 \\
 \hline
 00000000 \\
 1111100_ \\
 000000_ \\
 + 11100_ \\
 \hline
 11011000 = -40
 \end{array}$$

Place Value: -8
Multiply by -1

With Sign Extension and w/
consideration of MSB...
Correct Answer!

$$\begin{array}{r}
 1100 = -4 \\
 * \textcircled{1}010 = -6 \\
 \hline
 00000000 \\
 1111100_ \\
 000000_ \\
 + 00100_ \\
 \hline
 00011000 = +24
 \end{array}$$

Main Point: Signed and Unsigned Multiplication require
different techniques...Thus different instructions

Binary Representation Systems

- Integer Systems
 - Unsigned
 - Unsigned (Normal) binary
 - Signed
 - Signed Magnitude
 - 2's complement
 - *1's complement**
 - *Excess-N**
- Floating Point
 - For very large and small (fractional) numbers
- Codes
 - Text
 - ASCII / Unicode
 - Decimal Codes
 - BCD (Binary Coded Decimal) / (8421 Code)

Binary Codes

- Using binary we can represent any kind of information by coming up with a code
- Using n bits we can represent 2^n distinct items

Colors of the rainbow:

- Red = 000
- Orange = 001
- Yellow = 010
- Green = 100
- Blue = 101
- Purple = 111

Letters:

- 'A' = 00000
- 'B' = 00001
- 'C' = 00010
- .
- .
- .
- 'Z' = 11001

ASCII Code

- Used for representing text characters
- Originally 7-bits but usually stored as 8-bits = 1- byte in a computer
- Example:
 - `printf("Hello\n");`
 - Each character is converted to ASCII equivalent
 - 'H' = 0x48, 'e' = 0x65, ...
 - \n = newline character is represented by either one or two ASCII character
 - LF (0x0A) = line feed (moves cursor down a line)
 - CR (0x0D) = carriage return character (moves cursor to start of current line)
 - Newline for Unix / Mac = LF only
 - Newline for Windows = CR + LF

ASCII Table

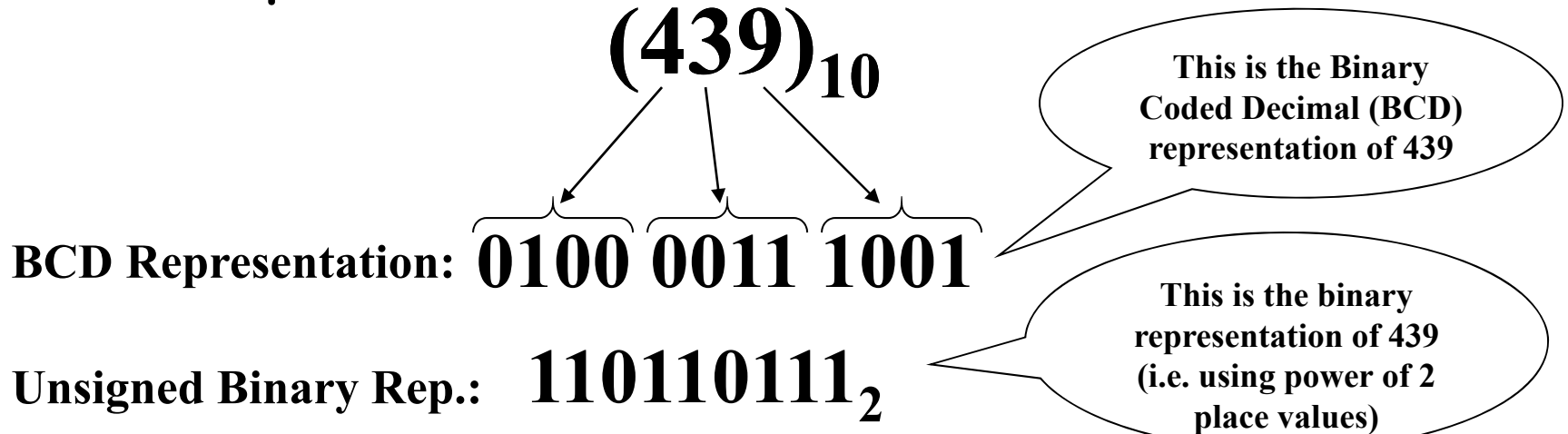
LSD/MSD	0	1	2	3	4	5	6	7
0	NULL	DLW	SPACE	0	@	P	`	p
1	SOH	DC1	!	1	A	Q	a	q
2	STX	DC2	“	2	B	R	b	r
3	ETX	DC3	#	3	C	S	c	s
4	EOT	DC4	\$	4	D	T	d	t
5	ENQ	NAK	%	5	E	U	e	u
6	ACK	SYN	&	6	F	V	f	v
7	BEL	ETB	‘	7	G	W	g	w
8	BS	CAN	(8	H	X	h	x
9	TAB	EM)	9	I	Y	i	y
A	LF	SUB	*	:	J	Z	j	z
B	VT	ESC	+	;	K	[k	{
C	FF	FS	,	<	L	\	l	
D	CR	GS	-	=	M]	m	}
E	SO	RS	.	>	N	^	n	~
F	SI	US	/	?	O	_	o	DEL

Unicode

- **Supersedes ASCII and incorporates most alphabets/characters of other languages**
- **16-bit (2-byte) character code**
- **Used by Java**

BCD

- Rather than convert a decimal number to binary which may lose some precision (i.e. 0.1_{10} = infinite binary fraction), BCD represents each decimal digit as a separate group of bits (exact decimal precision)
 - Each digit is represented as a separate 4-bit number (using place values 8,4,2,1 for each dec. digit)
 - Often used in financial and other applications where decimal precision is needed



Important: Some processors have specific instructions to operate on #'s represented in BCD